

*CSEE 4840 Final Report - Spring 2026*

# Systolic Array based CNN Accelerator

## **Instructor:**

Prof. Stephen A. Edwards

## **Group Members:**

Cheng-wen Chu (cc5397)

Chengcheng Xu (cx2355)

Harvey Lu (hl3999)

Linxiao Wu (lw3227)

Mingyuan Zheng (mz3143)

May 2026

## **Abstract**

This project implements a three-layer INT8 CNN accelerator on the DE1-SoC for classifying  $64 \times 64$  grayscale ASL digit gestures. The design uses an offline training, quantization, and golden-model flow to generate hardware preload data, then executes inference on an FPGA datapath with reused on-chip SRAM buffers, systolic-array convolution, requantization, pooling, and a fully connected classifier. The accelerator is controlled by the HPS through an MMIO interface and demonstrated through a Flask web application. The compact 3,810-parameter CNN achieves 96.77% validation accuracy with no measured INT8 quantization loss. The FPGA build meets the 50 MHz target clock, reaches a 75.72 MHz accelerator Fmax, and uses 86 of 87 DSP blocks.

# Contents

<b>1</b>	<b>Introduction</b>	<b>5</b>
<b>2</b>	<b>Algorithm</b>	<b>6</b>
2.1	Algorithm: CNN development and training . . . . .	6
<b>3</b>	<b>System Block Diagram</b>	<b>9</b>
3.1	Software Blocks . . . . .	10
3.2	Host and Runtime Path . . . . .	10
3.3	MMIO Wrapper . . . . .	11
3.4	Hardware Blocks . . . . .	12
3.5	Quantization internal computation . . . . .	13
3.6	Pooling pipeline . . . . .	14
3.7	Memory Resource Allocation . . . . .	15
3.8	Memory-to-Systolic-Array Datapath . . . . .	17
3.9	Connection Path . . . . .	20
<b>4</b>	<b>The Hardware/Software Interface</b>	<b>21</b>
4.1	Platform Designer And Bus Integration . . . . .	21
4.2	Signal-Level MMIO Interface . . . . .	21
4.3	Address Space . . . . .	22
4.4	Register Map . . . . .	23
4.5	Register Fields . . . . .	24
4.6	Scratchpad Layout . . . . .	24
4.7	Model Preload Storage and Replay . . . . .	25
4.8	Image Storage And Inference Replay . . . . .	25
4.9	Transaction Sequence . . . . .	26
4.10	Readback Path . . . . .	26
<b>5</b>	<b>Internal Memory And Data Reuse</b>	<b>26</b>
<b>6</b>	<b>Performance And Resource Results</b>	<b>27</b>
6.1	Timing . . . . .	27

6.2	Resource Utilization . . . . .	28
<b>7</b>	<b>Demo And Validation</b>	<b>29</b>
7.1	Web Demo Screenshot . . . . .	29
7.2	Board Evidence . . . . .	31
<b>8</b>	<b>Team Contributions</b>	<b>33</b>
<b>A</b>	<b>Web-To-Board Runtime Listings</b>	<b>35</b>
A.1	Focused Source Tree . . . . .	35
A.2	Included Resource Overview . . . . .	38
A.3	Web And Host Python Runtime . . . . .	41
A.3.1	Flask Web Entry Point . . . . .	41
A.3.2	Web Demo Configuration . . . . .	45
A.3.3	CPU-vs-FPGA Comparison Flow . . . . .	48
A.3.4	FPGA Board Service . . . . .	53
A.3.5	SSH Transport . . . . .	55
A.3.6	Board Transport Interface . . . . .	57
A.3.7	Inference Case Export . . . . .	58
A.3.8	Image Preprocessing . . . . .	59
A.3.9	Preprocess Mode Selection . . . . .	63
A.3.10	INT8 Quantization Helpers . . . . .	65
A.3.11	Host CPU Classifier . . . . .	65
A.3.12	Web Front-End Template . . . . .	67
A.3.13	Web Front-End Script . . . . .	73
A.3.14	Web Front-End Styles . . . . .	81
A.4	HPS Board C Runtime . . . . .	96
A.4.1	MMIO Register ABI . . . . .	96
A.4.2	HPS MMIO Helper API . . . . .	99
A.4.3	Common MMIO Implementation . . . . .	101
A.4.4	Status Probe Command . . . . .	107
A.4.5	Model Load Command . . . . .	108
A.4.6	FPGA Prediction Command . . . . .	109

A.4.7	Board CPU Baseline . . . . .	111
A.5	Golden Model And Export Code . . . . .	117
A.5.1	Golden Batch Entry Point . . . . .	117
A.5.2	Golden Case Exporter . . . . .	119
A.5.3	Hardware-Aligned Forward Pass . . . . .	122
A.5.4	SRAM Preload Packer . . . . .	127
A.5.5	Golden Parameter Loader . . . . .	130
A.5.6	Golden TXT Writer . . . . .	132
A.5.7	Flatten Helper . . . . .	134
A.5.8	Golden Max-Pool Helper . . . . .	135
A.5.9	TFLite Multiplier Helper . . . . .	136
A.5.10	TFLite Parameter Exporter . . . . .	137
A.5.11	Digit CNN Retraining Script . . . . .	142
A.6	Quartus And Platform Designer Setup . . . . .	146
A.6.1	Platform Designer Component Descriptor . . . . .	146
A.6.2	Quartus Build Wrapper . . . . .	149
A.6.3	Quartus Project Script . . . . .	157
A.6.4	Platform Designer System . . . . .	158
A.6.5	Timing Constraint File . . . . .	171
A.7	FPGA Board Boundary And Accelerator Top . . . . .	171
A.7.1	DE1-SoC Board Wrapper . . . . .	171
A.7.2	Avalon-MM MMIO Wrapper . . . . .	175
A.7.3	Accelerator Top . . . . .	184
A.7.4	Network-Level FSM . . . . .	193
A.7.5	Layer Runner FSM . . . . .	201
A.7.6	Convolution Data Adapter . . . . .	206
A.7.7	Weight Prepad Inserter . . . . .	209
A.8	Convolution Core RTL . . . . .	211
A.8.1	Convolution Top . . . . .	211
A.8.2	Convolution Engine Controller . . . . .	217
A.8.3	Convolution Buffer . . . . .	221

A.8.4	Line Buffer . . . . .	230
A.8.5	Input Row Aligner . . . . .	232
A.8.6	Systolic-Array Skew Feeder . . . . .	234
A.8.7	Systolic Array Top . . . . .	237
A.8.8	Weight Buffer . . . . .	243
A.9	Quantization, Pooling, And FC RTL . . . . .	246
A.9.1	Convolution-Quantization-Pooling Wrapper . . . . .	246
A.9.2	Convolution-Quantization Adapter . . . . .	252
A.9.3	Quantization Parameter Loader . . . . .	256
A.9.4	Quantization-Pooling Adapter . . . . .	259
A.9.5	Quantization Processing Element . . . . .	260
A.9.6	Quantization Top . . . . .	264
A.9.7	Pooling Core . . . . .	266
A.9.8	Pooling Stream Top . . . . .	272
A.9.9	Fully-Connected Top . . . . .	273
A.9.10	FC Bias Loader . . . . .	274
A.9.11	FC Data Adapter . . . . .	276
A.9.12	FC MAC . . . . .	279
A.10	SRAM Controller RTL . . . . .	280
A.10.1	SRAM Address Generator . . . . .	280
A.10.2	FC Weight Preload Packer . . . . .	281
A.10.3	SRAM A Controller . . . . .	283
A.10.4	SRAM A Wrapper . . . . .	287
A.10.5	SRAM B Controller . . . . .	288
A.10.6	SRAM B Wrapper . . . . .	290
A.10.7	SRAM FCW Wrapper . . . . .	291
A.10.8	Top SRAM A Integration . . . . .	292
A.10.9	Top SRAM B Integration . . . . .	298

# 1 Introduction

This project presents a three-layer integer-quantized convolutional neural network (CNN) accelerator implemented on the DE1-SoC platform. The accelerator classifies hand signs for the digits 0–9 from the American Sign Language (ASL) Digits dataset, using  $64 \times 64$  grayscale images as input, then performs end-to-end inference entirely in INT8 arithmetic. The network achieves 96.77% top-1 validation accuracy, and post-training quantization introduces zero accuracy degradation relative to the floating-point baseline.

The system is partitioned into two domains: an *offline software toolchain* that trains the model, applies full-integer post-training quantization, and packages the parameters into binary load streams; and a *synthesizable RTL pipeline* that receives those streams through a 32-bit streaming port, stores parameters in on-chip SRAM, and executes a pipelined convolution→quantization→pooling datapath followed by a fully-connected classifier and argmax reduction.

At a user level, the system works as a gesture-recognition demo. A user uploads or captures a hand image in the browser, the host software converts it into the same  $64 \times 64$  INT8 input format used during training, and the DE1-SoC board returns the predicted digit class. Figure 1 shows the ten target gesture classes used by the demo.



Figure 1: American Sign Language (ASL) Digits

The remainder of the report follows the path of the project from model design to board demonstration. We first describe the gesture dataset, compact CNN architecture, training setup, and INT8 quantization flow. We then move from the algorithm into the deployed system, explaining how the host software, HPS Linux tools, Platform Designer interconnect, MMIO wrapper, and RTL accelerator fit together. The middle sections give the detailed hardware/software contract: where model and image data are stored, how the HPS writes and reads the accelerator, and how the internal FPGA memories are reused across CNN layers. The final sections summarize resource usage, timing, web and board validation results, and the operating constraints of the current im-

plementation.

## 2 Algorithm

### 2.1 Algorithm: CNN development and training

#### Dataset

We use the open *ASL Digits 0–9* dataset published on Kaggle by Victor Anthony which contains RGB photographs of American Sign Language hand gestures for the digits 0–9. To match the compact input format expected by the hardware accelerator, every image is converted to single-channel grayscale and resampled to  $64 \times 64$  using `cv2.INTER_AREA`, giving an input tensor of shape  $(64, 64, 1)$  with pixel values normalized to  $[0, 1]$  via a  $1/255$  rescaling.

The 2,062 processed images are partitioned per class into a 70%/15%/15% train / validation / test split, producing 1,444 training, 310 validation, and 308 test images. Class balance is preserved across all three splits (roughly 144 training and 31 validation images per digit).

#### Network architecture

The classifier is a deliberately small convolutional network chosen so that every layer maps cleanly onto the accelerator’s MAC array and on-chip memory budget. It contains three  $3 \times 3$ , stride-1, valid-padding convolutions with ReLU activation, each followed by a  $2 \times 2$  max-pooling stage, a flatten, and a single fully connected classification head (Table 1). The resulting model has only 3,810 trainable parameters, which makes it tractable for FPGA deployment.

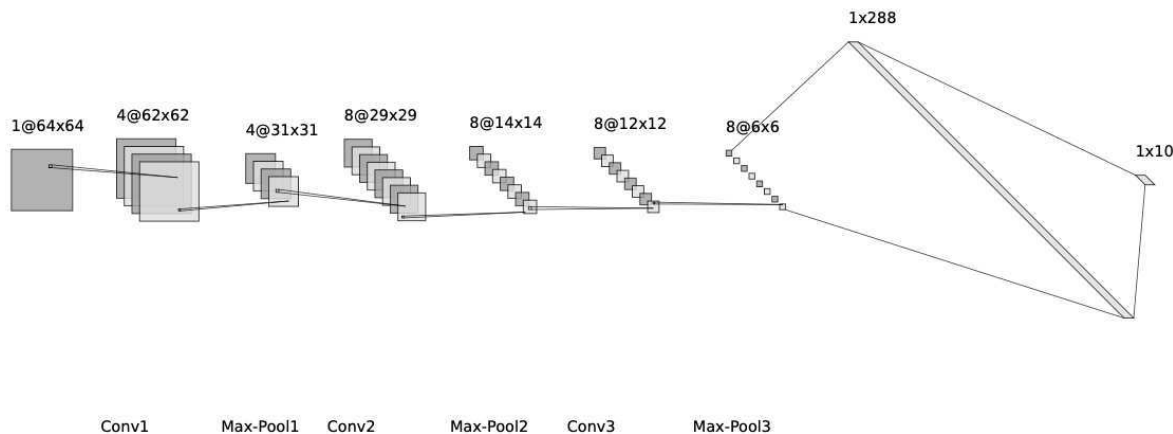


Figure 2: CNN Architecture

Table 1: CNN architecture for  $64 \times 64$  grayscale digit gestures.

Layer	Kernel / Size	Output shape	Params
Input	–	$64 \times 64 \times 1$	0
Conv2D + ReLU	$3 \times 3$ , $C_{\text{out}} = 4$	$62 \times 62 \times 4$	40
MaxPool	$2 \times 2$	$31 \times 31 \times 4$	0
Conv2D + ReLU	$3 \times 3$ , $C_{\text{out}} = 8$	$29 \times 29 \times 8$	296
MaxPool	$2 \times 2$	$14 \times 14 \times 8$	0
Conv2D + ReLU	$3 \times 3$ , $C_{\text{out}} = 8$	$12 \times 12 \times 8$	584
MaxPool	$2 \times 2$	$6 \times 6 \times 8$	0
Flatten	–	288	0
Dense + Softmax	$288 \rightarrow 10$	10	2,890
<b>Total</b>			<b>3,810</b>

### Training procedure

The network is trained in TensorFlow/Keras with categorical cross-entropy loss and the Adam optimizer (default learning rate  $10^{-3}$ ). We use a batch size of 32 and train for up to 200 epochs. Each epoch iterates over  $\lceil 1444/32 \rceil = 46$  training steps; validation uses all 310 images. After convergence the network reaches 96.77% top-1 accuracy on the validation split, with per-class F1 scores at or above 0.90 for every digit (Table 2).

Table 2: Per-class validation metrics for the float32 model (support = 31 per class).

Digit	Precision	Recall	F1	Support
0	0.939	1.000	0.969	31
1	0.969	1.000	0.984	31
2	1.000	0.968	0.984	31
3	0.968	0.968	0.968	31
4	0.968	0.968	0.968	31
5	1.000	1.000	1.000	31
6	0.903	0.903	0.903	31
7	0.968	0.968	0.968	31
8	0.966	0.903	0.933	31
9	1.000	1.000	1.000	31
<b>Accuracy</b>	<b>0.9677 (310 images)</b>			

### Post-training quantization

Floating-point MAC units are costly to implement on FPGA — they require wide mantissa/exponent datapaths, normalization and rounding logic, and significantly more DSP and logic resources than their integer counterparts. To make the network hardware-friendly, every weight, activation,

and accumulator in the inference graph must therefore be reduced to fixed-width integer arithmetic. We achieve this with *full-integer post-training quantization* (PTQ) to `int8` using the TensorFlow Lite converter.

After PTQ, we evaluated the `int8` model on the full 310-image validation split end-to-end in `int8` gives 96.77% top-1 accuracy, matching the `float32` baseline ( $\Delta_{\text{acc}} = 0.00$ ) with a float/`int8` prediction agreement of 98.39%, confirming that the move to integer arithmetic costs essentially no accuracy for this network. The quantized weights, biases, and per-layer scales are then exported and used by the MATLAB hardware-aligned reference model and the downstream FPGA bitstream.

### 3 System Block Diagram

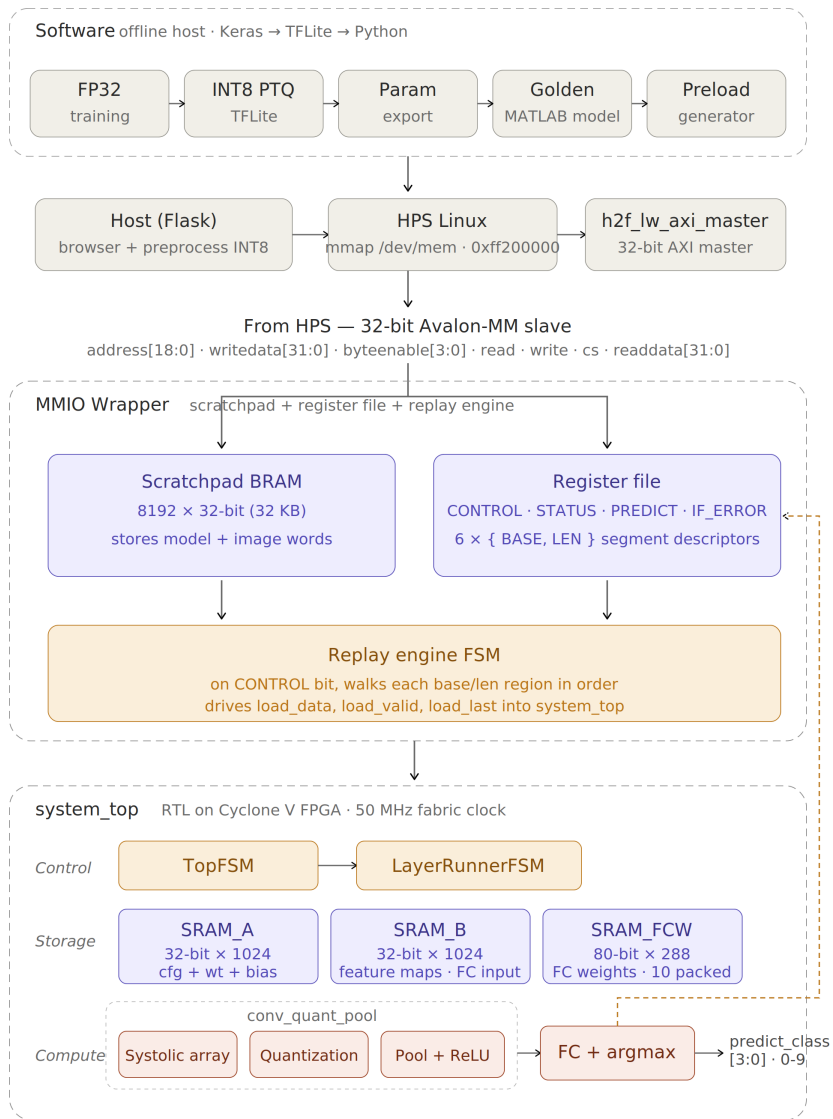


Figure 3: Block Diagram

Figure 3 shows the block diagram of our project, a three-layer INT8 CNN accelerator that performs end-to-end inference from a  $64 \times 64 \times 1$  grayscale image to a 0–9 gesture class label. The system is partitioned into an *offline software toolchain* that runs on a development host (top region) and a *synthesizable RTL pipeline* that runs on the FPGA (bottom region). The HPS writes staged payloads through the MMIO wrapper, and the wrapper replays those 32-bit words into the accelerator’s streaming load port.

### 3.1 Software Blocks

The software side is a sequential offline pipeline. We first train a floating-point model, quantize it to INT8, export the quantized parameters, run a bit-exact golden reference, and finally pack everything into the binary streams that the hardware expects.

**FP32 training.** A three-layer CNN, consisting of three convolution/pooling stages followed by a fully-connected classifier, is trained in a standard deep-learning framework on a 0–9 hand-gesture dataset. The output of this stage is a floating-point model checkpoint.

**INT8 quantization.** The trained floating-point model is converted into a fully INT8 model through post-training quantization. Weights use per-channel symmetric quantization and activations use per-tensor asymmetric quantization. For every layer this stage also produces the requantization multiplier and right-shift value that the hardware will later apply to each accumulator before saturating it back to the INT8 range.

**Parameter export.** The quantized model is parsed by a set of host scripts that extract the per-channel weights, biases, scales, and zero-points, and emit them as plain numerical parameter files that downstream stages can consume.

**Golden reference model.** A bit-exact software re-implementation of the quantized network is run on the host to produce per-layer reference tensors. These tensors are compared against the hardware outputs during RTL simulation so that every pipeline stage can be verified independently.

**SRAM preload generator.** The final software stage takes the exported parameters and the input images and reformats them for the hardware. Convolution weights are reordered into the column-interleaved layout expected by the systolic array, and the fully-connected weights are packed into the wide slot layout expected by the dedicated FC weight memory. The output is a sequence of 32-bit words that forms the load stream.

**Software-internal communication.** The software blocks are sequential host programs; they communicate strictly through files. There is no run-time communication between them.

### 3.2 Host and Runtime Path

The offline software pipeline described above only produces files; the runtime path is what carries one inference from a user click in the browser down to a single MMIO write seen by the FPGA. It spans three distinct machines in three distinct languages.

**Browser front end.** A small JavaScript page running in the user’s browser handles image upload and webcam capture, posts the image to the host as a PNG/JPG blob, and renders the predicted class, timing breakdown, and confidence scores returned by the backend. The page does not emulate the accelerator: it never touches the model itself.

**Flask host runtime.** A Python Flask application on the development host receives the browser request, preprocesses the image into the same  $64 \times 64 \times 1$  INT8 tensor used during training, and exports a temporary inference case directory. The host then transfers the case to the DE1-SoC board over SCP and triggers the board-side tool through SSH.

**HPS Linux tool.** A C program, `hps_mmio_predict`, runs on the embedded Linux system on the Cyclone V ARM core. It opens `/dev/mem` and `mmaps` the accelerator’s 4 MiB Avalon-MM window at HPS physical address `0xff200000`, packs the model and image payloads into 32-bit words, writes them into the wrapper’s scratchpad memory, programs the base/length descriptors, and finally sets the `CONTROL` start bit. After triggering inference it polls the `STATUS` register and reads back `PREDICT`, `IF_ERROR`, and the per-stage profile counters.

**HPS-to-FPGA bridge.** Every store from the C program is transported across the Cyclone V die by the `h2f_lw_axi_master` port of the HPS hard processor system, which exposes a 32-bit lightweight AXI master interface. Platform Designer automatically inserts the required AXI-to-Avalon-MM adapter between this master and our custom Avalon-MM slave, so the RTL only ever sees an ordinary Avalon-MM transaction.

**Host-to-board communication.** The host and the board run in different operating systems and on different CPUs. They communicate exclusively through SSH and SCP; there is no shared memory, no socket protocol, and no streaming RPC. This keeps the host code portable and the board code minimal, at the cost of a per-inference startup overhead that is reflected in the difference between the pure RTL cycle count and the HPS-observed wait time reported in Section 6.

### 3.3 MMIO Wrapper

The MMIO wrapper is the entire RTL interface between the Avalon-MM bus and the streaming load port that drives `system_top`. It contains a write-once scratchpad memory, a small register file, and a replay engine that turns the staged payload into the streaming protocol the accelerator core expects. Its purpose is to decouple HPS write timing from the deterministic per-cycle streaming required by the convolution datapath.

**Address space split.** The wrapper exposes a single 19-bit Avalon-MM word address. Bit `address[18]` selects between two logical spaces: `address[18] = 0` addresses the scratchpad BRAM,

and `address[18] = 1` addresses the configuration and status register file. From the Linux side this is a flat memory-mapped region; the partition is entirely an internal decoding decision.

**Scratchpad BRAM.** An  $8192 \times 32$ -bit on-chip BRAM (32 KB) stores the four model preload segments and the input image before inference begins. The HPS tool writes every word into this scratchpad through the Avalon-MM interface; the wrapper later replays them into the accelerator core word-by-word. Staging the payload in BRAM, rather than streaming it directly from the HPS, removes any timing dependence between the (asynchronous) software writes and the (strictly cycle-deterministic) streaming load protocol.

**Register file.** Fourteen 16-bit registers expose the control surface to the HPS. `CONTROL` kicks off model preload or inference; `STATUS` reports busy, model-loaded, inference-done, and the predicted class; `PREDICT` latches the final class and `IF_ERROR` reports interface-level errors such as a start request while the engine is still busy or an out-of-range scratchpad access. Six base/length pairs (`{CONV_CFG, CONV_WT, FC_BIAS, FCW, IMAGE}` each with `BASE` and `LEN`) describe where each preload segment lives in the scratchpad and how many 32-bit words it spans. The complete register map and bit-field definitions are listed in Section ??.

**Replay engine FSM.** When the HPS sets `CONTROL[0]` or `CONTROL[1]`, the replay engine walks each configured base/length region in order, reads one 32-bit word per cycle from the scratchpad, and drives the three-signal streaming protocol (`load_data[31:0]`, `load_valid`, `load_last`) into `system_top`. During model load it walks the four model regions; during inference it walks only the image region, in which case `system_top` consumes the stream directly as the first convolution layer's input. When the accelerator signals completion, the replay engine returns to idle and sets the corresponding done flag in `STATUS`.

**Readback path.** Inference results travel back through the same wrapper. The `system_top` output `predict_class[3:0]` and the per-stage cycle counters are latched into the `STATUS/PREDICT/PROFILE` registers when the accelerator asserts `predict_done`. The HPS tool reads these registers through the Avalon-MM `readdata` channel and forwards them to the host, which renders them in the browser.

### 3.4 Hardware Blocks

**Systolic Array (SA) structure** The systolic array is organized as a 2-D grid of processing elements (PEs), where each PE performs a local MAC operation on streaming activations and weights [1]. Each PE receives input activations from one direction (from the left) and weights from an orthogonal direction (from the top), multiplies them, accumulates the result into a local register, and then forwards the activation and weight to its neighboring PEs in the next cycle.



by a learned or pre-defined factor  $s$ , then rounded and saturated into an 8-bit range. Concretely, the computation can be described as:

$$y_{\text{float}} = s \cdot x, \quad y_{\text{int}} = \text{round}(y_{\text{float}}), \quad y_{\text{q}} = \text{clip}(y_{\text{int}}, -128, 127)$$

where  $y_{\text{q}}$  is the int8 activation written back to memory or passed to the next layer.

In hardware, this is implemented as a shift-and-add or multiplier stage for the scaling factor, followed by a rounding unit and a saturating adder/comparer block, which is shown as Fig. ?? . The QNT unit is fully pipelined so it can accept one MAC result per cycle from the systolic array, ensuring that quantization does not become a performance bottleneck for the accelerator dataflow.

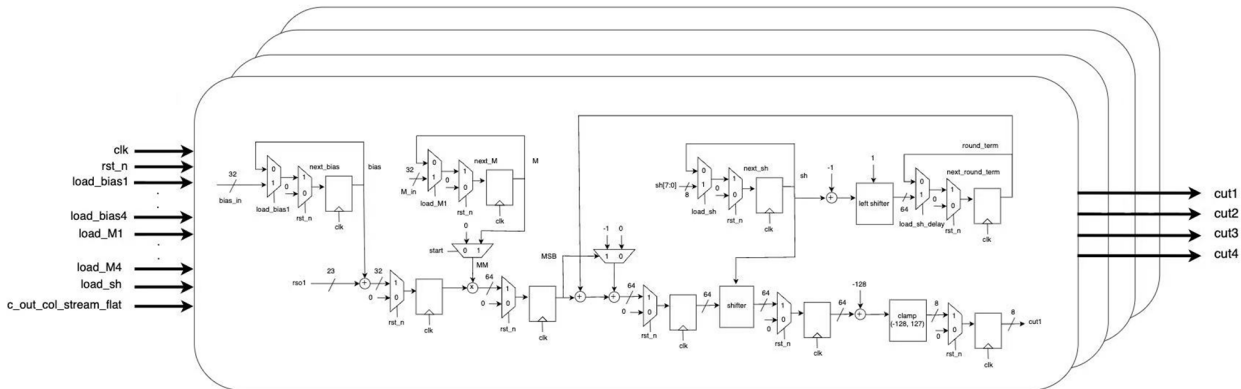


Figure 6: Quantization: 32-bit accumulator  $\rightarrow$  8-bit activation

### 3.6 Pooling pipeline

The pooling pipeline processes the quantized feature maps to reduce spatial resolution and increase invariance. After a convolution layer is completed and its outputs are quantized, the feature map is streamed into the pooling unit, which forms  $2 \times 2$  spatial windows and computes the maximum value over each window.

This pooling operation is also pipelined: each stage performs part of the comparison reduction, so that a new window can enter the pipeline once the required inputs are available. The pooled outputs are then written back to on-chip memory and reused as inputs to the next convolution layer, forming a continuous conv  $\rightarrow$  quantization  $\rightarrow$  max-pooling pipeline across the 3-layer CNN.

**Load port.** The only data entry point of the chip. The host streams 32-bit beats with a valid/ready/last handshake. A mode signal selects between model-load traffic (configuration, convolution weights, FC bias, and FC weights) and image-load traffic, which also triggers inference.

**Top-level FSM.** A network-level sequencer. During model-load it routes the four preload segments to their destinations in the on-chip memories. During inference it drives the fixed layer schedule of the three convolutional/pooling stages followed by the fully-connected layer, and finally an argmax reduction over the ten output accumulators that produces the predicted class.

**Layer-runner FSM.** A per-layer transaction controller that loads the current layer’s configuration and weights, streams data through the convolution datapath, and waits for the frame to finish. It decouples the network-level schedule from the per-layer streaming protocol and generates the control signals for the on-chip memories.

**SRAM\_A.** A 32-bit memory that stores the convolution configuration words, the convolution weights, and the fully-connected bias vector. It also doubles as the writeback buffer for the second convolution layer, using a stride-2 interleaved address region.

**SRAM\_B.** A 32-bit memory that holds the pooled feature maps between layers. The output of the first pooling stage and the output of the last pooling stage share the same physical memory; the later write simply overwrites the earlier data once it is no longer needed.

**SRAM\_FCW.** A dedicated wide memory for the fully-connected weights. Each slot packs the per-output-channel weights for one flattened input position, so that all ten output channels can be fed in parallel in a single read. On the FPGA the internal storage can be mapped onto a block-memory IP.

All three memories receive control signals from the FSMs and expose streaming interfaces for data access.

### 3.7 Memory Resource Allocation

The memory subsystem supports both model parameter storage and intermediate activation buffering while enabling data reuse across CNN layers. The design adopts three logical on-chip memories: **SRAM\_A**, **SRAM\_B**, and **SRAM\_FCW**. Although these modules follow SRAM-style interfaces, they are implemented in the FPGA using Verilog register arrays rather than physical SRAM macros.

**SRAM\_A (parameter and activation storage).** **SRAM\_A** is a 32-bit  $\times$  1024-word memory, giving a total capacity of 4 KB. During model preload, it stores convolution configuration words, convolution weights, and fully-connected bias values. During inference, it is also reused as an intermediate activation buffer. After the second convolution and pooling stage, the  $14 \times 14 \times 8 = 1568$  byte output feature map is written into **SRAM\_A** and later read as the input to the third convolution layer.

**SRAM\_B (activation buffer).** SRAM\_B is also a  $32\text{-bit} \times 1024\text{-word}$  memory. It is used mainly for intermediate feature maps between convolution layers. The first pooling stage writes a  $31 \times 31 \times 4 = 3844$  byte feature map into SRAM\_B, which is then reused as the input to the second convolution layer. Later, the third pooling stage writes the  $6 \times 6 \times 8 = 288$  byte final feature map back into SRAM\_B. This final feature vector is then consumed by the fully-connected classifier.

**SRAM\_FCW (fully-connected weight memory).** SRAM\_FCW is a dedicated wide memory for fully-connected weights. Each entry is 80 bits wide and contains ten INT8 weights, one for each output class. The implemented memory has 512 physical entries, with 288 entries used by the current  $6 \times 6 \times 8$  feature vector. During execution, one 80-bit word is read per cycle, allowing one input activation to update all ten class accumulators in parallel.

**Streaming input and dataflow.** The original input image is not first stored in SRAM\_A or SRAM\_B. Instead, the HPS stages the packed image in the MMIO scratchpad, and the MMIO wrapper replays those 32-bit words directly into the first convolution layer. After that point, only pooled INT8 feature maps are written back to the local SRAM buffers for reuse by later layers.

Across layers, the feature map sizes are:

- Input:  $64 \times 64 \times 1 = 4096$  bytes
- Layer 1 output:  $31 \times 31 \times 4 = 3844$  bytes
- Layer 2 output:  $14 \times 14 \times 8 = 1568$  bytes
- Layer 3 output:  $6 \times 6 \times 8 = 288$  bytes

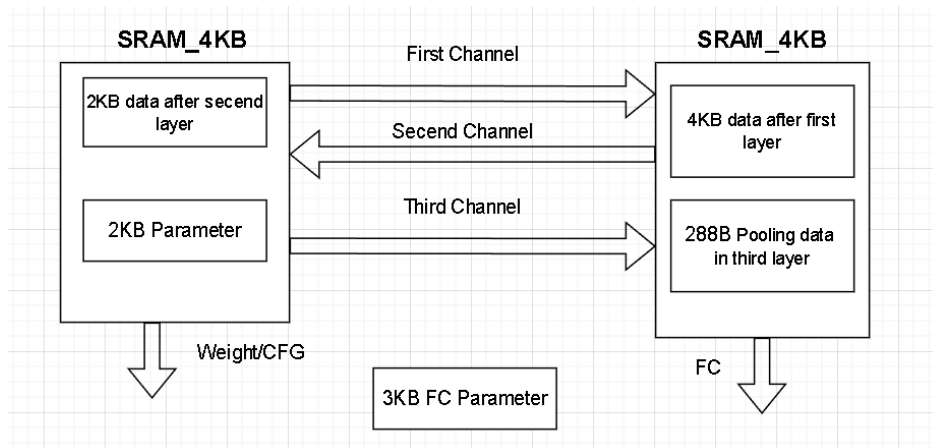


Figure 7: Memory organization and data reuse across CNN layers.

The overall dataflow follows the reuse pattern:

Input -> L1 -> SRAM\_B -> L2 -> SRAM\_A -> L3 -> SRAM\_B -> FC

**Design rationale.** The memory architecture is optimized for three objectives: (1) minimizing on-chip memory usage through buffer reuse, (2) reducing external memory bandwidth by keeping intermediate feature maps on chip, and (3) improving fully-connected throughput through wide parallel weight reads. This balance allows the full CNN pipeline to fit within the limited FPGA memory budget.

### 3.8 Memory-to-Systolic-Array Datapath

After the feature maps are stored or replayed through the memory system, the convolution core consumes them as a packed 32-bit stream. This section describes how that stream is unpacked, buffered, converted into a  $3 \times 3$  sliding window, and scheduled into the systolic array.

The activation source depends on the active layer. For layer 1, the input image is replayed directly from the MMIO scratchpad. For layer 2, the input feature map is read from SRAM\_B. For layer 3, the input feature map is read from SRAM\_A. In all three cases, the convolution core receives a 32-bit word stream and interprets the packed bytes according to the active layer configuration.

The line-buffer byte grouping changes because the number of channels per spatial position changes after each pooling stage. Layer 1 consumes a  $64 \times 64 \times 1$  input, so each spatial position contains one INT8 byte. Layer 2 consumes a  $31 \times 31 \times 4$  feature map, so each spatial position contains four INT8 channel values. Layer 3 consumes a  $14 \times 14 \times 8$  feature map, so each spatial position contains eight INT8 channel values. The SRAM access width remains fixed at 32 bits, while the input adapter and line buffers reinterpret the byte stream for the current layer[2].

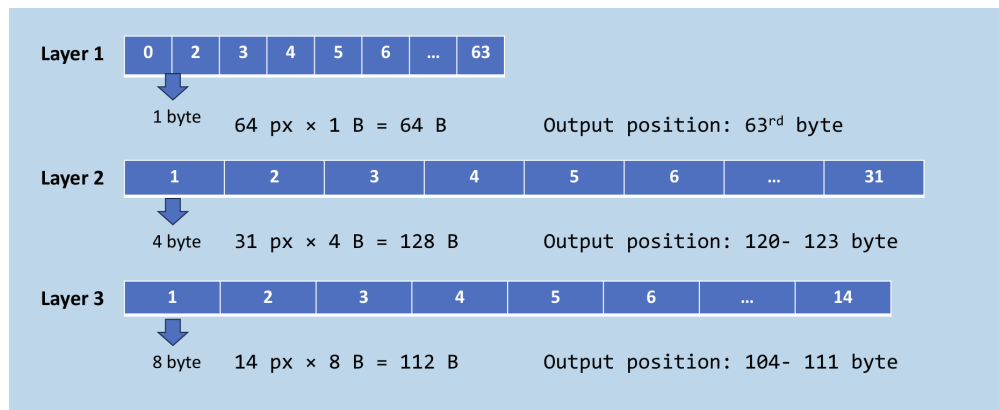


Figure 8: Layer-dependent row layout for the line-buffer input stream.

Figures 9 and 10 illustrate the sliding-window formation using a simplified  $8 \times 8$ , single-channel input. The same buffering principle is used in the actual accelerator, but the real row width and byte grouping are selected from the active layer configuration shown in Figure 8.

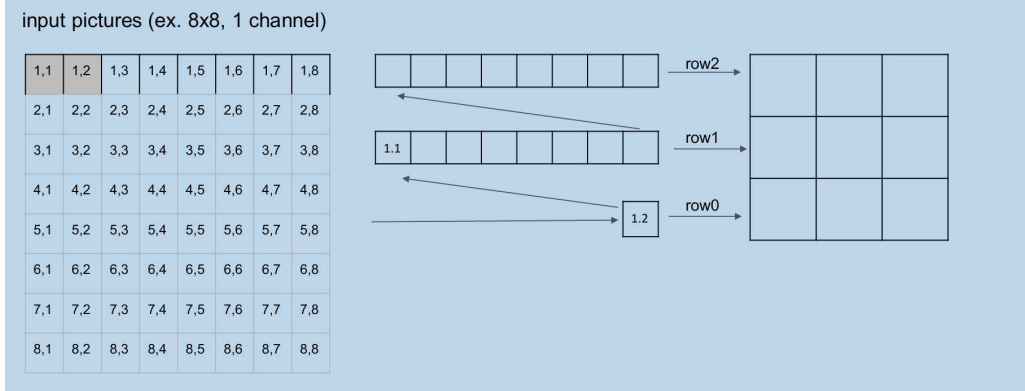


Figure 9: Initial line-buffer fill before a complete  $3 \times 3$  convolution window is valid.

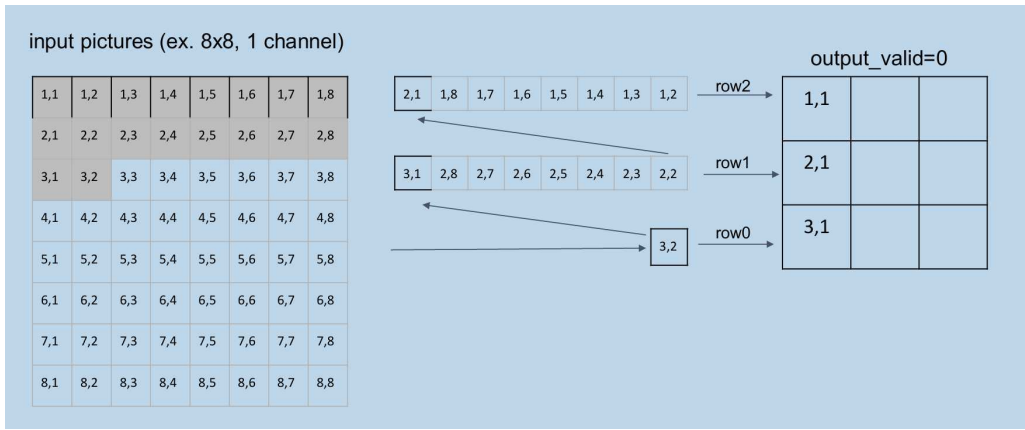


Figure 10: Formation of the three aligned rows used by the  $3 \times 3$  sliding convolution window.

During the warm-up phase, the line buffers consume the packed memory stream and delay earlier rows so that three adjacent image rows can be presented at the same time. The convolution window is not marked valid until the datapath has received enough rows and columns to populate all nine positions of the  $3 \times 3$  window. Once the window is valid, the input aligner and skew feeder schedule the INT8 activation values into the systolic array, where they are multiplied with reordered INT8 weights and accumulated into wider partial sums.

**Convolution datapath (conv-quant-pool).** The output of the systolic array continues through a three-stage streaming compute path:

- **Systolic array** — forms MAC results from the  $3 \times 3$  sliding window using parallel  $\text{INT8} \times \text{INT8}$  multiply-accumulate operations.
- **Quantization stage** — rescales each partial sum using the per-channel multiplier, shift, and effective bias, then saturates the result back to INT8.

- **Pooling stage** — performs  $2 \times 2$  stride-2 max-pooling and fuses the ReLU behavior before writeback.

The three stages are connected by streaming handshakes, so intermediate convolution sums are not written back to memory. Only the compact INT8 pooled feature maps are stored in SRAM\_A or SRAM\_B for reuse by the next layer.

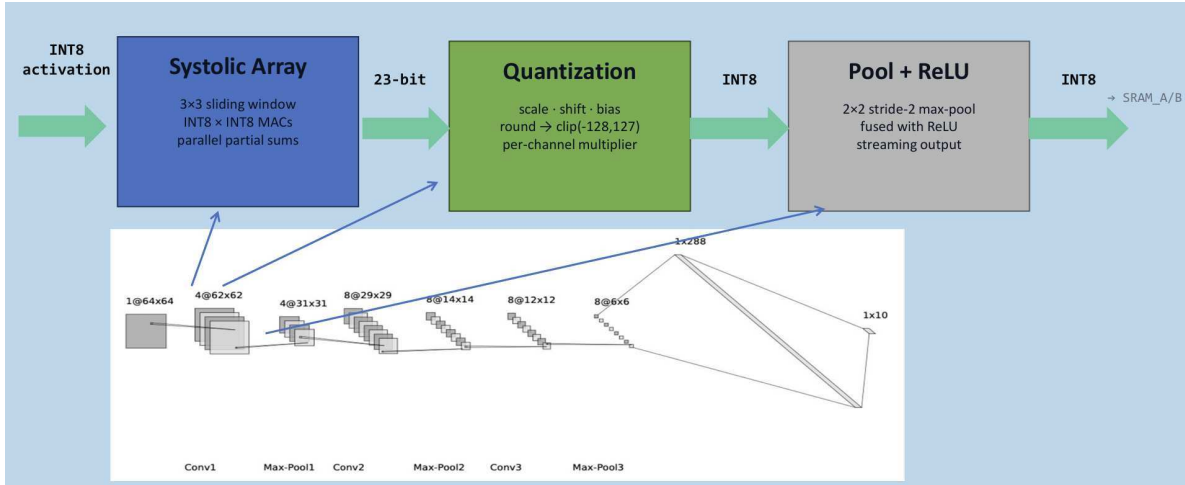


Figure 11: Streaming convolution datapath from the systolic array through quantization and pooling.

**FC and argmax.** A ten-way parallel multiply-accumulate array that computes all ten class scores simultaneously, reading the broadcast pixel from the feature-map memory, the ten per-channel weights from the wide FC weight memory, and the bias vector from the configuration memory. The ten final accumulators are reduced by the top-level FSM into the 4-bit predicted class, together with a one-cycle valid pulse on the output pins.

Memory	Width	Depth	Size	Function
SRAM_A	32-bit	1024	4 KB	Stores convolution configuration, convolution weights, FC bias, and layer-2 output feature maps
SRAM_B	32-bit	1024	4 KB	Stores intermediate feature maps (layer-1 and layer-3 outputs) and serves as FC input buffer
SRAM_FCW	80-bit	512 physical, 288 used	5 KB physical, 2.88 KB used	Stores fully-connected weights; each entry contains 10 INT8 weights for parallel computation

Table 3: On-chip memory allocation and functionality

### 3.9 Connection Path

One FPGA inference request crosses the host, HPS Linux, Platform Designer, and RTL boundaries shown in Table 4.

Table 4: End-to-end connection path for one FPGA inference request.

Step	Actor	Action	Data or interface
1	Browser	Upload or capture one gesture image.	PNG/JPG image bytes.
2	Host Flask runtime	Preprocess the image into the model input format.	$64 \times 64 \times 1$ INT8 tensor.
3	Host runtime	Export a temporary inference case and copy it to the board.	Case directory sent by SCP.
4	HPS Linux tool	Run <code>hps_mmio_predict</code> .	Board-side command over SSH.
5	HPS Linux tool	Map the accelerator address range.	<code>/dev/mem</code> at <code>0xff200000</code> .
6	HPS Linux tool	Write packed image/model words into the MMIO scratchpad.	32-bit H2F MMIO writes.
7	HPS Linux tool	Write base/length and control registers.	Avalon-MM config register space.
8	MMIO wrapper	Replay scratchpad words into the streaming accelerator core.	<code>load_data</code> , <code>load_valid</code> , <code>load_last</code> .
9	CNN accelerator	Run the fixed CNN schedule and latch outputs.	Predicted class and profile counters.
10	HPS Linux tool	Read back status, prediction, error, and profile registers.	MMIO reads returned to Flask/browser.

## 4 The Hardware/Software Interface

### 4.1 Platform Designer And Bus Integration

Use	Connections	Name	Description	Export	Clock	Base	End	...	Tags	Opcode Name
✓		<b>clk_0</b>	Clock Source							
		clk_in	Clock Input	clk	exporte					
		clk_in_reset	Reset Input	reset						
		clk	Clock Output		clk_0					
		clk_reset	Reset Output							
✓		<b>cnn_mmio_...</b>	cnn_mmio_interface							
		clock	Clock Input		clk_0					
		reset	Reset Input		[clock]					
		avalon_sla...	Avalon Memory Mapp...		[clock]	# 0x0000_0000	0x0000_7fff			
✓		<b>hps_0</b>	Arria V/Cyclone V ...							
		memory	Conduit	memory						
		hps_io	Conduit	hps						
		h2f_reset	Reset Output							
		f2h_sdram0...	Clock Input		clk_0					
		f2h_sdram0...	AXI Slave		[f2h_... #					
		h2f_axi_clock	Clock Input		clk_0					
		h2f_axi_ma...	AXI Master		[h2f_...					
		f2h_axi_clock	Clock Input		clk_0					
		f2h_axi_slave	AXI Slave		[f2h_... #					
		h2f_lw_axi...	Clock Input		clk_0					
		h2f_lw_axi...	AXI Master		[h2f_...					

Figure 12: Screenshot of Platform Designer

As shown in Fig .12, the Platform Designer system is centered on three blocks: `clk_0`, `hps_0`, and `cnn_mmio_interface`. The `clk_0` block exports the board clock and reset to the rest of the system. The `hps_0` block is the Cyclone V hard processor system, including the ARM processor, HPS DDR interface, board I/O conduits, and HPS-to-FPGA bridge ports.

The active software path uses the lightweight HPS-to-FPGA bridge. In the HPS port name `h2f_lw_axi_master`, `h2f` means that Linux on the HPS initiates the transaction and the FPGA component responds as a peripheral. The custom CNN wrapper is mapped at offset `0x0000_0000` in this bridge window, so the HPS software accesses it at physical address `0xff200000`.

The CNN wrapper itself is packaged as a 32-bit Avalon-MM slave rather than as an AXI slave. Platform Designer automatically inserts the required interconnect and protocol adapter between the HPS AXI-side master and the custom Avalon-MM slave. This lets the RTL implement a simple memory-mapped register and scratchpad interface while Linux still sees an ordinary HPS physical address.

The `f2h` ports shown on the HPS block are the reverse direction: FPGA-to-HPS. They are available on the HPS IP, but this accelerator does not use them for inference. Data is pushed into the accelerator by HPS software over H2F, and intermediate feature maps remain in FPGA on-chip memories.

### 4.2 Signal-Level MMIO Interface

The custom IP exposes a 32-bit Avalon-MM slave. The key signals are shown in Table 5.

Table 5: Custom IP Avalon-MM signal interface.

Signal	Direction	Width	Meaning
address	input	13	Word address; <code>address[12]</code> selects scratchpad memory space.
writedata	input	32	Data written by the HPS/interconnect.
byteenable	input	4	Byte write enables.
write	input	1	Write transaction strobe.
read	input	1	Read transaction strobe.
chipselct	input	1	Selected transaction qualifier.
readdata	output	32	Data returned to the HPS/interconnect.

### 4.3 Address Space

The wrapper uses `address[12]` as a space selector. The low logical space is the 16 KiB scratchpad memory. The high logical space is the config/status register file. Software uses word addresses `0x0000..0x0FFF` for scratchpad writes and `0x1000..0x101F` for the 32 active register slots.

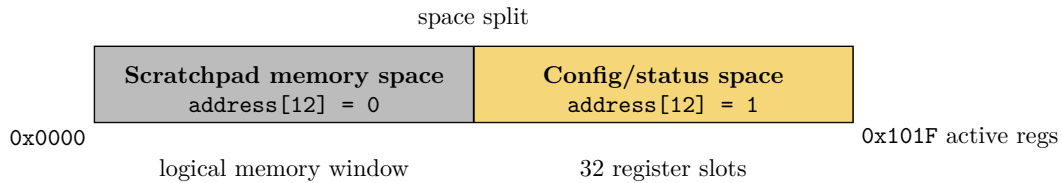


Figure 13: MMIO address-space split in the custom wrapper.

There are three sizes to distinguish:

Table 6: MMIO map sizes.

Layer	Size	Purpose
HPS userspace mapping	20 KiB	Convenient <code>/dev/mem</code> mapping window.
Scratchpad namespace	<code>0x0000..0x0FFF</code> word indices	4096 32-bit words selected by <code>address[12]=0</code> .
Config/status registers	<code>0x1000..0x101F</code> word indices	32 active register slots selected by <code>address[12]=1</code> .
Implemented scratchpad BRAM	4096 32-bit words, 16 KiB	Actual storage in <code>cnn_mmio_interface</code> .

## 4.4 Register Map

The Avalon-MM slave returns 32-bit read data. The compact control/status and base/length fields use the low 16 bits; profile counters are exposed as 32-bit values through the profile register slots. Table 7 summarizes the register interface exposed by the MMIO wrapper.

Register	Index	Access	Description
CONTROL	0	WO	start model load, start inference, clear status
STATUS	1	RO	busy flag, model-loaded flag, done flag, predicted class
CONV_CFG_BASE	2	RW	base word address of convolution configuration segment
CONV_CFG_LEN	3	RW	number of 32-bit words in convolution configuration segment
CONV_WT_BASE	4	RW	base word address of convolution weight segment
CONV_WT_LEN	5	RW	number of 32-bit words in convolution weight segment
FC_BIAS_BASE	6	RW	base word address of FC bias segment
FC_BIAS_LEN	7	RW	number of 32-bit words in FC bias segment
FCW_BASE	8	RW	base word address of FC weight segment
FCW_LEN	9	RW	number of 32-bit words in FC weight segment
IMAGE_BASE	10	RW	base word address of input image segment
IMAGE_LEN	11	RW	number of 32-bit words in input image segment
PREDICT	12	RO	latched predicted class
IF_ERROR	13	RO	interface error flags
PROFILE_*	14-29	RO	per-stage cycle counters; the low slots are reserved and the high slots hold the 32-bit counts
LAST_WRITE	30	RO	debug readback of the most recent write transaction
MAGIC	31	RO	debug constant 0x434E4E32

Table 7: MMIO register map.

## 4.5 Register Fields

Table 8 lists the meaning of each implemented control and status bit.

Register	Bits	Meaning
CONTROL	bit 0	start model preload replay
CONTROL	bit 1	start inference replay
CONTROL	bit 2	clear <code>predict_done</code> and <code>IF_ERROR</code>
CONTROL	bits 15:3	reserved, write zero
STATUS	bit 0	reserved
STATUS	bit 1	replay engine busy
STATUS	bit 2	model preload complete
STATUS	bit 3	inference complete
STATUS	bits 7:4	predicted class
STATUS	bits 15:8	reserved
PREDICT	bits 3:0	predicted class
PREDICT	bits 15:4	reserved
IF_ERROR	bit 0	model load requested while busy
IF_ERROR	bit 1	inference requested while busy
IF_ERROR	bit 2	inference requested before model load completed
IF_ERROR	bit 3	scratchpad replay read out of range
IF_ERROR	bits 15:4	reserved

Table 8: Implemented control and status bit definitions.

## 4.6 Scratchpad Layout

The HPS writes all model and image payloads into the scratchpad before triggering the replay engine. All base addresses are 32-bit word addresses.

Table 9: Default scratchpad layout for the active 10-class model.

Region	Base word	Words	Bytes	Content
<code>conv_cfg</code>	0	45	180	layer/pass configuration words
<code>conv_wt</code>	45	225	900	packed convolution weights
<code>fc_bias</code>	270	10	40	10 FC bias words
<code>fcw</code>	280	864	3456	packed fully connected weights
<code>image</code>	1144	1024	4096	packed $64 \times 64$ INT8 image



Figure 14: Contiguous 32-bit word layout used by the HPS tools. The figure is drawn for readability; exact word counts are listed in Table 9.

## 4.7 Model Preload Storage and Replay

The model preload directory contains four payloads:

- `preload_conv_cfg_45w.txt`: 45 signed 32-bit configuration words.
- `preload_conv_wt_225w_bytes.txt`: 225 packed words containing convolution weights.
- `preload_fc_bias_10w.txt`: 10 signed 32-bit FC bias words.
- `preload_fcw_864w.txt`: 864 signed 32-bit words carrying packed FC weights.

During model load, the HPS tool programs the default base/length registers, writes the four model regions into scratchpad memory, and writes `CONTROL[0]`. The wrapper then walks the configured regions in order, loads one 32-bit word at a time from scratchpad BRAM, and drives `load_valid`, `load_data`, and `load_last` into `system_top`. When the internal accelerator reports model-load completion, the wrapper sets `STATUS[2]`.

## 4.8 Image Storage And Inference Replay

The inference case directory contains the packed input source and expected class metadata. The important file for the FPGA path is `tb_conv1_in_i8_64x64x1.txt`. The board-side HPS code packs every four signed INT8 pixels into one 32-bit word. A  $64 \times 64$  image therefore occupies 4096 bytes, or 1024 32-bit scratchpad words.

For inference, the HPS tool writes the image words into the image region, writes `CONTROL[1]`, and polls `STATUS[3]`. The wrapper replays the image region into the internal load stream with inference mode selected. Layer 1 consumes the image stream, later layers operate from the accelerator’s on-chip SRAM dataflow, and the final argmax result is latched for MMIO readback.

## 4.9 Transaction Sequence

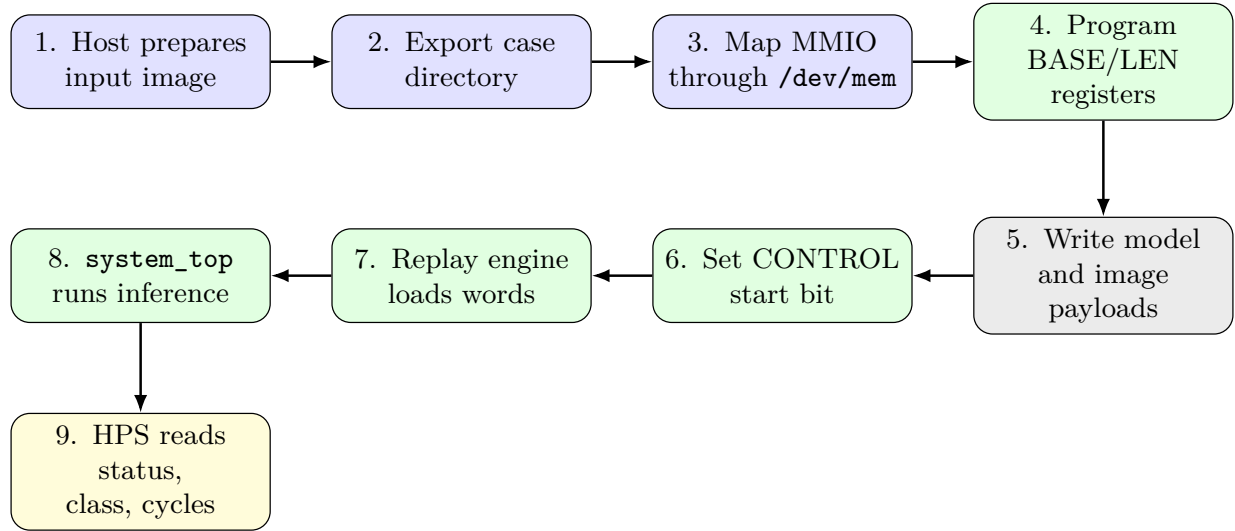


Figure 15: Software-triggered model-load and inference transaction state flow.

## 4.10 Readback Path

When inference completes, `system_top` provides the predicted class and profile counters to the MMIO wrapper. The wrapper exposes them through register reads:

- `STATUS[3]` tells software that prediction is done.
- `STATUS[7:4]` and `PREDICT[3:0]` contain the class.
- `IF_ERROR` reports incorrect command ordering or replay range errors.
- `PROFILE_*` registers report per-stage cycle counts.

The board tool prints these fields as text. The host parses the output and returns it to the Flask app, which then renders the result in the browser.

## 5 Internal Memory And Data Reuse

The accelerator avoids external DDR for intermediate feature maps. Once payloads have entered the FPGA wrapper, the CNN datapath uses local on-chip storage and deterministic streaming.

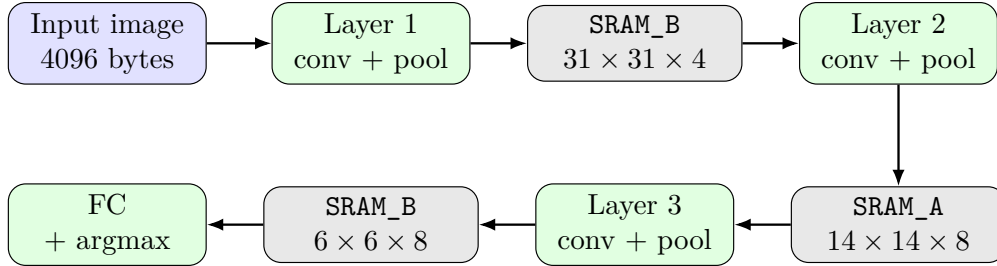


Figure 16: Internal feature-map memory reuse across CNN layers.

Table 10: Logical on-chip memories.

Memory	Width/depth	Main role	Reuse
SRAM_A	32-bit by 1024	configuration, conv weights, FC bias	layer-2 output buffer
SRAM_B	32-bit by 1024	activation feature maps	layer-1 and layer-3 pooled outputs
SRAM_FCW	80-bit by 512, 288 used	fully connected weights	ten class weights per FC input

The final feature vector has  $6 \times 6 \times 8 = 288$  INT8 values. Each FC input position reads one packed 80-bit FC weight entry containing one INT8 weight for each of the ten classes. This lets the FC datapath update all ten class accumulators in parallel.

## 6 Performance And Resource Results

### 6.1 Timing

The target fabric clock is 50 MHz. The current build meets that target with positive slack on `clock_50`.

Table 11: Current timing summary for the user fabric clock.

Corner	Metric	Value	Meaning
Slow 1100 mV 85 C	setup slack	6.297 ns	Positive at 50 MHz.
Slow 1100 mV 85 C	hold slack	0.220 ns	Positive hold margin.
Slow 1100 mV 0 C	setup slack	6.661 ns	Positive at 50 MHz.
Slow 1100 mV 85 C	Fmax	72.98 MHz	Useful headroom above 50 MHz.
Slow 1100 mV 0 C	Fmax	74.97 MHz	Useful headroom above 50 MHz.

The full SoC timing report can still show HPS DDR-related paths as the whole-design worst paths because Linux and the HPS subsystem remain in the project. For accelerator-centric timing, the relevant clock is `clock_50`. The evidence of Fmax is shown in Fig. 17.

```

+-----+
; Slow 1100mV 85C Model Fmax Summary
+-----+
; Fmax          ; Restricted Fmax ; Clock Name
+-----+
; 72.98 MHz    ; 72.98 MHz      ; clock_50
; 1184.83 MHz ; 717.36 MHz     ; soc_system:u_soc_system|soc_system_hps_0:hps_0|soc_
+-----+
This panel reports FMAX for every clock in the design, regardless of the user-specifi

```

Figure 17: Evidence of Maximum Frequency

## 6.2 Resource Utilization

The current post-fit resource usage is shown as Fig.18. The design fits the Cyclone V device, but DSP usage is the tightest resource.

```

+-----+
; Fitter Summary
+-----+
; Fitter Status          ; Successful - Mon May 11 18:26:35 2026
; Quartus Prime Version ; 21.1.0 Build 842 10/21/2021 SJ Lite Edition
; Revision Name         ; soc_system
; Top-level Entity Name ; soc_system_top
; Family                ; Cyclone V
; Device                ; 5CSEMA5F31C6
; Timing Models         ; Final
; Logic utilization (in ALMs) ; 18,357 / 32,070 ( 57 % )
; Total registers       ; 26916
; Total pins            ; 164 / 457 ( 36 % )
; Total virtual pins    ; 0
; Total block memory bits ; 241,664 / 4,065,280 ( 6 % )
; Total RAM Blocks      ; 29 / 397 ( 7 % )
; Total DSP Blocks      ; 86 / 87 ( 99 % )
; Total HSSI RX PCSs    ; 0
; Total HSSI PMA RX Deserializers ; 0
; Total HSSI TX PCSs    ; 0
; Total HSSI PMA TX Serializers ; 0
; Total PLLs            ; 0 / 6 ( 0 % )
; Total DLLs            ; 1 / 4 ( 25 % )
+-----+

```

Figure 18: Resource Utilization from Quartus Report

# 7 Demo And Validation

## 7.1 Web Demo Screenshot

### BOARD CPU VS FPGA GESTURE COMPARISON

## Upload One Gesture Image, Run Two Inference Paths

The same browser-uploaded image is preprocessed once, then compared across board-side HPS CPU inference and the DE1-SoC FPGA accelerator path.

#### Input

Upload an image from your computer or capture one from your browser webcam. Wrong predictions can be saved as correction samples for later retraining.

**Select Image**  
PNG, JPG, or other browser supported formats.

Preprocess Mode: plain, pred 2, margin 255.0, top 127.0, leaders 2:127.0, 0: 128.0 (chosen) | crop: pred 9, margin 239.0, top 118.0, leaders 9:118.0, 5: 121.0

**Browser Webcam** | Start Camera | Capture Frame

**Upload Preview** | Image selected from file upload

Camera is off

This uses your PC/browser camera only. The captured frame is sent through the same pipeline as upload mode.

Model Input (84x64) | Run Comparison | Clear

#### Correction Capture

Use this when the prediction is wrong.

Correct Gesture ID: Note

Select | Optional note about lighting, background, hand angle...

Save Correction Sample

If this result is wrong, select the correct label and save the sample for retraining.

#### Results

This panel compares the board-side HPS CPU software baseline against the existing board-side HPS + MMIO + FPGA path for the active 10-class gesture model file.

Predicted Gesture ID: **2** | Board CPU Gesture ID: **2**

Board CPU Time: **5.20 ms** | FPGA Wait Time: **1.12 ms**

FPGA RTL Time: **0.325 ms** | Speedup: **4.64x**

RTL 325.4 us from on-chip counters. Board wait 1.12 ms. Remote request 329.11 ms. Board total 4.91 ms, load 3.32 ms, program 0.42 ms.

#### CPU Confidence And Channel Scores

Top Class	Normalized Confidence	Margin
2	100.0%	255.0

2: raw 127.0, conf 100.0% | 0: raw 128.0, conf 0.0% | 1: raw 128.0, conf 0.0%

Class	Score	Confidence
0	128.0	0.0%
1	128.0	0.0%
2	127.0	100.0%
3	128.0	0.0%
4	128.0	0.0%
5	128.0	0.0%
6	128.0	0.0%
7	128.0	0.0%
8	128.0	0.0%
9	128.0	0.0%

#### FPGA RTL Profile

50.00 MHz fabric clock

Stage	Time	Time
L1	7201 cyc / 144.0 us	L2 P0: 3022 cyc / 60.4 us
L2 P1	3022 cyc / 60.4 us	L3 P0: 1202 cyc / 24.0 us
L3 P1	1202 cyc / 24.0 us	FC: 611 cyc / 12.2 us
Argmax	9 cyc / 0.18 us	
Total	16269 cyc / 325.4 us	

#### Execution Time Comparison

Board CPU: 5.20 ms

FPGA: 1.12 ms

Done: Board CPU predicted gesture ID 2, FPGA predicted gesture ID 2. Preprocess: plain.

#### Built-In Sanity Suite

Run 0-8 Suite

Runs the current web path on the repo's built-in 'digit 0' test to 'digit 9' test sample images so we can verify whether the deployed model file is internally consistent.

Not run yet.

Figure 19: Screenshot of Web Demo

The web demo is the user-facing control and visualization layer of the final system. It is implemented as a Flask application running on the development host, with browser-side JavaScript for image selection, webcam capture, result rendering, and feedback collection. The page does not emulate the FPGA result: when the user presses **Run Comparison**, the backend creates a hardware-aligned inference case, transfers it to the DE1-SoC over SSH/SCP, runs the HPS CPU baseline and the FPGA MMIO path on the board, and returns the measured results to the browser.

Table 12: Main functions provided by the web demo.

Function	User-visible behavior	Backend/system behavior
Image upload	Select a PNG/JPG gesture image from the local computer.	The browser sends the file to <code>POST /api/infer</code> ; the host preprocesses it into a $64 \times 64 \times 1$ INT8 tensor.
Webcam capture	Start the browser camera and capture one frame as the input image.	The captured frame is converted to a PNG blob and submitted through the same upload path as a file image.
Preprocess selection	Choose <b>Auto</b> , <b>Plain Resize</b> , or <b>Crop + Background Removal</b> .	The host evaluates the requested preprocessing branch. In auto mode, it compares candidate branches using the CPU model confidence and selects the stronger one.
Run comparison	Display board CPU prediction, FPGA prediction, timing, and speedup for the same input.	The backend exports a temporary case, copies it to the board, runs <code>hps_cpu_reference</code> and <code>hps_mmio_predict</code> , then parses their text output.
Input visualization	Show the original image and the exact $64 \times 64$ model input.	The backend returns a base64 preview of the final INT8 image that was packed into the FPGA scratchpad.
CPU confidence report	Show top classes, normalized confidence, raw channel scores, and margin.	The host TFLite model scores the same INT8 input and reports whether the top classes are close or tied.
FPGA profile report	Show per-stage cycle counts for L1, L2 passes, L3 passes, FC, argmax, and total cycles.	The HPS reads the MMIO <code>PROFILE_*</code> registers after inference and converts cycle counts to time using the configured 50 MHz fabric clock.
Built-in sanity suite	Run <code>digit_0_test</code> through <code>digit_9_test</code> from the repository and show pass/fail rows.	The server repeatedly invokes the same web comparison path on known sample images to check model/software/FPGA consistency.
Correction capture	Save a mispredicted input with the selected correct label and an optional note.	The Flask app stores the raw image and metadata under <code>web_demo/feedback_samples</code> as debugging evidence for the demonstrated system.

This interface is useful during bring-up because it puts several independent signals on one page: the uploaded image, the actual model input, the HPS CPU prediction, the FPGA prediction, the

HPS-observed accelerator wait time, the pure RTL cycle-counter time, and interface error flags. A mismatch between CPU and FPGA predictions usually points to a packing, preload, or MMIO replay problem, while a high `IF_ERROR` value points to command-ordering or scratchpad-range issues. The page therefore works as both a live demo and a debug dashboard for the hardware/software boundary.

## 7.2 Board Evidence

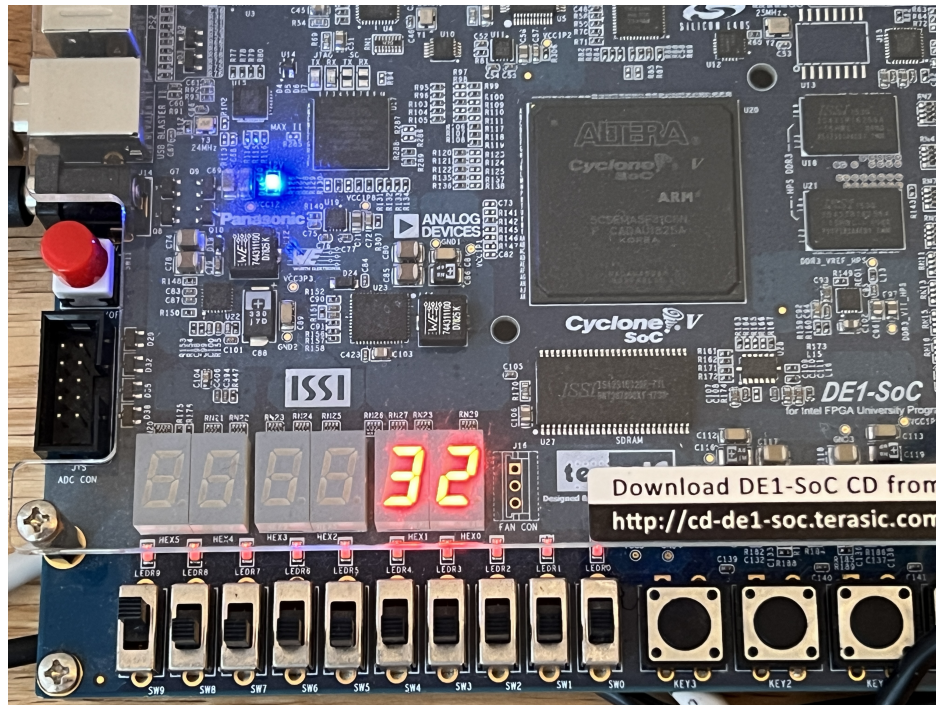


Figure 20: Status of MMIO and Results on Board

The DE1-SoC board also mirrors a small amount of accelerator state on the physical seven-segment displays. This is only a debug readout; the authoritative software result is still read from the MMIO `STATUS`, `PREDICT`, `IF_ERROR`, and `PROFILE_*` registers. The display is useful during demonstrations because it gives an immediate visual indication that the model was loaded and that inference completed.

The two main digits after a normal inference are `HEX1` and `HEX0`. `HEX0` shows the predicted gesture class, and remains blank until the accelerator asserts `predict_done`. `HEX1` shows a compact status code formed from the two debug bits `model_loaded` and `predict_done`.

Table 13: Meaning of the on-board seven-segment debug digits.

Display	Value source	Meaning
HEX0	<code>predict_class[3:0]</code>	Final predicted gesture ID, 0–9, shown after <code>predict_done=1</code> .
HEX1	<code>{model_loaded, predict_done}</code>	Status code: 0 = idle, 2 = model loaded, 3 = model loaded and inference done.
HEX2	<code>interface_error[3:0]</code>	Lowest 4 bits of the interface error code, shown only if an interface error is latched.
HEX5	constant E	Error marker, shown only if <code>interface_error != 0</code> .

For example, after a successful inference on a digit-2 input, the board shows 3 on HEX1 and 2 on HEX0. The 3 means that both `model_loaded` and `predict_done` are high, while the 2 is the accelerator’s predicted class. If an interface error occurs, HEX5 displays E and HEX2 displays the low nibble of the error word.

## 8 Team Contributions

Table 14 lists the main implementation and integration responsibilities for the project. The work was collaborative, but the table identifies the areas where each member took the lead.

Table 14: Team contribution summary.

Team member	Main area	Work completed
Cheng-wen Chu	Quantization and pooling RTL	From theory to RTL for Quantization and pooling. Implemented the INT8 requantization path after convolution, including multiplier/shift handling, bias application, saturation back to INT8, and the pooling datapath. Helped verify the convolution-to-quantization-to-pooling streaming behavior against the hardware-aligned reference outputs.
Chengcheng Xu	SRAM subsystem	Implemented and debugged the SRAM-related hardware, including <code>SRAM_A</code> , <code>SRAM_B</code> , <code>SRAM_FCW</code> , address generation, read/write control, FC-weight packing support, activation-buffer reuse, and the memory layout used across the three convolution layers and FC stage.
Harvey Lu	HW/SW interface, system flow and integration	Implemented the board-facing software and hardware/-software interface flow: the MMIO access utilities, scratchpad programming, register-map usage, model-load replay, image-case replay, readback, and host-side SSH/SCP orchestration. Implemented a C-language version CNN model for timing measurement. Also integrated the Flask web demo with the HPS tools, brought up the full host-to-board-to-FPGA path, debugged interface and packing issues, and improved the runtime flow used for final validation.
Linxiao Wu	Golden model and accelerator RTL	Built the training/export lane for the compact CNN model and the hardware-aligned golden module, including reference tensors and SRAM preload artifacts. Also implemented and integrated major accelerator RTL pieces, including layer buffering, FC datapath support, systolic-array-related modules, top-level hardware integration, debug support, slides, and report documentation.
Mingyuan Zheng	Systolic-array, CNN Implementation on Board, and demo UI	Contributed to systolic-array design and simulation. Implemented the C language version CNN model for timing measurement and designed the user-facing web-demo layout used to present image upload, board execution, and prediction results during the final demonstration.

## References

- [1] N. P. Jouppi, D. H. Yoon, M. Ashcraft, M. Gottscho, T. B. Jablin, G. Kurian, J. Laudon, S. Li, P. Ma, X. Ma, T. Norrie, N. Patil, S. Prasad, C. Young, Z. Zhou, and D. A. Patterson, “Ten lessons from three generations shaped google’s tpuv4i: Industrial product,” in *2021 ACM/IEEE 48th Annual International Symposium on Computer Architecture (ISCA)*, 2021, pp. 1–14.
- [2] C. Farabet, C. Poulet, J. Y. Han, and Y. LeCun, “Cnp: An fpga-based processor for convolutional networks,” in *2009 International Conference on Field Programmable Logic and Applications*, 2009, pp. 32–37.

## A Web-To-Board Runtime Listings

This appendix lists the source files directly used by the implemented web-to-board demonstration path and the hardware-aligned golden model used to verify it. It focuses on the Python host flow, the HPS C commands executed on the DE1-SoC, the Verilog/SystemVerilog accelerator hierarchy, and the MATLAB / PyTorch resources that generate the reference tensors and SRAM preload bundles. Generated Quartus output, generated Platform Designer output, waveform/debug dumps, collected web feedback samples, and git history are intentionally excluded from both the printed listings and the machine-readable submission archive. Auxiliary board bring-up scripts, raw register debug tools, standalone batch runners, and alternate dataset-ingestion utilities are also excluded from the printed listings because they are not required for the final web-to-board runtime or the golden-reference verification path.

### A.1 Focused Source Tree

```
code/
|-- web_demo/
|   |-- app.py
|   |-- config.py
|   |-- static/
|       |-- app.js
|       |-- styles.css
|   |-- templates/
|       |-- index.html
|-- gesture_runtime/
|   |-- board_transport.py
|   |-- case_export.py
|   |-- cpu_inference.py
|   |-- demo_compare.py
|   |-- fpga_service.py
|   |-- preprocess.py
|   |-- preprocess_selection.py
|   |-- quantization.py
|   |-- ssh_transport.py
|-- include/
|   |-- cnn_mmio_host.h
|   |-- cnn_mmio_regs.h
|-- tools/
|   |-- cnn_mmio_host.c
```

```

| |-- hps_cpu_reference.c
| |-- hps_mmio_load_model.c
| |-- hps_mmio_predict.c
| '-- hps_mmio_status.c
|-- platform_designer/
| '-- cnn_mmio_interface_hw.tcl
|-- de1_soc/
| |-- build_soc_system.py
| |-- soc_system_project.tcl
| |-- soc_system.qsys
| |-- soc_system.sdc
| '-- soc_system_top.sv
|-- Golden-Module/
| |-- models/
| | |-- v1.int8.tflite
| | '-- v1.int8.params.mat
| |-- matlab/
| | |-- digit_0_test.png ... digit_9_test.png
| | '-- hardware_aligned/
| | |-- run_all.m
| | |-- export_case.m
| | |-- hw_forward.m
| | |-- gen_sram_preload.m
| | |-- load_params.m
| | |-- dump_txt.m
| | |-- flatten_nhwc_int8.m
| | |-- maxpool2x2_int8.m
| | |-- tflite_quantize_multiplier.m
| | '-- debug/
| | |-- sram_preload/
| | | '-- digit_0_test/ ... digit_9_test/
| | '-- txt_cases/
| | '-- digit_0_test/ ... digit_9_test/
|-- pytorch/
| |-- export_tflite_params_mat.py
| '-- retrain_digits_cnn.py
'-- input/RTL/
    |-- system_top.v
    |-- interface/

```

```

|   '-- cnn_mmio_interface.v
|-- fc/
|   |-- FC.v
|   |-- fc_bias_loader.v
|   |-- fc_data_adapter.v
|   '-- mac.v
|-- fsm/
|   |-- conv_data_adapter.v
|   |-- layer_runner_fsm.v
|   |-- top_fsm.v
|   '-- wt_prepad_inserter.v
|-- conv_core/
|   |-- Conv_Buffer.v
|   |-- Line_Buffer.v
|   |-- conv_engine_ctrl.v
|   |-- conv_quant_pool.v
|   |-- conv_top.v
|   |-- input_row_aligner.v
|   |-- sa_skew_feeder.v
|   |-- systolic_array_top.v
|   '-- weight_buffer.v
|-- quant_pool/
|   |-- integration/
|   |   |-- conv_quant_adapter.v
|   |   |-- quant_param_loader.v
|   |   '-- quant_pool_adapter.v
|   |-- quant/
|   |   |-- Quantization_PE.v
|   |   '-- Quantization_Top.v
|   '-- pool/
|       |-- pool_core.v
|       '-- pool_stream_top.v
'-- SRAM/
    |-- Addr_Gen.v
    |-- fcw_preload_packer.v
    |-- sram_A_controller.v
    |-- sram_A_wrapper.v
    |-- sram_B_controller.v
    |-- sram_B_wrapper.v

```

```

|-- sram_FCW_wrapper.v
|-- top_sram_A.v
'-- top_sram_B.v

```

## A.2 Included Resource Overview

Table 15: Source and resource groups printed in this appendix.

Group	Main language / format	Role in the web-to-board path
Web and host runtime	Python / HTML / CSS / JavaScript	Receives browser uploads, preprocesses images, exports board cases, chooses preprocessing mode, runs host CPU comparison, and calls the FPGA board service.
HPS board runtime	C	Maps the FPGA control window through <code>/dev/mem</code> , writes preload and image data, triggers model load or inference, polls status, and prints results for the web server.
FPGA board boundary	SystemVerilog / Verilog	Connects the HPS Platform Designer system to the custom MMIO wrapper and passes the replayed model/image stream into the accelerator top-level control logic.
Quartus and Platform Designer setup	Tcl / Qsys / SDC / Python	Contains the custom Platform Designer component descriptor, the Qsys system, and build wrapper scripts needed to regenerate and compile the board project.
Accelerator datapath RTL	Verilog	Implements the convolution engine, quantization/pooling path, fully-connected classifier, SRAM wrappers, and layer-control FSMs driven by the board boundary.
Golden model and export resources	MATLAB / Python / TFLite / MAT / TXT	Stores the quantized model, exported parameter file, hardware-aligned MATLAB golden forward pass, reference tensors, and SRAM preload bundles used by RTL simulation, HPS commands, and web-demo comparison.

Table 16: Golden module resources used by verification and board preload.

Resource	Function
<code>code/Golden-Module/models/v1.in</code> <code>t8.tflite</code>	Deployed INT8 TFLite model used by the host-side CPU classifier and as the source for exported hardware parameters.
<code>code/Golden-Module/models/v1.in</code> <code>t8.params.mat</code>	Quantized weights, biases, zero-points, scales, multipliers, and shifts imported by the MATLAB hardware-aligned flow.
<code>code/Golden-Module/matlab/hardware_aligned/run_all.m</code>	Batch entry point that regenerates all digit reference cases and preload bundles.
<code>code/Golden-Module/matlab/hardware_aligned/export_case.m</code>	Exports one image into layer-by-layer <code>tb_*.txt</code> reference tensors plus a manifest.
<code>code/Golden-Module/matlab/hardware_aligned/hw_forward.m</code>	Bit-exact hardware-style forward pass matching the RTL data ordering and INT8/INT32 arithmetic assumptions.
<code>code/Golden-Module/matlab/hardware_aligned/gen_sram_preload.m</code>	Packs convolution configuration, convolution weights, FC weights, and FC bias into the 32-bit preload streams consumed by the HPS loader.
<code>code/Golden-Module/matlab/hardware_aligned/debug/txt_cases/</code>	Per-digit input tensors and expected intermediate/output tensors used by RTL testbenches and CPU/FPGA comparison.
<code>code/Golden-Module/matlab/hardware_aligned/debug/sram_preload/</code>	Per-digit preload bundles copied to the board before model-load replay.
<code>code/Golden-Module/pytorch/export_tflite_params_mat.py</code>	Converts the TFLite model parameters into the MAT file consumed by the golden MATLAB flow.
<code>code/Golden-Module/pytorch/retrain_digits_cnn.py</code>	Training / retraining script for the compact digit CNN model family.

Table 17: Python files used by the web and host-side runtime.

Source file	Function
<code>code/web_demo/app.py</code>	Flask routes, service construction, upload endpoint, inference endpoint, and JSON response path.
<code>code/web_demo/config.py</code>	Board SSH settings, CSR base address, model paths, label names, and fabric-clock configuration.
<code>code/gesture_runtime/demo_compare.py</code>	End-to-end comparison flow combining preprocessing, host CPU scoring, HPS CPU baseline, FPGA run, and result formatting.
<code>code/gesture_runtime/fpga_service.py</code>	Board orchestration layer that copies cases, loads the model when needed, runs HPS commands, and parses command output.
<code>code/gesture_runtime/ssh_transport.py</code>	SSH/SCP command execution and directory transfer implementation.
<code>code/gesture_runtime/board_transport.py</code>	Transport interface shared by the FPGA service and any concrete board transport.
<code>code/gesture_runtime/case_export.py</code>	Exports the preprocessed tensor into the case directory consumed by board-side commands.
<code>code/gesture_runtime/preprocess.py</code>	Converts uploaded image bytes into the 64-by-64 INT8 input tensor and preview image.
<code>code/gesture_runtime/preprocess_selection.py</code>	Selects the preprocessing branch used by automatic mode.
<code>code/gesture_runtime/quantization.py</code>	INT8 helper routines used by preprocessing and export.
<code>code/gesture_runtime/cpu_inference.py</code>	Host-side TFLite inference used for confidence reporting and preprocessing selection.

Table 18: C files used by the HPS board runtime.

Source file	Function
<code>code/include/cnn_mmio_regs.h</code>	Register ABI shared by C and RTL: offsets, bit fields, control flags, status flags, and scratchpad layout.
<code>code/include/cnn_mmio_host.h</code>	Shared C structs and function declarations for board-side MMIO commands.
<code>code/tools/cnn_mmio_host.c</code>	Common implementation for memory mapping, scratchpad writes, control writes, polling, and profile readback.
<code>code/tools/hps_mmio_status.c</code>	Lightweight health check used by the web service before deciding whether to reload the model.
<code>code/tools/hps_mmio_load_model.c</code>	Loads the preload bundle into the scratchpad and starts model replay through the MMIO control register.
<code>code/tools/hps_mmio_predict.c</code>	Writes one image case, starts FPGA inference, waits for completion, and prints machine-readable prediction output.
<code>code/tools/hps_cpu_reference.c</code>	HPS ARM software baseline used for web-demo comparison against FPGA output.

Table 19: Verilog/SystemVerilog files at the FPGA board boundary.

Source file	Function
code/de1_soc/soc_system_top.sv	DE1-SoC board wrapper connecting the generated HPS system, board pins, clocks, resets, and seven-segment display.
code/input/RTL/interface/cnn_mio_interface.v	Avalon-MM slave wrapper containing the software-visible register file, scratchpad, replay FSM, status readback, and profile readback.
code/input/RTL/system_top.v	Accelerator top-level connecting model load, CNN layer sequencing, SRAM datapath, FC classifier, argmax, and profiling outputs.
code/input/RTL/fsm/top_fsm.v	Network-level FSM for model load and inference scheduling across convolution, FC, and argmax phases.
code/input/RTL/fsm/layer_runner_fsm.v	Per-layer FSM that sequences layer configuration, weights, input data, output writes, and completion.

## A.3 Web And Host Python Runtime

### A.3.1 Flask Web Entry Point

*Defines the web routes, service construction, image upload endpoint, and JSON response path.*

code/web\_demo/app.py

```

1  #!/usr/bin/env python3
2  """Flask app for the upload-first CPU-vs-FPGA comparison demo."""
3
4  from __future__ import annotations
5
6  import os
7  import sys
8  import json
9  from datetime import datetime
10 from pathlib import Path
11 from uuid import uuid4
12
13 from flask import Flask, jsonify, render_template, request
14
15 REPO_ROOT = Path(__file__).resolve().parents[1]
16 if str(REPO_ROOT) not in sys.path:
17     sys.path.insert(0, str(REPO_ROOT))
18
19 from web_demo.config import load_demo_config
20 from gesture_runtime.cpu_inference import TFLiteCPUClassifier
21 from gesture_runtime.demo_compare import ComparisonDemoService, DemoLabels
22 from gesture_runtime.fpga_service import BoardCPUReferenceService, FPGABoardService, FPGAServiceConfig
23 from gesture_runtime.ssh_transport import SshBoardTransport
24
25
26 def build_service() -> ComparisonDemoService:

```

```

27     cfg = load_demo_config()
28     labels = DemoLabels(cfg.labels) if cfg.labels else DemoLabels.default_digits()
29     cpu_classifier = TFLiteCPUClassifier(cfg.cpu_model)
30     transport = SshBoardTransport(
31         host=cfg.board_host,
32         user=cfg.board_user,
33         port=cfg.board_port,
34         identity_file=cfg.ssh_key,
35         known_hosts_file=cfg.ssh_known_hosts,
36         bind_address=cfg.ssh_bind_address,
37     )
38     fpga_service = FPGABoardService(
39         transport=transport,
40         config=FPGAServiceConfig(
41             remote_repo=cfg.remote_repo,
42             csr_base=cfg.csr_base,
43             remote_preload_root=cfg.remote_preload_root,
44             remote_reference_case_root=cfg.remote_reference_case_root,
45         ),
46     )
47     board_cpu_service = BoardCPUReferenceService(
48         transport=transport,
49         config=FPGAServiceConfig(
50             remote_repo=cfg.remote_repo,
51             csr_base=cfg.csr_base,
52             remote_preload_root=cfg.remote_preload_root,
53             remote_reference_case_root=cfg.remote_reference_case_root,
54         ),
55     )
56     return ComparisonDemoService(
57         cpu_classifier=cpu_classifier,
58         board_cpu_service=board_cpu_service,
59         fpga_service=fpga_service,
60         labels=labels,
61         fabric_mhz=cfg.fabric_mhz,
62     )
63
64
65 app = Flask(__name__, template_folder="templates", static_folder="static")
66 SERVICE: ComparisonDemoService | None = None
67 SAMPLE_IMAGE_DIR = REPO_ROOT / "Golden-Module" / "matlab"
68 FEEDBACK_DIR = REPO_ROOT / "web_demo" / "feedback_samples"
69
70
71 def get_service() -> ComparisonDemoService:
72     global SERVICE
73     if SERVICE is None:
74         SERVICE = build_service()
75     return SERVICE
76
77
78 @app.route("/")
79 def index():
80     return render_template("index.html")

```

```

81
82
83 @app.get("/api/health")
84 def health():
85     cfg = load_demo_config()
86     return jsonify(
87         {
88             "ok": True,
89             "board_host": cfg.board_host,
90             "board_user": cfg.board_user,
91             "board_port": cfg.board_port,
92             "ssh_bind_address": cfg.ssh_bind_address,
93             "remote_repo": cfg.remote_repo,
94             "remote_preload_root": cfg.remote_preload_root,
95             "remote_reference_case_root": cfg.remote_reference_case_root,
96             "cpu_model": str(cfg.cpu_model),
97             "fabric_mhz": cfg.fabric_mhz,
98             "labels": cfg.labels,
99             "model_lane": cfg.model_lane,
100            "model_warning": cfg.model_warning,
101            "task_semantics": "gesture-classification",
102            "dataset_name": cfg.dataset_name,
103            "dataset_source_hint": cfg.dataset_source_hint,
104        }
105    )
106
107
108 @app.post("/api/infer")
109 def infer():
110     uploaded = request.files.get("image")
111     if uploaded is None or not uploaded.filename:
112         return jsonify({"error": "image upload is required"}), 400
113
114     preprocess_mode = request.form.get("preprocess_mode", "auto").strip() or "auto"
115     try:
116         cfg = load_demo_config()
117         result = get_service().compare_image_bytes(uploaded.read(), preferred_mode=preprocess_mode)
118     except Exception as exc: # pragma: no cover - integration path
119         return jsonify({"error": str(exc)}), 500
120     result["model_lane"] = cfg.model_lane
121     result["model_warning"] = cfg.model_warning
122     result["task_semantics"] = "gesture-classification"
123     result["dataset_name"] = cfg.dataset_name
124     result["dataset_source_hint"] = cfg.dataset_source_hint
125     return jsonify(result)
126
127
128 @app.post("/api/feedback")
129 def save_feedback():
130     uploaded = request.files.get("image")
131     corrected_label = request.form.get("corrected_label", "").strip()
132     predicted_label = request.form.get("predicted_label", "").strip()
133     preprocess_mode = request.form.get("preprocess_mode", "").strip()
134     note = request.form.get("note", "").strip()

```

```

135
136 if uploaded is None or not uploaded.filename:
137     return jsonify({"error": "image upload is required"}), 400
138 if corrected_label not in {str(i) for i in range(10)}:
139     return jsonify({"error": "corrected_label must be one of 0..9"}), 400
140
141 stamp = datetime.now().strftime("%Y%m%d_%H%M%S")
142 sample_id = f"{stamp}_{uuid4().hex[:8]}"
143 sample_dir = FEEDBACK_DIR / sample_id
144 sample_dir.mkdir(parents=True, exist_ok=True)
145
146 filename = Path(uploaded.filename).name or "capture.png"
147 raw_ext = Path(filename).suffix or ".png"
148 raw_path = sample_dir / f"input{raw_ext}"
149 uploaded.save(raw_path)
150
151 metadata = {
152     "id": sample_id,
153     "captured_at": datetime.now().isoformat(),
154     "source_filename": filename,
155     "raw_image": str(raw_path.relative_to(REPO_ROOT)),
156     "corrected_label": corrected_label,
157     "predicted_label": predicted_label,
158     "preprocess_mode": preprocess_mode,
159     "note": note,
160 }
161 metadata_path = sample_dir / "metadata.json"
162 metadata_path.write_text(json.dumps(metadata, indent=2), encoding="utf-8")
163
164 return jsonify(
165     {
166         "ok": True,
167         "sample_id": sample_id,
168         "saved_dir": str(sample_dir.relative_to(REPO_ROOT)),
169         "corrected_label": corrected_label,
170     }
171 )
172
173
174 @app.get("/api/sample-suite")
175 def sample_suite():
176     try:
177         cfg = load_demo_config()
178         sample_paths = sorted(SAMPLE_IMAGE_DIR.glob("digit_*_test.png"))
179         if not sample_paths:
180             return jsonify({"error": f"no sample images found under {SAMPLE_IMAGE_DIR}"}), 500
181
182         rows = []
183         passed = 0
184         service = get_service()
185         for sample_path in sample_paths:
186             stem = sample_path.stem
187             expected = int(stem.split("_")[1])
188             result = service.compare_image_path(sample_path)

```

```

189         cpu_ok = int(result["cpu_predicted_class"]) == expected
190         fpga_ok = int(result["predicted_class"]) == expected
191         both_ok = cpu_ok and fpga_ok
192         if both_ok:
193             passed += 1
194         rows.append(
195             {
196                 "case": stem,
197                 "expected_class": expected,
198                 "cpu_predicted_class": int(result["cpu_predicted_class"]),
199                 "fpga_predicted_class": int(result["predicted_class"]),
200                 "cpu_time_ms": result["cpu_time_ms"],
201                 "fpga_time_ms": result["fpga_time_ms"],
202                 "preprocess_mode": result.get("preprocess_mode", ""),
203                 "cpu_confidence_margin": result.get("cpu_confidence_margin"),
204                 "cpu_ok": cpu_ok,
205                 "fpga_ok": fpga_ok,
206                 "pass": both_ok,
207             }
208         )
209
210     return jsonify(
211         {
212             "ok": True,
213             "dataset_name": cfg.dataset_name,
214             "model_lane": cfg.model_lane,
215             "task_semantics": "gesture-classification",
216             "total": len(rows),
217             "passed": passed,
218             "rows": rows,
219         }
220     )
221     except Exception as exc: # pragma: no cover - integration path
222         return jsonify({"error": str(exc)}), 500
223
224
225 if __name__ == "__main__":
226     cfg = load_demo_config()
227     app.run(host=cfg.local_host, port=cfg.local_port, debug=cfg.debug)

```

### A.3.2 Web Demo Configuration

*Loads board SSH settings, CSR base address, model paths, labels, and fabric-clock configuration.*

```
code/web_demo/config.py
```

```

1 """Configuration loading for the web comparison demo."""
2
3 from __future__ import annotations
4
5 import json

```

```

6 import os
7 from dataclasses import dataclass
8 from pathlib import Path
9 from typing import List, Optional
10
11
12 REPO_ROOT = Path(__file__).resolve().parents[1]
13 DEFAULT_MODEL = REPO_ROOT / "Golden-Module" / "models" / "v1.int8.tflite"
14 DEFAULT_CONFIG_PATH = Path(
15     os.environ.get("CNN_ACC_DEMO_CONFIG", str(REPO_ROOT / "web_demo" / "demo_config.json"))
16 )
17
18
19 @dataclass
20 class DemoConfig:
21     cpu_model: Path
22     labels: List[str]
23     board_host: str
24     board_user: str
25     board_port: int
26     ssh_key: Optional[Path]
27     ssh_known_hosts: Optional[Path]
28     ssh_bind_address: Optional[str]
29     remote_repo: str
30     remote_preload_root: str
31     remote_reference_case_root: str
32     csr_base: str
33     fabric_mhz: float
34     local_host: str
35     local_port: int
36     debug: bool
37
38     @property
39     def model_lane(self) -> str:
40         preload = self.remote_preload_root.lower()
41         labels_are_numeric = self.labels and all(label.isdigit() for label in self.labels)
42         if "digit_" in preload or labels_are_numeric:
43             return "gesture-10class"
44         if any(token in preload for token in ("paper", "rock", "scissors", "gesture")):
45             return "gesture"
46         return "unknown"
47
48     @property
49     def model_warning(self) -> Optional[str]:
50         if self.model_lane == "unknown":
51             return (
52                 "Current demo config does not clearly map to a gesture-specific deployed model. "
53                 "Predictions may not match the intended gesture task."
54             )
55         return None
56
57     @property
58     def dataset_name(self) -> str:
59         if self.model_lane == "gesture-10class":

```

```

60         return "Sign Language Digits Dataset"
61     if self.model_lane == "gesture":
62         return "Gesture dataset"
63     return "Unknown dataset"
64
65     @property
66     def dataset_source_hint(self) -> Optional[str]:
67         if self.model_lane == "gesture-10class":
68             return "Kaggle sign-language digits dataset family; local training root expected as
dataset_wlx/"
69         return None
70
71
72     def _read_json_config(path: Path) -> dict:
73         if not path.is_file():
74             return {}
75         return json.loads(path.read_text(encoding="utf-8"))
76
77
78     def _resolve_repo_path(value: Optional[str]) -> Optional[Path]:
79         if not value:
80             return None
81         raw = Path(value)
82         if raw.is_absolute():
83             return raw
84         return (REPO_ROOT / raw).resolve()
85
86
87     def load_demo_config(path: Optional[Path] = None) -> DemoConfig:
88         cfg_path = path or DEFAULT_CONFIG_PATH
89         raw = _read_json_config(cfg_path)
90
91         labels = raw.get("labels") or os.environ.get("CNN_ACC_LABELS", "")
92         label_list = labels if isinstance(labels, list) else [x.strip() for x in labels.split(",") if x.strip
()]
93         if not label_list:
94             label_list = [str(i) for i in range(10)]
95
96         ssh_key = raw.get("ssh_key") or os.environ.get("CNN_ACC_SSH_KEY")
97         ssh_known_hosts = raw.get("ssh_known_hosts") or os.environ.get("CNN_ACC_SSH_KNOWN_HOSTS")
98         ssh_bind_address = raw.get("ssh_bind_address") or os.environ.get("CNN_ACC_SSH_BIND_ADDRESS")
99
100         return DemoConfig(
101             cpu_model=_resolve_repo_path(
102                 raw.get("cpu_model") or os.environ.get("CNN_ACC_CPU_MODEL", str(DEFAULT_MODEL))
103             )
104             or DEFAULT_MODEL,
105             labels=label_list,
106             board_host=raw.get("board_host") or os.environ.get("CNN_ACC_BOARD_HOST", "192.168.0.2"),
107             board_user=raw.get("board_user") or os.environ.get("CNN_ACC_BOARD_USER", "root"),
108             board_port=int(raw.get("board_port") or os.environ.get("CNN_ACC_BOARD_PORT", "22")),
109             ssh_key=_resolve_repo_path(ssh_key),
110             ssh_known_hosts=_resolve_repo_path(ssh_known_hosts),
111             ssh_bind_address=ssh_bind_address or None,

```

```

112     remote_repo=raw.get("remote_repo") or os.environ.get("CNN_ACC_REMOTE_REPO", "/root/cnn_acc_hps"),
113     remote_preload_root=raw.get("remote_preload_root")
114     or os.environ.get(
115         "CNN_ACC_REMOTE_PRELOAD_ROOT",
116         "Golden-Module/matlab/hardware_aligned/debug/sram_preload/digit_0_test",
117     ),
118     remote_reference_case_root=raw.get("remote_reference_case_root")
119     or os.environ.get(
120         "CNN_ACC_REMOTE_REFERENCE_CASE_ROOT",
121         "Golden-Module/matlab/hardware_aligned/debug/txt_cases/digit_0_test",
122     ),
123     csr_base=raw.get("csr_base") or os.environ.get("CNN_ACC_CSR_BASE", "0xff200000"),
124     fabric_mhz=float(raw.get("fabric_mhz") or os.environ.get("CNN_ACC_FABRIC_MHZ", "50")),
125     local_host=raw.get("local_host") or os.environ.get("HOST", "127.0.0.1"),
126     local_port=int(raw.get("local_port") or os.environ.get("PORT", "5000")),
127     debug=bool(raw.get("debug", os.environ.get("FLASK_DEBUG", "1") not in ("0", "false", "False"))),
128 )

```

### A.3.3 CPU-vs-FPGA Comparison Flow

*Coordinates preprocessing, CPU scoring, board CPU reference execution, FPGA execution, and result formatting.*

code/gesture\_runtime/demo\_compare.py

```

1  """Unified CPU-vs-FPGA comparison flow for the web demo."""
2
3  from __future__ import annotations
4
5  import math
6  import tempfile
7  from dataclasses import dataclass
8  from pathlib import Path
9  from typing import Dict, List
10
11 from .case_export import write_inference_case
12 from .cpu_inference import TFLiteCPUClassifier
13 from .fpga_service import BoardCPUReferenceService, FPGABoardService
14 from .preprocess import preprocess_image_bytes_variants
15 from .preprocess_selection import choose_best_variant
16
17
18 @dataclass
19 class DemoLabels:
20     labels: List[str]
21
22     @classmethod
23     def default_digits(cls) -> "DemoLabels":
24         return cls(labels=[str(i) for i in range(10)])
25

```

```

26     def label_for(self, idx: int) -> str:
27         if 0 <= idx < len(self.labels):
28             return self.labels[idx]
29         return str(idx)
30
31
32 class ComparisonDemoService:
33     def __init__(
34         self,
35         cpu_classifier: TFLiteCPUClassifier,
36         board_cpu_service: BoardCPUReferenceService,
37         fpga_service: FPGABoardService,
38         labels: DemoLabels,
39         fabric_mhz: float = 25.0,
40     ):
41         self.cpu_classifier = cpu_classifier
42         self.board_cpu_service = board_cpu_service
43         self.fpga_service = fpga_service
44         self.labels = labels
45         self.fabric_mhz = fabric_mhz
46
47     def _cycles_to_us(self, cycles: int) -> float:
48         if self.fabric_mhz <= 0:
49             return 0.0
50         return float(cycles) / self.fabric_mhz
51
52     def _score_variants(self, variants):
53         scored = []
54         for variant in variants:
55             cpu_result = self.cpu_classifier.predict_int8_image(variant.values)
56             scored.append((variant, cpu_result))
57         return scored
58
59     def _softmax_probabilities(self, scores: List[float]) -> List[float]:
60         if not scores:
61             return []
62         max_score = max(scores)
63         exp_scores = [math.exp(score - max_score) for score in scores]
64         total = sum(exp_scores)
65         if total <= 0:
66             return [0.0 for _ in scores]
67         return [value / total for value in exp_scores]
68
69     def _build_cpu_score_report(self, cpu_result) -> Dict[str, object]:
70         probabilities = self._softmax_probabilities(cpu_result.scores)
71         channels = []
72         for idx, (score, probability) in enumerate(zip(cpu_result.scores, probabilities)):
73             channels.append(
74                 {
75                     "class_id": idx,
76                     "label": self.labels.label_for(idx),
77                     "raw_score": float(score),
78                     "probability": float(probability),
79                     "probability_pct": float(probability * 100.0),

```

```

80         }
81     )
82
83     top_channels = sorted(channels, key=lambda row: row["raw_score"], reverse=True)[:3]
84     note = ""
85     if len(top_channels) >= 2 and abs(top_channels[0]["raw_score"] - top_channels[1]["raw_score"]) <
1e-6:
86         note = "Top classes are tied, so this sample is currently low-confidence."
87     elif float(cpu_result.margin) < 5.0:
88         note = "Top classes are close, so this sample is currently low-confidence."
89
90     return {
91         "predicted_class": int(cpu_result.predicted_class),
92         "predicted_label": self.labels.label_for(int(cpu_result.predicted_class)),
93         "top_score": float(cpu_result.top_score),
94         "margin": float(cpu_result.margin),
95         "confidence": float(top_channels[0]["probability"]) if top_channels else 0.0,
96         "confidence_pct": float(top_channels[0]["probability_pct"]) if top_channels else 0.0,
97         "top_channels": top_channels,
98         "channels": channels,
99         "note": note,
100     }
101
102 def _build_branch_diagnostics(self, scored) -> List[Dict[str, object]]:
103     diagnostics = []
104     for variant, cpu_result in scored:
105         top_channels = self._build_cpu_score_report(cpu_result)["top_channels"]
106         diagnostics.append(
107             {
108                 "mode": variant.mode,
109                 "predicted_class": int(cpu_result.predicted_class),
110                 "predicted_label": self.labels.label_for(int(cpu_result.predicted_class)),
111                 "margin": float(cpu_result.margin),
112                 "top_score": float(cpu_result.top_score),
113                 "top_channels": top_channels,
114             }
115         )
116     return diagnostics
117
118 def _build_fpga_profile(self, fpga_result: Dict[str, object]) -> Dict[str, object]:
119     raw_profile = fpga_result.get("profile", {})
120     stage_specs = [
121         ("L1", "l1_cycles"),
122         ("L2 P0", "l2_p0_cycles"),
123         ("L2 P1", "l2_p1_cycles"),
124         ("L3 P0", "l3_p0_cycles"),
125         ("L3 P1", "l3_p1_cycles"),
126         ("FC", "fc_cycles"),
127         ("Argmax", "argmax_cycles"),
128     ]
129     stages = []
130     for label, key in stage_specs:
131         cycles = int(raw_profile.get(key, 0))
132         stages.append(

```

```

133         {
134             "label": label,
135             "cycles": cycles,
136             "time_us": self._cycles_to_us(cycles),
137         }
138     )
139     total_cycles = int(raw_profile.get("total_cycles", 0))
140     return {
141         "fabric_mhz": self.fabric_mhz,
142         "total_cycles": total_cycles,
143         "rtl_time_us": self._cycles_to_us(total_cycles),
144         "rtl_time_ms": self._cycles_to_us(total_cycles) / 1000.0,
145         "stages": stages,
146     }
147
148 def _compare_preprocessed_image(self, image_values: List[int], preview_b64: str) -> Dict[str, object
149 ]:
150     cpu_result = self.cpu_classifier.predict_int8_image(image_values)
151
152     with tempfile.TemporaryDirectory(prefix="cnn_acc_case_") as tmp_dir:
153         case_root = Path(tmp_dir) / "gesture_upload_case"
154         write_inference_case(case_root, image_values, case_name="gesture_upload", expected_class=-1)
155         board_cpu_result = self.board_cpu_service.predict_case_dir(case_root)
156         fpga_result = self.fpga_service.predict_case_dir(case_root)
157
158     cpu_ms = float(board_cpu_result["board_cpu_infer_ms"])
159     fpga_kernel_ms = float(fpga_result["board_infer_wait_ms"])
160     fpga_request_ms = float(fpga_result["request_elapsed_ms"])
161     fpga_profile = self._build_fpga_profile(fpga_result)
162     speedup = (cpu_ms / fpga_kernel_ms) if fpga_kernel_ms > 0 else None
163
164     cpu_class = int(board_cpu_result["predict_class"])
165     fpga_class = int(fpga_result["predict_class"])
166     return {
167         "predicted_class": fpga_class,
168         "predicted_label": self.labels.label_for(fpga_class),
169         "cpu_predicted_class": cpu_class,
170         "cpu_predicted_label": self.labels.label_for(cpu_class),
171         "cpu_time_ms": cpu_ms,
172         "fpga_time_ms": fpga_kernel_ms,
173         "fpga_rtl_time_ms": fpga_profile["rtl_time_ms"],
174         "fpga_rtl_time_us": fpga_profile["rtl_time_us"],
175         "fpga_request_time_ms": fpga_request_ms,
176         "fpga_board_total_time_ms": float(fpga_result["board_total_ms"]),
177         "fpga_board_program_time_ms": float(fpga_result["board_program_ms"]),
178         "fpga_board_case_load_time_ms": float(fpga_result["board_case_load_ms"]),
179         "speedup": speedup,
180         "fpga_status": fpga_result["status"],
181         "fpga_error": fpga_result["error"],
182         "chart": {
183             "labels": ["CPU", "FPGA wait"],
184             "values_ms": [cpu_ms, fpga_kernel_ms],
185         },
186         "fpga_profile": fpga_profile,

```

```

186         "preprocessed_preview_url": f"data:image/png;base64,{preview_b64}",
187     }
188
189 def compare_image_bytes(self, image_bytes: bytes, preferred_mode: str = "auto") -> Dict[str, object]:
190     variants = preprocess_image_bytes_variants(image_bytes)
191     scored = self._score_variants(variants)
192
193     if preferred_mode != "auto":
194         matches = [(variant, cpu_result) for variant, cpu_result in scored if variant.mode ==
preferred_mode]
195         if not matches:
196             raise ValueError(f"unknown preprocess mode: {preferred_mode}")
197         best_variant, best_cpu_result = matches[0]
198     else:
199         best_variant, best_cpu_result = choose_best_variant(scored)
200
201     cpu_score_report = self._build_cpu_score_report(best_cpu_result)
202
203     with tempfile.TemporaryDirectory(prefix="cnn_acc_case_") as tmp_dir:
204         case_root = Path(tmp_dir) / "gesture_upload_case"
205         write_inference_case(case_root, best_variant.values, case_name="gesture_upload",
expected_class=-1)
206         board_cpu_result = self.board_cpu_service.predict_case_dir(case_root)
207         fpga_result = self.fpga_service.predict_case_dir(case_root)
208
209     cpu_ms = float(board_cpu_result["board_cpu_infer_ms"])
210     fpga_kernel_ms = float(fpga_result["board_infer_wait_ms"])
211     fpga_request_ms = float(fpga_result["request_elapsed_ms"])
212     fpga_profile = self._build_fpga_profile(fpga_result)
213     speedup = (cpu_ms / fpga_kernel_ms) if fpga_kernel_ms > 0 else None
214
215     cpu_class = int(board_cpu_result["predict_class"])
216     fpga_class = int(fpga_result["predict_class"])
217     return {
218         "predicted_class": fpga_class,
219         "predicted_label": self.labels.label_for(fpga_class),
220         "cpu_predicted_class": cpu_class,
221         "cpu_predicted_label": self.labels.label_for(cpu_class),
222         "cpu_time_ms": cpu_ms,
223         "fpga_time_ms": fpga_kernel_ms,
224         "fpga_rtl_time_ms": fpga_profile["rtl_time_ms"],
225         "fpga_rtl_time_us": fpga_profile["rtl_time_us"],
226         "fpga_request_time_ms": fpga_request_ms,
227         "fpga_board_total_time_ms": float(fpga_result["board_total_ms"]),
228         "fpga_board_program_time_ms": float(fpga_result["board_program_ms"]),
229         "fpga_board_case_load_time_ms": float(fpga_result["board_case_load_ms"]),
230         "speedup": speedup,
231         "fpga_status": fpga_result["status"],
232         "fpga_error": fpga_result["error"],
233         "chart": {
234             "labels": ["CPU", "FPGA wait"],
235             "values_ms": [cpu_ms, fpga_kernel_ms],
236         },
237         "fpga_profile": fpga_profile,

```

```

238         "branch_diagnostics": self._build_branch_diagnostics(scored),
239         "preprocessed_preview_url": f"data:image/png;base64,{best_variant.preview_b64}",
240         "preprocess_mode": best_variant.mode,
241         "cpu_confidence_margin": best_cpu_result.margin,
242         "cpu_score_report": cpu_score_report,
243     }
244
245     def compare_image_path(self, image_path: Path) -> Dict[str, object]:
246         image_bytes = Path(image_path).read_bytes()
247         return self.compare_image_bytes(image_bytes)

```

### A.3.4 FPGA Board Service

Uses SSH/SCP to load the model when needed, copy an exported case to the board, run HPS tools, and parse key-value output.

code/gesture\_runtime/fpga\_service.py

```

1  """Board-side FPGA inference service wrappers."""
2
3  from __future__ import annotations
4
5  import shlex
6  import tempfile
7  from dataclasses import dataclass
8  from pathlib import Path
9  from time import perf_counter
10 from typing import Dict
11
12 from .board_transport import BoardTransport
13
14
15 def parse_key_value_output(text: str) -> Dict[str, str]:
16     result: Dict[str, str] = {}
17     for raw_line in text.splitlines():
18         line = raw_line.strip()
19         if not line or "=" not in line:
20             continue
21         key, value = line.split("=", 1)
22         result[key.strip()] = value.strip()
23     return result
24
25
26 @dataclass
27 class FPGAServiceConfig:
28     remote_repo: str = "/root/cnn_acc_hps"
29     csr_base: str = "0xff200000"
30     remote_preload_root: str = "Golden-Module/matlab/hardware_aligned/debug/sram_preload/digit_0_test"
31     remote_reference_case_root: str = "Golden-Module/matlab/hardware_aligned/debug/txt_cases/digit_0_test"

```

```

32     remote_case_parent: str = "/tmp"
33
34
35 class FPGABoardService:
36     def __init__(self, transport: BoardTransport, config: FPGAServiceConfig):
37         self.transport = transport
38         self.config = config
39         self._model_loaded = False
40
41     def _read_status(self) -> Dict[str, str]:
42         cmd = (
43             f"cd {shlex.quote(self.config.remote_repo)} && "
44             f"./tools/hps_mmio_status {shlex.quote(self.config.csr_base)}"
45         )
46         output = self.transport.run(cmd, timeout_s=20.0)
47         return parse_key_value_output(output)
48
49     def ensure_model_loaded(self, force: bool = False) -> None:
50         if not force and self._model_loaded:
51             try:
52                 status = self._read_status()
53                 if status.get("model_loaded") == "1":
54                     return
55             except Exception:
56                 # Fall back to reloading the model if a status probe fails.
57                 pass
58         cmd = (
59             f"cd {shlex.quote(self.config.remote_repo)} && "
60             f"./tools/hps_mmio_load_model {shlex.quote(self.config.csr_base)} "
61             f"{shlex.quote(self.config.remote_preload_root)}"
62         )
63         self.transport.run(cmd, timeout_s=60.0)
64         self._model_loaded = True
65
66     def predict_case_dir(self, local_case_dir: Path) -> Dict[str, object]:
67         self.ensure_model_loaded()
68         self.transport.put_dir(local_case_dir, self.config.remote_case_parent, timeout_s=60.0)
69
70         remote_case_dir = f"{self.config.remote_case_parent}/{local_case_dir.name}"
71         cmd = (
72             f"cd {shlex.quote(self.config.remote_repo)} && "
73             f"./tools/hps_mmio_predict {shlex.quote(self.config.csr_base)} "
74             f"{shlex.quote(remote_case_dir)}"
75         )
76
77         started = perf_counter()
78         try:
79             output = self.transport.run(cmd, timeout_s=60.0)
80         except Exception as exc:
81             message = str(exc)
82             if "timeout waiting for predict_done" not in message:
83                 raise
84             self.ensure_model_loaded(force=True)
85             output = self.transport.run(cmd, timeout_s=60.0)

```

```

86     elapsed_ms = (perf_counter() - started) * 1000.0
87     parsed = parse_key_value_output(output)
88     profile = {
89         "l1_cycles": int(parsed.get("l1_cycles", "0") or 0),
90         "l2_p0_cycles": int(parsed.get("l2_p0_cycles", "0") or 0),
91         "l2_p1_cycles": int(parsed.get("l2_p1_cycles", "0") or 0),
92         "l3_p0_cycles": int(parsed.get("l3_p0_cycles", "0") or 0),
93         "l3_p1_cycles": int(parsed.get("l3_p1_cycles", "0") or 0),
94         "fc_cycles": int(parsed.get("fc_cycles", "0") or 0),
95         "argmax_cycles": int(parsed.get("argmax_cycles", "0") or 0),
96         "total_cycles": int(parsed.get("total_cycles", "0") or 0),
97     }
98     return {
99         "predict_class": int(parsed.get("predict_class", "0")),
100        "status": parsed.get("status", ""),
101        "error": parsed.get("error", ""),
102        "request_elapsed_ms": elapsed_ms,
103        "board_case_load_ms": float(parsed.get("board_case_load_ms", "0") or 0),
104        "board_program_ms": float(parsed.get("board_program_ms", "0") or 0),
105        "board_infer_wait_ms": float(parsed.get("board_infer_wait_ms", "0") or 0),
106        "board_total_ms": float(parsed.get("board_total_ms", "0") or 0),
107        "profile": profile,
108        "raw_output": output,
109    }
110
111
112 class BoardCPUReferenceService:
113     def __init__(self, transport: BoardTransport, config: FPGAServiceConfig):
114         self.transport = transport
115         self.config = config
116
117     def predict_case_dir(self, local_case_dir: Path) -> Dict[str, object]:
118         self.transport.put_dir(local_case_dir, self.config.remote_case_parent, timeout_s=60.0)
119
120         remote_case_dir = f"{self.config.remote_case_parent}/{local_case_dir.name}"
121         cmd = (
122             f"cd {shlex.quote(self.config.remote_repo)} && "
123             f"./tools/hps_cpu_reference {shlex.quote(remote_case_dir)} "
124             f"{shlex.quote(self.config.remote_reference_case_root)}"
125         )
126         output = self.transport.run(cmd, timeout_s=60.0)
127         parsed = parse_key_value_output(output)
128         return {
129             "predict_class": int(parsed.get("predict_class", "0")),
130             "board_cpu_infer_ms": float(parsed.get("board_cpu_infer_ms", "0") or 0),
131             "raw_output": output,
132         }

```

### A.3.5 SSH Transport

*Implements the command execution and directory transfer layer used by the board service.*

```

1  """SSH transport for host-to-board execution."""
2
3  from __future__ import annotations
4
5  import subprocess
6  from pathlib import Path
7  from typing import List, Optional
8
9  from .board_transport import BoardTransport, TransportError
10
11
12 class SshBoardTransport(BoardTransport):
13     def __init__(
14         self,
15         host: str,
16         user: str = "root",
17         port: int = 22,
18         identity_file: Optional[Path] = None,
19         known_hosts_file: Optional[Path] = None,
20         bind_address: Optional[str] = None,
21     ):
22         self.host = host
23         self.user = user
24         self.port = port
25         self.identity_file = Path(identity_file) if identity_file else None
26         self.known_hosts_file = Path(known_hosts_file) if known_hosts_file else None
27         self.bind_address = bind_address
28
29     @property
30     def target(self) -> str:
31         return f"{self.user}@{self.host}"
32
33     def _base_args(self) -> List[str]:
34         args = [
35             "-p",
36             str(self.port),
37             "-o",
38             "StrictHostKeyChecking=no",
39         ]
40         if self.identity_file:
41             args.extend(["-i", str(self.identity_file)])
42         if self.known_hosts_file:
43             args.extend(["-o", f"UserKnownHostsFile={self.known_hosts_file}"])
44         if self.bind_address:
45             args.extend(["-o", f"BindAddress={self.bind_address}"])
46         return args
47
48     def run(self, command: str, timeout_s: float = 30.0) -> str:
49         proc = subprocess.run(
50             ["ssh", *self._base_args(), self.target, command],
51             capture_output=True,
52             text=True,

```

```

53         encoding="utf-8",
54         errors="replace",
55         timeout=timeout_s,
56     )
57     if proc.returncode != 0:
58         raise TransportError(
59             "ssh command failed\n"
60             f"command: {command}\n"
61             f"stdout: \n{proc.stdout}\n"
62             f"stderr: \n{proc.stderr}"
63         )
64     return proc.stdout
65
66 def put_dir(self, local_dir: Path, remote_parent: str, timeout_s: float = 30.0) -> None:
67     proc = subprocess.run(
68         [
69             "scp",
70             "-r",
71             "-P",
72             str(self.port),
73             "-o",
74             "StrictHostKeyChecking=no",
75             *([ "-i", str(self.identity_file)] if self.identity_file else []),
76             *([ "-o", f"UserKnownHostsFile={self.known_hosts_file}"] if self.known_hosts_file else []),
77             *([ "-o", f"BindAddress={self.bind_address}"] if self.bind_address else []),
78             str(local_dir),
79             f"{self.target}:{remote_parent}",
80         ],
81         capture_output=True,
82         text=True,
83         encoding="utf-8",
84         errors="replace",
85         timeout=timeout_s,
86     )
87     if proc.returncode != 0:
88         raise TransportError(
89             "scp failed\n"
90             f"local_dir: {local_dir}\n"
91             f"remote_parent: {remote_parent}\n"
92             f"stdout: \n{proc.stdout}\n"
93             f"stderr: \n{proc.stderr}"
94     )

```

### A.3.6 Board Transport Interface

*Defines the abstract transport contract shared by SSH and other possible board transports.*

```
code/gesture_runtime/board_transport.py
```

```
1 """Transport abstraction for host-to-board execution."""
2
3 from __future__ import annotations
4
5 from abc import ABC, abstractmethod
6 from pathlib import Path
7 from typing import Optional
8
9
10 class BoardTransport(ABC):
11     @abstractmethod
12     def run(self, command: str, timeout_s: float = 30.0) -> str:
13         raise NotImplementedError
14
15     @abstractmethod
16     def put_dir(self, local_dir: Path, remote_parent: str, timeout_s: float = 30.0) -> None:
17         raise NotImplementedError
18
19
20 class TransportError(RuntimeError):
21     pass
```

### A.3.7 Inference Case Export

*Writes a preprocessed image into the hardware-aligned case directory consumed by the HPS tools.*

```
code/gesture_runtime/case_export.py
```

```
1 """Helpers for exporting one uploaded image into a board-consumable case folder."""
2
3 from pathlib import Path
4 from typing import Iterable
5
6 from .preprocess import write_txt_int8_image
7
8
9 def write_inference_case(
10     case_root: Path,
11     image_values: Iterable[int],
12     case_name: str = "gesture_upload",
13     expected_class: int = -1,
14 ) -> Path:
15     case_root.mkdir(parents=True, exist_ok=True)
16
17     image_path = case_root / "tb_conv1_in_i8_64x64x1.txt"
18     manifest_path = case_root / "manifest.txt"
19
20     write_txt_int8_image(image_path, image_values)
21
```

```

22     manifest_lines = [
23         f"case={case_name}",
24         "image_file=uploaded_image",
25         f"conv1_input={image_path.name}",
26         f"predict_class={expected_class}",
27     ]
28     manifest_path.write_text("\n".join(manifest_lines) + "\n", encoding="utf-8")
29     return case_root

```

### A.3.8 Image Preprocessing

*Converts uploaded image bytes into the 64-by-64 INT8 input tensor and preview image used by the demo.*

code/gesture\_runtime/preprocess.py

```

1  """Input preprocessing helpers.
2
3  The current RTL expects a quantized 64x64x1 INT8 tensor. For hardware-aligned
4  tests we support direct TXT loading; for image files we expose Pillow-based
5  helpers that work both from on-disk paths and in-memory upload bytes.
6  """
7
8  import base64
9  from dataclasses import dataclass
10 from io import BytesIO
11 from pathlib import Path
12 from typing import Iterable, List
13
14 from .quantization import quantize_u8_to_i8
15
16
17 @dataclass
18 class PreprocessedImage:
19     mode: str
20     values: List[int]
21     preview_b64: str
22
23
24 def load_txt_int8_image(path: Path) -> List[int]:
25     with path.open("r", encoding="utf-8") as fp:
26         return [int(line.strip()) for line in fp if line.strip()]
27
28
29 def _encode_image_b64(image) -> str:
30     buf = BytesIO()
31     image.save(buf, format="PNG")
32     return base64.b64encode(buf.getvalue()).decode("ascii")
33
34

```

```

35 def _otsu_threshold(arr) -> int:
36     import numpy as np
37
38     hist = np.bincount(arr.reshape(-1), minlength=256).astype(np.float64)
39     total = arr.size
40     cumulative = np.cumsum(hist)
41     cumulative_mean = np.cumsum(hist * np.arange(256))
42     global_mean = cumulative_mean[-1] / max(total, 1)
43
44     numerator = (global_mean * cumulative - cumulative_mean) ** 2
45     denominator = cumulative * (total - cumulative)
46     with np.errstate(divide="ignore", invalid="ignore"):
47         score = numerator / denominator
48     score[~np.isfinite(score)] = 0.0
49     return int(score.argmax())
50
51
52 def _select_foreground_mask(arr):
53     import numpy as np
54
55     threshold = _otsu_threshold(arr)
56     dark_mask = arr <= threshold
57     light_mask = arr >= threshold
58     border = np.concatenate([arr[0, :], arr[-1, :], arr[:, 0], arr[:, -1]])
59     bg_hint = int(np.median(border))
60
61     def score(mask):
62         area = float(mask.mean())
63         border_frac = float(
64             np.concatenate([mask[0, :], mask[-1, :], mask[:, 0], mask[:, -1]]).mean()
65         )
66         area_penalty = 0.0
67         if area < 0.01:
68             area_penalty += 1.0
69         if area > 0.85:
70             area_penalty += 1.0
71         return border_frac + area_penalty
72
73     if bg_hint >= threshold:
74         preferred = dark_mask
75         alternate = light_mask
76     else:
77         preferred = light_mask
78         alternate = dark_mask
79
80     return preferred if score(preferred) <= score(alternate) else alternate
81
82
83 def _estimate_background_level(arr) -> int:
84     import numpy as np
85
86     border = np.concatenate([arr[0, :], arr[-1, :], arr[:, 0], arr[:, -1]])
87     return int(np.median(border))
88

```

```

89
90 def _canonical_background_level(bg_level: int) -> int:
91     # Snap the background to a clean binary pole so upload previews and model
92     # input do not retain curtain / wall gradients after foreground extraction.
93     return 255 if bg_level >= 128 else 0
94
95
96 def _largest_connected_component(mask):
97     import numpy as np
98
99     height, width = mask.shape
100    visited = np.zeros((height, width), dtype=bool)
101    best_component = []
102
103    for y in range(height):
104        for x in range(width):
105            if not mask[y, x] or visited[y, x]:
106                continue
107
108            stack = [(y, x)]
109            visited[y, x] = True
110            component = []
111
112            while stack:
113                cy, cx = stack.pop()
114                component.append((cy, cx))
115                for ny in range(max(0, cy - 1), min(height, cy + 2)):
116                    for nx in range(max(0, cx - 1), min(width, cx + 2)):
117                        if visited[ny, nx] or not mask[ny, nx]:
118                            continue
119                        visited[ny, nx] = True
120                        stack.append((ny, nx))
121
122                if len(component) > len(best_component):
123                    best_component = component
124
125    if not best_component:
126        return mask
127
128    refined = np.zeros_like(mask, dtype=bool)
129    for y, x in best_component:
130        refined[y, x] = True
131    return refined
132
133
134 def _refine_foreground_mask(mask):
135     import numpy as np
136     from PIL import Image, ImageFilter
137
138     mask_img = Image.fromarray((mask.astype(np.uint8) * 255), mode="L")
139     # Close small gaps first, then trim isolated speckles.
140     mask_img = mask_img.filter(ImageFilter.MaxFilter(5))
141     mask_img = mask_img.filter(ImageFilter.MinFilter(5))
142     mask_img = mask_img.filter(ImageFilter.MaxFilter(3))

```

```

143     refined = np.array(mask_img, dtype=np.uint8) >= 128
144     return _largest_connected_component(refined)
145
146
147 def _crop_and_square_image(image, width: int, height: int):
148     import numpy as np
149     from PIL import Image, ImageOps
150
151     gray = ImageOps.autocontrast(image.convert("L"))
152     arr = np.array(gray, dtype=np.uint8)
153     bg_level = _canonical_background_level(_estimate_background_level(arr))
154     mask = _refine_foreground_mask(_select_foreground_mask(arr))
155
156     ys, xs = np.nonzero(mask)
157     if len(xs) == 0 or len(ys) == 0:
158         cropped = gray
159     else:
160         x0, x1 = xs.min(), xs.max()
161         y0, y1 = ys.min(), ys.max()
162         box_w = x1 - x0 + 1
163         box_h = y1 - y0 + 1
164         pad_x = max(2, int(box_w * 0.12))
165         pad_y = max(2, int(box_h * 0.12))
166         x0 = max(0, x0 - pad_x)
167         y0 = max(0, y0 - pad_y)
168         x1 = min(arr.shape[1] - 1, x1 + pad_x)
169         y1 = min(arr.shape[0] - 1, y1 + pad_y)
170
171         masked_arr = np.full_like(arr, bg_level)
172         masked_arr[mask] = arr[mask]
173         masked = Image.fromarray(masked_arr, mode="L")
174         cropped = masked.crop((int(x0), int(y0), int(x1 + 1), int(y1 + 1)))
175
176         side = max(cropped.size)
177         square = Image.new("L", (side, side), color=bg_level)
178         offset = ((side - cropped.size[0]) // 2, (side - cropped.size[1]) // 2)
179         square.paste(cropped, offset)
180         return square.resize((width, height))
181
182
183 def _plain_resize_image(image, width: int, height: int):
184     from PIL import ImageOps
185
186     return ImageOps.autocontrast(image.convert("L")).resize((width, height))
187
188
189 def _image_to_int8(image) -> List[int]:
190     return quantize_u8_to_i8(list(image.getdata()), zero_point=-128)
191
192
193 def _preprocess_image(image, width: int, height: int, mode: str) -> PreprocessedImage:
194     if mode == "plain":
195         processed = _plain_resize_image(image, width, height)
196     elif mode == "crop":

```

```

197     processed = _crop_and_square_image(image, width, height)
198 else:
199     raise ValueError(f"unknown preprocess mode: {mode}")
200 return PreprocessedImage(mode=mode, values=_image_to_int8(processed), preview_b64=_encode_image_b64(
    processed))
201
202
203 def preprocess_image_bytes_variants(image_bytes: bytes, width: int = 64, height: int = 64) -> List[
    PreprocessedImage]:
204     try:
205         from PIL import Image
206     except ImportError as exc: # pragma: no cover - optional dependency
207         raise RuntimeError("Pillow is required for image preprocessing") from exc
208
209     with Image.open(BytesIO(image_bytes)) as img:
210         return [
211             _preprocess_image(img, width, height, "plain"),
212             _preprocess_image(img, width, height, "crop"),
213         ]
214
215
216 def render_preprocessed_preview_b64(image_bytes: bytes, width: int = 64, height: int = 64) -> str:
217     return preprocess_image_bytes_variants(image_bytes, width, height)[1].preview_b64
218
219
220 def load_image_bytes_to_int8(image_bytes: bytes, width: int = 64, height: int = 64) -> List[int]:
221     return preprocess_image_bytes_variants(image_bytes, width, height)[1].values
222
223
224 def load_image_to_int8(path: Path, width: int = 64, height: int = 64) -> List[int]:
225     try:
226         from PIL import Image
227     except ImportError as exc: # pragma: no cover - optional dependency
228         raise RuntimeError("Pillow is required for image preprocessing") from exc
229
230     with Image.open(path) as img:
231         return _preprocess_image(img, width, height, "crop").values
232
233
234 def write_txt_int8_image(path: Path, values: Iterable[int]) -> None:
235     with path.open("w", encoding="utf-8") as fp:
236         for value in values:
237             fp.write(f"{int(value)}\n")

```

### A.3.9 Preprocess Mode Selection

*Chooses the preprocessing branch used by auto mode based on CPU-side scoring.*

```

1  """Heuristics for choosing between preprocessing variants."""
2
3  from __future__ import annotations
4
5  from typing import Any, Iterable, Tuple
6
7
8  def choose_best_variant(scored: Iterable[Tuple[Any, Any]]) -> Tuple[Any, Any]:
9      """Choose a preprocess branch using a conservative disagreement policy.
10
11     The background-removed crop branch can occasionally become overconfident on
12     the wrong class. We therefore keep the old "pick by margin" behavior when
13     both branches agree, but when they disagree we only trust 'crop' if it is
14     *obviously* stronger while the 'plain' branch looks weak.
15     """
16
17     scored = list(scored)
18     if not scored:
19         raise ValueError("scored variants cannot be empty")
20
21     if len(scored) == 1:
22         return scored[0]
23
24     ordered = sorted(
25         scored,
26         key=lambda item: (item[1].margin, item[1].top_score),
27         reverse=True,
28     )
29     best_variant, best_result = ordered[0]
30     second_variant, second_result = ordered[1]
31
32     if best_result.predicted_class == second_result.predicted_class:
33         return best_variant, best_result
34
35     plain_pair = next(((v, r) for v, r in scored if getattr(v, "mode", "") == "plain"), None)
36     crop_pair = next(((v, r) for v, r in scored if getattr(v, "mode", "") == "crop"), None)
37     if plain_pair is None or crop_pair is None:
38         return best_variant, best_result
39
40     plain_variant, plain_result = plain_pair
41     crop_variant, crop_result = crop_pair
42
43     # If crop saturates near the int8 ceiling while plain is weak, trust crop.
44     if crop_result.top_score >= 120.0 and plain_result.top_score < 60.0:
45         return crop_variant, crop_result
46
47     # If plain is already reasonably confident, prefer its more literal view.
48     if plain_result.margin >= 60.0:
49         return plain_variant, plain_result
50
51     # Otherwise fall back to the stronger-scoring branch.
52     return best_variant, best_result

```

### A.3.10 INT8 Quantization Helpers

*Contains the small quantization helpers used by preprocessing and export.*

```
code/gesture_runtime/quantization.py

1  """Shared quantization helpers for host/runtime code.
2
3  The current RTL datapath is INT8 activation/weight with INT32 accumulation.
4  The reference Gesture design used a simple Q7.8-style short transport for some
5  paths; we keep those helpers here for interoperability, but the project-wide
6  default remains the RTL-native INT8/INT32 format.
7  """
8
9  from typing import Iterable, List
10
11
12  Q7_8_SCALE = 128.0
13
14
15  def saturate_int8(value: int) -> int:
16      return max(-128, min(127, int(value)))
17
18
19  def saturate_int16(value: int) -> int:
20      return max(-32768, min(32767, int(value)))
21
22
23  def saturate_int32(value: int) -> int:
24      return max(-(2**31), min(2**31 - 1, int(value)))
25
26
27  def float_to_fixed_q7_8(value: float) -> int:
28      return saturate_int16(round(value * Q7_8_SCALE))
29
30
31  def fixed_q7_8_to_float(value: int) -> float:
32      return int(value) / Q7_8_SCALE
33
34
35  def quantize_u8_to_i8(samples: Iterable[int], zero_point: int = -128) -> List[int]:
36      return [saturate_int8(int(v) + zero_point) for v in samples]
37
38
39  def dequantize_q7_8_buffer(samples: Iterable[int]) -> List[float]:
40      return [fixed_q7_8_to_float(v) for v in samples]
```

### A.3.11 Host CPU Classifier

*Runs the TFLite CPU model on the host side for comparison and preprocessing selection.*

```

1  """CPU-side TFLite inference helpers for the web comparison demo."""
2
3  from dataclasses import dataclass
4  from pathlib import Path
5  from time import perf_counter
6  from typing import List, Optional
7
8
9  def _load_interpreter(model_path: Path):
10     try:
11         from tflite_runtime.interpreter import Interpreter # type: ignore
12     except ImportError:
13         try:
14             import tensorflow as tf # type: ignore
15
16             Interpreter = tf.lite.Interpreter
17         except ImportError as exc: # pragma: no cover - optional dependency
18             raise RuntimeError(
19                 "A TFLite interpreter is required. Install tflite-runtime or tensorflow."
20             ) from exc
21
22     return Interpreter(model_path=str(model_path))
23
24
25 @dataclass
26 class CPUInferenceResult:
27     predicted_class: int
28     elapsed_ms: float
29     scores: List[float]
30     top_score: float
31     margin: float
32
33
34 class TFLiteCPUClassifier:
35     def __init__(self, model_path: Path):
36         self.model_path = Path(model_path)
37         self.interpreter = _load_interpreter(self.model_path)
38         self.interpreter.allocate_tensors()
39         self.input_details = self.interpreter.get_input_details()[0]
40         self.output_details = self.interpreter.get_output_details()[0]
41
42         import numpy as np
43
44         self.np = np
45
46     def _prepare_input(self, image_values: List[int]):
47         np = self.np
48         input_shape = self.input_details["shape"]
49         input_dtype = self.input_details["dtype"]
50         quant = self.input_details.get("quantization", (0.0, 0))
51
52         arr = np.array(image_values, dtype=np.int16).reshape((1, 64, 64, 1))

```

```

53
54     if input_dtype == np.int8:
55         return arr.astype(np.int8)
56     if input_dtype == np.uint8:
57         return (arr + 128).astype(np.uint8)
58     if input_dtype in (np.float32, np.float64):
59         return ((arr.astype(np.float32) + 128.0) / 255.0).astype(input_dtype)
60
61     scale, zero_point = quant
62     if scale:
63         return ((arr.astype(np.float32) / scale) + zero_point).astype(input_dtype)
64     return arr.astype(input_dtype)
65
66 def predict_int8_image(self, image_values: List[int]) -> CPUInferenceResult:
67     np = self.np
68     input_tensor = self._prepare_input(image_values)
69
70     started = perf_counter()
71     self.interpreter.set_tensor(self.input_details["index"], input_tensor)
72     self.interpreter.invoke()
73     output = self.interpreter.get_tensor(self.output_details["index"])
74     elapsed_ms = (perf_counter() - started) * 1000.0
75
76     scores = np.array(output).reshape(-1).astype(np.float32).tolist()
77     predicted_class = int(np.argmax(scores))
78     ordered = sorted(scores, reverse=True)
79     top_score = float(ordered[0]) if ordered else 0.0
80     second_score = float(ordered[1]) if len(ordered) > 1 else float("-inf")
81     margin = top_score - second_score if ordered else 0.0
82     return CPUInferenceResult(
83         predicted_class=predicted_class,
84         elapsed_ms=elapsed_ms,
85         scores=scores,
86         top_score=top_score,
87         margin=margin,
88     )

```

### A.3.12 Web Front-End Template

*Defines the browser-side upload, preview, status, result, and feedback UI structure.*

```
code/web_demo/templates/index.html
```

```

1 <!doctype html>
2 <html lang="en">
3   <head>
4     <meta charset="utf-8">
5     <meta name="viewport" content="width=device-width, initial-scale=1">
6     <title>CNN_ACC Gesture Demo</title>
7     <link rel="stylesheet" href="{% url_for('static', filename='styles.css') %}">
8   </head>

```

```

9 <body>
10 <main class="page">
11 <section class="hero">
12 <p class="eyebrow">Board CPU vs FPGA Gesture Comparison</p>
13 <h1>Upload One Gesture Image, Run Two Inference Paths</h1>
14 <p class="lede">
15     The same browser-uploaded image is preprocessed once, then compared across
16     board-side HPS CPU inference and the DE1-SoC FPGA accelerator path.
17 </p>
18 <div id="model-warning" class="model-warning" hidden></div>
19 </section>
20
21 <section class="layout">
22 <div class="panel input-panel">
23 <h2>Input</h2>
24 <p class="panel-copy">Upload an image from your computer or capture one from your browser
25 webcam. Wrong predictions can be saved as correction samples for later retraining.</p>
26
27 <label class="upload-box" for="image-input">
28 <input id="image-input" type="file" accept="image/*">
29 <span>Select Image</span>
30 <small>PNG, JPG, or other browser-supported formats</small>
31 </label>
32
33 <div class="preprocess-card">
34 <label class="preprocess-field">
35 <span>Preprocess Mode</span>
36 <select id="preprocess-mode">
37 <option value="auto" selected>Auto</option>
38 <option value="plain">Plain Resize</option>
39 <option value="crop">Crop + Background Removal</option>
40 </select>
41 </label>
42 <div id="branch-diagnostics" class="branch-diagnostics">Branch diagnostics will appear after
43 one run.</div>
44 </div>
45
46 <div class="input-top">
47 <div class="webcam-card">
48 <div class="webcam-head">
49 <h3>Browser Webcam</h3>
50 <div class="webcam-actions">
51 <button id="start-camera-btn" class="secondary">Start Camera</button>
52 <button id="capture-btn" class="secondary" disabled>Capture Frame</button>
53 </div>
54 </div>
55 <div class="camera-stage">
56 <video id="camera-preview" autoplay playsinline muted></video>
57 <div id="camera-placeholder" class="camera-placeholder">Camera is off</div>
58 </div>
59 <p class="camera-copy">This uses your PC/browser camera only. The captured frame is sent
60 through the same pipeline as upload mode.</p>
61 </div>

```

```

60     <div class="preview-card">
61         <div class="preview-head">
62             <h3>Upload Preview</h3>
63             <span class="preview-note">Image selected from file upload</span>
64         </div>
65         <div class="preview-wrap">
66             <img id="preview" alt="Uploaded preview">
67             <div id="preview-placeholder" class="preview-placeholder">No image selected yet</div>
68         </div>
69     </div>
70
71     <div class="processed-wrap">
72         <span class="processed-label">Model Input (64x64)</span>
73         <img id="processed-preview" alt="Preprocessed preview">
74         <div id="processed-placeholder" class="processed-placeholder">Will appear after inference</
div>
75     <div class="actions">
76         <button id="run-btn" class="primary" disabled>Run Comparison</button>
77         <button id="reset-btn" class="secondary">Clear</button>
78     </div>
79 </div>
80 </div>
81
82 <div class="feedback-card">
83     <div class="feedback-head">
84         <h3>Correction Capture</h3>
85         <span class="feedback-note">Use this when the prediction is wrong.</span>
86     </div>
87     <div class="feedback-controls">
88         <label class="feedback-field">
89             <span>Correct Gesture ID</span>
90             <select id="correct-label" disabled>
91                 <option value="">Select</option>
92                 <option value="0">0</option>
93                 <option value="1">1</option>
94                 <option value="2">2</option>
95                 <option value="3">3</option>
96                 <option value="4">4</option>
97                 <option value="5">5</option>
98                 <option value="6">6</option>
99                 <option value="7">7</option>
100                <option value="8">8</option>
101                <option value="9">9</option>
102            </select>
103        </label>
104        <label class="feedback-field feedback-field-wide">
105            <span>Note</span>
106            <input id="feedback-note" type="text" placeholder="Optional note about lighting,
background, hand angle...">
107        </label>
108    </div>
109    <button id="save-feedback-btn" class="secondary" disabled>Save Correction Sample</button>
110    <div id="feedback-status" class="feedback-status">Run one inference first if you want to save
a corrected sample.</div>

```

```

111     </div>
112
113 </div>
114
115 <div class="side-column">
116   <div class="panel result-panel">
117     <h2>Results</h2>
118     <p class="panel-copy">This panel compares the board-side HPS CPU software baseline against
the existing board-side HPS + MMIO + FPGA path for the active 10-class gesture model line.</p>
119
120     <div class="prediction-strip">
121       <div class="metric">
122         <span class="metric-label">Predicted Gesture ID</span>
123         <strong id="predicted-class">-</strong>
124       </div>
125       <div class="metric">
126         <span class="metric-label">Board CPU Gesture ID</span>
127         <strong id="cpu-predicted-class">-</strong>
128       </div>
129     </div>
130
131     <div class="metrics">
132       <div class="metric">
133         <span class="metric-label">Board CPU Time</span>
134         <strong id="cpu-time">-</strong>
135       </div>
136       <div class="metric">
137         <span class="metric-label">FPGA Wait Time</span>
138         <strong id="fpga-time">-</strong>
139       </div>
140       <div class="metric">
141         <span class="metric-label">FPGA RTL Time</span>
142         <strong id="fpga-rtl-time">-</strong>
143       </div>
144       <div class="metric speedup">
145         <span class="metric-label">Speedup</span>
146         <strong id="speedup">-</strong>
147       </div>
148     </div>
149
150     <div class="timing-note">
151       <span id="timing-note">Board CPU and FPGA wait time are now same-platform metrics. Remote
request time is reported separately.</span>
152     </div>
153
154     <div class="confidence-card">
155       <div class="confidence-head">
156         <h3>CPU Confidence And Channel Scores</h3>
157         <span id="confidence-note" class="confidence-note">Waiting for one run.</span>
158       </div>
159       <div class="confidence-summary">
160         <div class="confidence-metric">
161           <span>Top Class</span>
162           <strong id="cpu-top-class">-</strong>

```

```

163         </div>
164         <div class="confidence-metric">
165             <span>Normalized Confidence</span>
166             <strong id="cpu-confidence"></strong>
167         </div>
168         <div class="confidence-metric">
169             <span>Margin</span>
170             <strong id="cpu-margin"></strong>
171         </div>
172     </div>
173     <div id="cpu-top-channels" class="cpu-top-channels">Top channels will appear after one run.
</div>

174     <div id="channel-grid" class="channel-grid"></div>
175 </div>
176
177 <div class="profile-card">
178     <div class="profile-head">
179         <h3>FPGA RTL Profile</h3>
180         <span id="fpga-profile-meta" class="profile-meta">Waiting for one run.</span>
181     </div>
182     <div id="profile-grid" class="profile-grid">
183         <div class="profile-row">
184             <span>L1</span>
185             <strong id="profile-l1"></strong>
186         </div>
187         <div class="profile-row">
188             <span>L2 P0</span>
189             <strong id="profile-l2-p0"></strong>
190         </div>
191         <div class="profile-row">
192             <span>L2 P1</span>
193             <strong id="profile-l2-p1"></strong>
194         </div>
195         <div class="profile-row">
196             <span>L3 P0</span>
197             <strong id="profile-l3-p0"></strong>
198         </div>
199         <div class="profile-row">
200             <span>L3 P1</span>
201             <strong id="profile-l3-p1"></strong>
202         </div>
203         <div class="profile-row">
204             <span>FC</span>
205             <strong id="profile-fc"></strong>
206         </div>
207         <div class="profile-row">
208             <span>Argmax</span>
209             <strong id="profile-argmax"></strong>
210         </div>
211         <div class="profile-row profile-row-total">
212             <span>Total</span>
213             <strong id="profile-total"></strong>
214         </div>
215     </div>

```

```

216     </div>
217
218     <div class="chart-card">
219         <h3>Execution Time Comparison</h3>
220         <div class="chart">
221             <div class="bar-group">
222                 <span>Board CPU</span>
223                 <div class="bar-track"><div id="cpu-bar" class="bar"></div></div>
224                 <span id="cpu-bar-label">-</span>
225             </div>
226             <div class="bar-group">
227                 <span>FPGA</span>
228                 <div class="bar-track"><div id="fpga-bar" class="bar bar-alt"></div></div>
229                 <span id="fpga-bar-label">-</span>
230             </div>
231         </div>
232     </div>
233
234     <div id="status-box" class="status-box">Waiting for an upload.</div>
235 </div>
236
237 <div class="panel suite-panel">
238     <div class="suite-head">
239         <h3>Built-In Sanity Suite</h3>
240         <button id="suite-btn" class="secondary suite-btn">Run 0-9 Suite</button>
241     </div>
242     <p class="suite-copy">Runs the current web path on the repo's built-in 'digit_0_test' to '
digit_9_test' sample images so we can verify whether the deployed model line is internally
consistent.</p>
243     <div id="suite-summary" class="suite-summary">Not run yet.</div>
244     <div id="suite-table-wrap" class="suite-table-wrap" hidden>
245         <table class="suite-table">
246             <thead>
247                 <tr>
248                     <th>Case</th>
249                     <th>Exp</th>
250                     <th>CPU</th>
251                     <th>FPGA</th>
252                     <th>Result</th>
253                 </tr>
254             </thead>
255             <tbody id="suite-table-body"></tbody>
256         </table>
257     </div>
258 </div>
259 </div>
260 </section>
261 </main>
262 <script src="{ url_for('static', filename='app.js') }"></script>
263 </body>
264 </html>

```

### A.3.13 Web Front-End Script

*Implements browser-side upload handling, inference requests, progress state, and result rendering.*

```
code/web_demo/static/app.js
```

```
1  const imageInput = document.getElementById("image-input");
2  const preview = document.getElementById("preview");
3  const previewPlaceholder = document.getElementById("preview-placeholder");
4  const processedPreview = document.getElementById("processed-preview");
5  const processedPlaceholder = document.getElementById("processed-placeholder");
6  const runBtn = document.getElementById("run-btn");
7  const resetBtn = document.getElementById("reset-btn");
8  const startCameraBtn = document.getElementById("start-camera-btn");
9  const captureBtn = document.getElementById("capture-btn");
10 const cameraPreview = document.getElementById("camera-preview");
11 const cameraPlaceholder = document.getElementById("camera-placeholder");
12 const statusBox = document.getElementById("status-box");
13 const modelWarningBox = document.getElementById("model-warning");
14 const suiteBtn = document.getElementById("suite-btn");
15 const suiteSummary = document.getElementById("suite-summary");
16 const suiteTableWrap = document.getElementById("suite-table-wrap");
17 const suiteTableBody = document.getElementById("suite-table-body");
18 const correctLabel = document.getElementById("correct-label");
19 const feedbackNote = document.getElementById("feedback-note");
20 const saveFeedbackBtn = document.getElementById("save-feedback-btn");
21 const feedbackStatus = document.getElementById("feedback-status");
22 const preprocessModeSelect = document.getElementById("preprocess-mode");
23 const branchDiagnosticsBox = document.getElementById("branch-diagnostics");
24
25 let selectedFile = null;
26 let cameraStream = null;
27 let lastResult = null;
28 let currentPreviewUrl = null;
29
30 function setStatus(message, isError = false) {
31   statusBox.textContent = message;
32   statusBox.classList.toggle("error", isError);
33 }
34
35 function setModelWarning(message) {
36   if (!message) {
37     modelWarningBox.hidden = true;
38     modelWarningBox.textContent = "";
39     return;
40   }
41   modelWarningBox.hidden = false;
42   modelWarningBox.textContent = message;
43 }
44
45 function resetResults() {
46   lastResult = null;
47   document.getElementById("predicted-class").textContent = "-";
48   document.getElementById("cpu-predicted-class").textContent = "-";
```

```

49 document.getElementById("cpu-time").textContent = "-";
50 document.getElementById("fpga-time").textContent = "-";
51 document.getElementById("fpga-rtl-time").textContent = "-";
52 document.getElementById("speedup").textContent = "-";
53 document.getElementById("timing-note").textContent =
54     "FPGA RTL time is derived from on-chip cycle counters. Board wait time is the HPS-observed elapsed
        time.";
55 document.getElementById("fpga-profile-meta").textContent = "Waiting for one run.";
56 document.getElementById("profile-l1").textContent = "-";
57 document.getElementById("profile-l2-p0").textContent = "-";
58 document.getElementById("profile-l2-p1").textContent = "-";
59 document.getElementById("profile-l3-p0").textContent = "-";
60 document.getElementById("profile-l3-p1").textContent = "-";
61 document.getElementById("profile-fc").textContent = "-";
62 document.getElementById("profile-argmax").textContent = "-";
63 document.getElementById("profile-total").textContent = "-";
64 document.getElementById("cpu-bar").style.width = "0%";
65 document.getElementById("fpga-bar").style.width = "0%";
66 document.getElementById("cpu-bar-label").textContent = "-";
67 document.getElementById("fpga-bar-label").textContent = "-";
68 document.getElementById("cpu-top-class").textContent = "-";
69 document.getElementById("cpu-confidence").textContent = "-";
70 document.getElementById("cpu-margin").textContent = "-";
71 document.getElementById("cpu-top-channels").textContent = "Top channels will appear after one run.";
72 document.getElementById("confidence-note").textContent = "Waiting for one run.";
73 document.getElementById("channel-grid").innerHTML = "";
74 branchDiagnosticsBox.textContent = "Branch diagnostics will appear after one run.";
75 processedPreview.removeAttribute("src");
76 processedPreview.style.display = "none";
77 processedPlaceholder.style.display = "grid";
78 correctLabel.value = "";
79 correctLabel.disabled = true;
80 saveFeedbackBtn.disabled = true;
81 feedbackStatus.textContent = "Run one inference first if you want to save a corrected sample.";
82 }
83
84 function renderBranchDiagnostics(rows, chosenMode) {
85     if (!rows || rows.length === 0) {
86         branchDiagnosticsBox.textContent = "Branch diagnostics unavailable.";
87         return;
88     }
89     const parts = rows.map((row) => {
90         const chosenTag = row.mode === chosenMode ? " [chosen]" : "";
91         const topSummary = (row.top_channels || [])
92             .slice(0, 2)
93             .map((entry) => `${entry.label}:${entry.raw_score.toFixed(1)}'`
94             .join(", ");
95         return `${row.mode}: pred ${row.predicted_label}, margin ${row.margin.toFixed(1)}, top ${row.
96             top_score.toFixed(1)}${topSummary ? ', leaders ${topSummary}' : ""}${chosenTag}`;
97     });
98     branchDiagnosticsBox.textContent = parts.join(" | ");
99 }
100 function renderCpuScoreReport(report) {

```

```

101   if (!report) {
102     document.getElementById("cpu-top-class").textContent = "-";
103     document.getElementById("cpu-confidence").textContent = "-";
104     document.getElementById("cpu-margin").textContent = "-";
105     document.getElementById("cpu-top-channels").textContent = "Top channels will appear after one run.";
106     document.getElementById("confidence-note").textContent = "Waiting for one run.";
107     document.getElementById("channel-grid").innerHTML = "";
108     return;
109   }
110
111   document.getElementById("cpu-top-class").textContent = report.predicted_label;
112   document.getElementById("cpu-confidence").textContent = `${report.confidence_pct.toFixed(1)}%`;
113   document.getElementById("cpu-margin").textContent = report.margin.toFixed(1);
114   document.getElementById("confidence-note").textContent = report.note || "";
115
116   const topSummary = (report.top_channels || [])
117     .map((entry) => `${entry.label}: raw ${entry.raw_score.toFixed(1)}, conf ${entry.probability_pct.
118       toFixed(1)}%`)
119     .join(" | ");
120   document.getElementById("cpu-top-channels").textContent =
121     topSummary || "Top channels unavailable.";
122
123   const grid = document.getElementById("channel-grid");
124   grid.innerHTML = "";
125   for (const channel of report.channels || []) {
126     const row = document.createElement("div");
127     row.className = "channel-row";
128     row.innerHTML = `
129       <span class="channel-label">${channel.label}</span>
130       <div class="channel-track"><div class="channel-fill" style="width:${channel.probability_pct.toFixed
131         (2)}%></div></div>
132       <span class="channel-score">${channel.raw_score.toFixed(1)}</span>
133       <span class="channel-prob">${channel.probability_pct.toFixed(1)}%</span>
134     `;
135     grid.appendChild(row);
136   }
137
138   function formatProfileRow(stage) {
139     if (!stage) {
140       return "-";
141     }
142     const timeUs =
143       stage.time_us < 0.1 ? stage.time_us.toFixed(2) :
144       stage.time_us < 1 ? stage.time_us.toFixed(2) :
145       stage.time_us.toFixed(1);
146     return `${stage.cycles} cyc / ${timeUs} us`;
147   }
148
149   function renderProfile(profile) {
150     if (!profile) {
151       return;
152     }
153     const stageByLabel = Object.fromEntries(profile.stages.map((stage) => [stage.label, stage]));

```

```

153 document.getElementById("fpga-profile-meta").textContent =
154     `${profile.fabric_mhz.toFixed(2)} MHz fabric clock`;
155 document.getElementById("profile-l1").textContent = formatProfileRow(stageByLabel["L1"]);
156 document.getElementById("profile-l2-p0").textContent = formatProfileRow(stageByLabel["L2 P0"]);
157 document.getElementById("profile-l2-p1").textContent = formatProfileRow(stageByLabel["L2 P1"]);
158 document.getElementById("profile-l3-p0").textContent = formatProfileRow(stageByLabel["L3 P0"]);
159 document.getElementById("profile-l3-p1").textContent = formatProfileRow(stageByLabel["L3 P1"]);
160 document.getElementById("profile-fc").textContent = formatProfileRow(stageByLabel["FC"]);
161 document.getElementById("profile-argmax").textContent = formatProfileRow(stageByLabel["Argmax"]);
162 document.getElementById("profile-total").textContent =
163     `${profile.total_cycles} cyc / ${profile.rtl_time_us.toFixed(1)} us`;
164 }
165
166 function setPreviewSource(url) {
167     if (currentPreviewUrl && currentPreviewUrl.startsWith("blob:")) {
168         URL.revokeObjectURL(currentPreviewUrl);
169     }
170     currentPreviewUrl = url;
171     preview.src = url;
172     preview.style.display = "block";
173     previewPlaceholder.style.display = "none";
174 }
175
176 async function startCamera() {
177     if (cameraStream) {
178         return;
179     }
180     const stream = await navigator.mediaDevices.getUserMedia({
181         video: {
182             facingMode: "user",
183             width: { ideal: 960 },
184             height: { ideal: 960 },
185         },
186         audio: false,
187     });
188     cameraStream = stream;
189     cameraPreview.srcObject = stream;
190     cameraPreview.style.display = "block";
191     cameraPlaceholder.style.display = "none";
192     captureBtn.disabled = false;
193     startCameraBtn.textContent = "Camera Ready";
194     startCameraBtn.disabled = true;
195 }
196
197 async function captureFrame() {
198     if (!cameraStream) {
199         return;
200     }
201     const width = cameraPreview.videoWidth || 640;
202     const height = cameraPreview.videoHeight || 640;
203     const canvas = document.createElement("canvas");
204     canvas.width = width;
205     canvas.height = height;
206     const context = canvas.getContext("2d");

```

```

207 context.drawImage(cameraPreview, 0, 0, width, height);
208 const blob = await new Promise((resolve) => canvas.toBlob(resolve, "image/png"));
209 if (!blob) {
210     throw new Error("Failed to capture webcam frame");
211 }
212 selectedFile = new File([blob], "webcam_capture.png", { type: "image/png" });
213 imageInput.value = "";
214 resetResults();
215 setPreviewSource(URL.createObjectURL(selectedFile));
216 runBtn.disabled = false;
217 setStatus("Webcam frame captured. Run comparison when you are ready.");
218 }
219
220 function renderSuite(payload) {
221     suiteSummary.textContent = `Built-in suite: ${payload.passed}/${payload.total} passed.`;
222     suiteTableBody.innerHTML = "";
223     for (const row of payload.rows) {
224         const tr = document.createElement("tr");
225         tr.innerHTML = `
226             <td>${row.case}</td>
227             <td>${row.expected_class}</td>
228             <td>${row.cpu_predicted_class}</td>
229             <td>${row.fpga_predicted_class}</td>
230             <td class="${row.pass ? "suite-pass" : "suite-fail"}">${row.pass ? "PASS" : "FAIL"}</td>
231         `;
232         suiteTableBody.appendChild(tr);
233     }
234     suiteTableWrap.hidden = false;
235 }
236
237 function renderChart(cpuMs, fpgaMs) {
238     const maxValue = Math.max(cpuMs, fpgaMs, 1);
239     const cpuWidth = (cpuMs / maxValue) * 100;
240     const fpgaWidth = (fpgaMs / maxValue) * 100;
241     document.getElementById("cpu-bar").style.width = `${cpuWidth}%`;
242     document.getElementById("fpga-bar").style.width = `${fpgaWidth}%`;
243     document.getElementById("cpu-bar-label").textContent = `${cpuMs.toFixed(2)} ms`;
244     document.getElementById("fpga-bar-label").textContent = `${fpgaMs.toFixed(2)} ms`;
245 }
246
247 async function loadHealth() {
248     try {
249         const response = await fetch("/api/health");
250         const payload = await response.json();
251         if (response.ok) {
252             setModelWarning(payload.model_warning || "");
253         }
254     } catch {
255         // Keep the page usable even if health probing fails.
256     }
257 }
258
259 imageInput.addEventListener("change", () => {
260     const [file] = imageInput.files;

```

```

261     selectedFile = file || null;
262     resetResults();
263
264     if (!selectedFile) {
265         preview.removeAttribute("src");
266         preview.style.display = "none";
267         previewPlaceholder.style.display = "grid";
268         runBtn.disabled = true;
269         setStatus("Waiting for an upload.");
270         return;
271     }
272
273     setPreviewSource(URL.createObjectURL(selectedFile));
274     runBtn.disabled = false;
275     setStatus("Image ready. Run comparison when you are ready.");
276 });
277
278 resetBtn.addEventListener("click", () => {
279     imageInput.value = "";
280     selectedFile = null;
281     preview.removeAttribute("src");
282     preview.style.display = "none";
283     previewPlaceholder.style.display = "grid";
284     runBtn.disabled = true;
285     resetResults();
286     setStatus("Waiting for an upload.");
287 });
288
289 startCameraBtn.addEventListener("click", async () => {
290     try {
291         await startCamera();
292         setStatus("Camera ready. Capture a frame when you are ready.");
293     } catch (error) {
294         setStatus(error.message || "Unable to start webcam.", true);
295     }
296 });
297
298 captureBtn.addEventListener("click", async () => {
299     try {
300         await captureFrame();
301     } catch (error) {
302         setStatus(error.message || "Unable to capture webcam frame.", true);
303     }
304 });
305
306 runBtn.addEventListener("click", async () => {
307     if (!selectedFile) {
308         return;
309     }
310
311     runBtn.disabled = true;
312     setStatus("Running CPU and FPGA inference...");
313
314     const formData = new FormData();

```

```

315   formData.append("image", selectedFile);
316   formData.append("preprocess_mode", preprocessModeSelect.value);
317
318   try {
319     const response = await fetch("/api/infer", {
320       method: "POST",
321       body: formData,
322     });
323     const payload = await response.json();
324     if (!response.ok) {
325       throw new Error(payload.error || "Inference failed");
326     }
327
328     lastResult = payload;
329     document.getElementById("predicted-class").textContent = payload.predicted_label;
330     document.getElementById("cpu-predicted-class").textContent = payload.cpu_predicted_label;
331     document.getElementById("cpu-time").textContent = `${payload.cpu_time_ms.toFixed(2)} ms`;
332     document.getElementById("fpga-time").textContent = `${payload.fpga_time_ms.toFixed(2)} ms`;
333     document.getElementById("fpga-rtl-time").textContent = `${payload.fpga_rtl_time_ms.toFixed(3)} ms`;
334     document.getElementById("speedup").textContent =
335       payload.speedup ? `${payload.speedup.toFixed(2)}x` : "-";
336     document.getElementById("timing-note").textContent =
337       `RTL ${payload.fpga_rtl_time_us.toFixed(1)} us from on-chip counters. ` +
338       `Board wait ${payload.fpga_time_ms.toFixed(2)} ms. Remote request ${payload.fpga_request_time_ms.toFixed(2)} ms. ` +
339       `Board total ${payload.fpga_board_total_time_ms.toFixed(2)} ms, load ${payload.fpga_board_case_load_time_ms.toFixed(2)} ms, program ${payload.fpga_board_program_time_ms.toFixed(2)} ms.`;
340
341     renderChart(payload.cpu_time_ms, payload.fpga_time_ms);
342     renderProfile(payload.fpga_profile);
343     renderBranchDiagnostics(payload.branch_diagnostics, payload.preprocess_mode);
344     renderCpuScoreReport(payload.cpu_score_report);
345     if (payload.preprocessed_preview_url) {
346       processedPreview.src = payload.preprocessed_preview_url;
347       processedPreview.style.display = "block";
348       processedPlaceholder.style.display = "none";
349     }
350     correctLabel.disabled = false;
351     saveFeedbackBtn.disabled = false;
352     feedbackStatus.textContent = "If this result is wrong, select the correct label and save the sample for retraining.";
353     setModelWarning(payload.model_warning || "");
354     const preprocessMode = payload.preprocess_mode ? `Preprocess: ${payload.preprocess_mode}.` : "";
355     setStatus(`Done. Board CPU predicted gesture ID ${payload.cpu_predicted_label}, FPGA predicted gesture ID ${payload.predicted_label}.${preprocessMode}`);
356   } catch (error) {
357     setStatus(error.message, true);
358   } finally {
359     runBtn.disabled = false;
360   }
361 });
362
363 saveFeedbackBtn.addEventListener("click", async () => {

```

```

364   if (!selectedFile || !lastResult) {
365     feedbackStatus.textContent = "Run one inference before saving feedback.";
366     return;
367   }
368   if (!correctLabel.value) {
369     feedbackStatus.textContent = "Choose the correct gesture ID first.";
370     return;
371   }
372
373   saveFeedbackBtn.disabled = true;
374   feedbackStatus.textContent = "Saving correction sample...";
375   const formData = new FormData();
376   formData.append("image", selectedFile);
377   formData.append("corrected_label", correctLabel.value);
378   formData.append("predicted_label", String(lastResult.predicted_label));
379   formData.append("preprocess_mode", String(lastResult.preprocess_mode || ""));
380   formData.append("note", feedbackNote.value || "");
381
382   try {
383     const response = await fetch("/api/feedback", {
384       method: "POST",
385       body: formData,
386     });
387     const payload = await response.json();
388     if (!response.ok) {
389       throw new Error(payload.error || "Unable to save correction sample");
390     }
391     feedbackStatus.textContent = `Saved correction sample ${payload.sample_id} under ${payload.saved_dir}.`;
392   } catch (error) {
393     feedbackStatus.textContent = error.message;
394   } finally {
395     saveFeedbackBtn.disabled = false;
396   }
397 });
398
399 suiteBtn.addEventListener("click", async () => {
400   suiteBtn.disabled = true;
401   suiteSummary.textContent = "Running built-in suite...";
402   try {
403     const response = await fetch("/api/sample-suite");
404     const payload = await response.json();
405     if (!response.ok) {
406       throw new Error(payload.error || "Suite failed");
407     }
408     renderSuite(payload);
409   } catch (error) {
410     suiteSummary.textContent = error.message;
411     suiteTableWrap.hidden = true;
412   } finally {
413     suiteBtn.disabled = false;
414   }
415 });
416

```

```
417 loadHealth();
```

### A.3.14 Web Front-End Styles

*Defines the responsive browser UI layout, controls, result panels, and visual states.*

```
code/web_demo/static/styles.css
```

```
1 :root {
2   --bg: #f5efe4;
3   --panel: #fffaf2;
4   --ink: #1b1711;
5   --muted: #655b4c;
6   --accent: #f5b000;
7   --accent-deep: #d89200;
8   --line: #e3d7c2;
9   --success: #173d2a;
10  --error: #8f2d23;
11 }
12
13 * {
14   box-sizing: border-box;
15 }
16
17 body {
18   margin: 0;
19   font-family: "Segoe UI", "Helvetica Neue", Arial, sans-serif;
20   background:
21     radial-gradient(circle at top right, rgba(245, 176, 0, 0.22), transparent 28%),
22     linear-gradient(180deg, #fbf8f1, var(--bg));
23   color: var(--ink);
24 }
25
26 .page {
27   max-width: 1400px;
28   margin: 0 auto;
29   padding: 32px 24px 48px;
30 }
31
32 .hero {
33   margin-bottom: 22px;
34 }
35
36 .model-warning {
37   margin-top: 16px;
38   max-width: 860px;
39   padding: 14px 16px;
40   border: 1px solid rgba(245, 176, 0, 0.35);
41   border-radius: 14px;
42   background: rgba(245, 176, 0, 0.12);
43   color: #6b4c00;
```

```

44   line-height: 1.45;
45 }
46
47 .eyebrow {
48   margin: 0 0 10px;
49   text-transform: uppercase;
50   letter-spacing: 0.12em;
51   font-size: 0.8rem;
52   color: var(--accent-deep);
53   font-weight: 700;
54 }
55
56 .hero h1 {
57   margin: 0 0 12px;
58   font-size: clamp(2rem, 4vw, 3.6rem);
59   line-height: 1.05;
60 }
61
62 .lede {
63   margin: 0;
64   max-width: 720px;
65   color: var(--muted);
66   font-size: 1.02rem;
67 }
68
69 .layout {
70   display: grid;
71   grid-template-columns: minmax(0, 1.08fr) minmax(380px, 0.92fr);
72   gap: 20px;
73   align-items: stretch;
74 }
75
76 .panel {
77   background: var(--panel);
78   border: 1px solid var(--line);
79   border-radius: 24px;
80   padding: 20px;
81   box-shadow: 0 18px 40px rgba(34, 25, 7, 0.08);
82 }
83
84 .input-panel {
85   display: grid;
86   grid-template-columns: 1fr;
87   gap: 16px;
88   align-content: start;
89   min-width: 0;
90   height: 100%;
91 }
92
93 .input-panel > h2,
94 .input-panel > .panel-copy,
95 .input-panel > .upload-box,
96 .input-panel > .preprocess-card,
97 .input-panel > .input-top {

```

```

98   grid-column: auto;
99 }
100
101 .result-panel {
102   display: grid;
103   gap: 16px;
104   align-content: start;
105   min-width: 0;
106 }
107
108 .side-column {
109   display: flex;
110   flex-direction: column;
111   gap: 16px;
112   min-width: 0;
113   height: 100%;
114 }
115
116 .suite-panel {
117   min-width: 0;
118   display: grid;
119   gap: 10px;
120   border-color: rgba(245, 176, 0, 0.32);
121   background:
122     linear-gradient(180deg, rgba(245, 176, 0, 0.08), rgba(255, 255, 255, 0.74));
123   margin-top: auto;
124 }
125
126 .feedback-card {
127   grid-column: auto;
128 }
129
130 .input-top {
131   display: grid;
132   grid-template-columns: repeat(2, minmax(0, 1fr));
133   grid-template-areas:
134     "webcam preview"
135     "processed processed";
136   gap: 16px;
137   align-items: start;
138   min-width: 0;
139 }
140
141 .input-top > * {
142   min-width: 0;
143 }
144
145 .panel h2,
146 .chart-card h3 {
147   margin-top: 0;
148   margin-bottom: 0;
149 }
150
151 .panel-copy {

```

```

152   color: var(--muted);
153   margin: 0;
154   line-height: 1.45;
155 }
156
157 .upload-box {
158   display: grid;
159   gap: 6px;
160   padding: 16px 18px;
161   border: 2px dashed var(--line);
162   border-radius: 18px;
163   background: rgba(245, 176, 0, 0.05);
164   cursor: pointer;
165 }
166
167 .upload-box input {
168   display: none;
169 }
170
171 .upload-box span {
172   font-weight: 700;
173 }
174
175 .upload-box small {
176   color: var(--muted);
177 }
178
179 .preprocess-card {
180   display: grid;
181   grid-template-columns: minmax(220px, 280px) minmax(0, 1fr);
182   gap: 12px;
183   padding: 14px 16px;
184   border: 1px solid var(--line);
185   border-radius: 18px;
186   background: rgba(255, 255, 255, 0.55);
187   align-items: end;
188 }
189
190 .preprocess-field {
191   display: grid;
192   gap: 8px;
193 }
194
195 .preprocess-field span {
196   color: var(--muted);
197   font-size: 0.9rem;
198 }
199
200 .preprocess-field select {
201   width: 100%;
202   border: 1px solid var(--line);
203   border-radius: 12px;
204   padding: 10px 12px;
205   background: #fff;

```

```

206   color: var(--ink);
207 }
208
209 .branch-diagnostics {
210   color: var(--muted);
211   font-size: 0.9rem;
212   line-height: 1.45;
213   min-width: 0;
214 }
215
216 .webcam-card,
217 .feedback-card,
218 .preview-card,
219 .preview-wrap,
220 .processed-wrap {
221   border: 1px solid var(--line);
222   border-radius: 18px;
223   padding: 16px;
224   background: rgba(255, 255, 255, 0.55);
225 }
226
227 .webcam-card {
228   display: flex;
229   flex-direction: column;
230   gap: 10px;
231   grid-area: webcam;
232 }
233
234 .preview-card {
235   display: grid;
236   gap: 10px;
237   grid-area: preview;
238 }
239
240 .webcam-head,
241 .feedback-head,
242 .preview-head {
243   display: flex;
244   align-items: flex-start;
245   justify-content: space-between;
246   gap: 10px;
247 }
248
249 .webcam-head h3,
250 .feedback-head h3,
251 .preview-head h3 {
252   margin: 0;
253 }
254
255 .preview-note {
256   color: var(--muted);
257   font-size: 0.84rem;
258   text-align: right;
259 }

```

```

260
261 .webcam-actions {
262     display: flex;
263     gap: 8px;
264     flex-wrap: wrap;
265 }
266
267 .camera-stage {
268     width: 100%;
269     aspect-ratio: 1 / 1;
270     max-width: 100%;
271     min-height: 280px;
272     border-radius: 14px;
273     overflow: hidden;
274     border: 1px solid var(--line);
275     background: #121212;
276     display: grid;
277     place-items: center;
278 }
279
280 #camera-preview {
281     width: 100%;
282     height: 100%;
283     object-fit: cover;
284     display: none;
285 }
286
287 .camera-placeholder {
288     color: #d6d6d6;
289     padding: 16px;
290     text-align: center;
291 }
292
293 .camera-copy,
294 .feedback-note {
295     color: var(--muted);
296     font-size: 0.9rem;
297     line-height: 1.45;
298 }
299
300 .preview-wrap {
301     overflow: hidden;
302     display: grid;
303     place-items: center;
304     min-width: 0;
305     width: 100%;
306     aspect-ratio: 1 / 1;
307     min-height: 280px;
308     background: #0c0c0c;
309     padding: 0;
310 }
311
312 #preview {
313     width: 100%;

```

```

314 height: 100%;
315 object-fit: cover;
316 display: none;
317 }
318
319 .preview-placeholder {
320 color: #d6d6d6;
321 padding: 16px;
322 text-align: center;
323 }
324
325 .processed-wrap {
326 grid-area: processed;
327 padding: 12px;
328 display: grid;
329 grid-template-columns: minmax(140px, 180px) minmax(0, 1fr);
330 grid-template-areas:
331 "label actions"
332 "preview actions";
333 align-items: start;
334 min-width: 0;
335 gap: 12px;
336 overflow: hidden;
337 }
338
339 .processed-label {
340 display: block;
341 grid-area: label;
342 margin-bottom: 0;
343 color: var(--muted);
344 font-size: 0.92rem;
345 }
346
347 #processed-preview,
348 .processed-placeholder {
349 grid-area: preview;
350 width: min(100%, 160px);
351 aspect-ratio: 1 / 1;
352 height: auto;
353 }
354
355 #processed-preview {
356 image-rendering: pixelated;
357 border-radius: 10px;
358 border: 1px solid var(--line);
359 display: none;
360 background: #111;
361 }
362
363 .processed-placeholder {
364 display: grid;
365 place-items: center;
366 border-radius: 10px;
367 border: 1px dashed var(--line);

```

```

368   color: var(--muted);
369   text-align: center;
370   font-size: 0.9rem;
371   padding: 10px;
372 }
373
374 .actions {
375   grid-area: actions;
376   display: flex;
377   flex-direction: column;
378   gap: 10px;
379   width: 100%;
380   align-self: stretch;
381   justify-self: end;
382   max-width: 320px;
383 }
384
385 .actions button {
386   width: 100%;
387 }
388
389 .feedback-controls {
390   display: grid;
391   grid-template-columns: 150px 1fr;
392   gap: 10px;
393   margin-top: 10px;
394 }
395
396 .feedback-field {
397   display: grid;
398   gap: 8px;
399 }
400
401 .feedback-field-wide {
402   min-width: 0;
403 }
404
405 .feedback-field span {
406   color: var(--muted);
407   font-size: 0.92rem;
408 }
409
410 .feedback-field select,
411 .feedback-field input {
412   width: 100%;
413   border: 1px solid var(--line);
414   border-radius: 12px;
415   padding: 10px 12px;
416   background: #fff;
417   color: var(--ink);
418 }
419
420 .feedback-status {
421   margin-top: 10px;

```

```
422     color: var(--muted);
423     font-size: 0.92rem;
424 }
425
426 .suite-head {
427     display: grid;
428     grid-template-columns: minmax(0, 1fr) auto;
429     align-items: start;
430     gap: 12px;
431 }
432
433 .suite-head h3 {
434     margin: 0;
435 }
436
437 .suite-copy {
438     color: var(--muted);
439     margin: 0;
440     line-height: 1.45;
441 }
442
443 .suite-summary {
444     font-weight: 700;
445 }
446
447 .suite-table-wrap {
448     margin-top: 8px;
449     overflow-x: auto;
450 }
451
452 .suite-table {
453     width: 100%;
454     border-collapse: collapse;
455     font-size: 0.92rem;
456 }
457
458 .suite-table th,
459 .suite-table td {
460     text-align: left;
461     padding: 8px 10px;
462     border-top: 1px solid var(--line);
463 }
464
465 .suite-pass {
466     color: #21553a;
467     font-weight: 700;
468 }
469
470 .suite-fail {
471     color: var(--error);
472     font-weight: 700;
473 }
474
475 button {
```

```

476   border: 0;
477   border-radius: 12px;
478   padding: 11px 16px;
479   font-weight: 700;
480   cursor: pointer;
481 }
482
483 button.primary {
484   background: var(--accent);
485   color: #201500;
486 }
487
488 button.secondary {
489   background: #e9e2d5;
490   color: #3f372b;
491 }
492
493 button.disabled {
494   opacity: 0.55;
495   cursor: not-allowed;
496 }
497
498 .prediction-strip {
499   display: grid;
500   grid-template-columns: repeat(2, minmax(0, 1fr));
501   gap: 12px;
502 }
503
504 .metrics {
505   display: grid;
506   grid-template-columns: repeat(2, minmax(0, 1fr));
507   gap: 12px;
508 }
509
510 .metric {
511   border: 1px solid var(--line);
512   border-radius: 18px;
513   padding: 14px;
514   background: rgba(255, 255, 255, 0.65);
515   min-width: 0;
516 }
517
518 .metric-label {
519   display: block;
520   color: var(--muted);
521   margin-bottom: 8px;
522 }
523
524 .metric strong {
525   display: block;
526   font-size: clamp(1.3rem, 1.9vw, 1.7rem);
527   min-width: 0;
528   overflow-wrap: anywhere;
529 }

```

```

530
531 .metric.speedup {
532   background: linear-gradient(135deg, rgba(245, 176, 0, 0.16), rgba(245, 176, 0, 0.03));
533 }
534
535 .chart-card {
536   border: 1px solid var(--line);
537   border-radius: 18px;
538   padding: 16px;
539   min-width: 0;
540 }
541
542 .profile-card {
543   border: 1px solid var(--line);
544   border-radius: 18px;
545   padding: 16px;
546   background: rgba(255, 255, 255, 0.65);
547   min-width: 0;
548 }
549
550 .confidence-card {
551   border: 1px solid var(--line);
552   border-radius: 18px;
553   padding: 16px;
554   background: rgba(255, 255, 255, 0.65);
555   min-width: 0;
556   display: grid;
557   gap: 12px;
558 }
559
560 .confidence-head {
561   display: flex;
562   align-items: baseline;
563   justify-content: space-between;
564   gap: 12px;
565 }
566
567 .confidence-head h3 {
568   margin: 0;
569 }
570
571 .confidence-note {
572   color: var(--muted);
573   font-size: 0.88rem;
574 }
575
576 .confidence-summary {
577   display: grid;
578   grid-template-columns: repeat(3, minmax(0, 1fr));
579   gap: 10px;
580 }
581
582 .confidence-metric {
583   padding: 10px 12px;

```

```

584   border-radius: 14px;
585   background: rgba(245, 176, 0, 0.07);
586 }
587
588 .confidence-metric span {
589   display: block;
590   color: var(--muted);
591   margin-bottom: 6px;
592 }
593
594 .confidence-metric strong {
595   font-size: 1rem;
596 }
597
598 .cpu-top-channels {
599   color: var(--muted);
600   font-size: 0.92rem;
601   line-height: 1.45;
602 }
603
604 .channel-grid {
605   display: grid;
606   gap: 8px;
607 }
608
609 .channel-row {
610   display: grid;
611   grid-template-columns: 32px minmax(0, 1fr) 64px 58px;
612   align-items: center;
613   gap: 10px;
614 }
615
616 .channel-label,
617 .channel-score,
618 .channel-prob {
619   font-size: 0.9rem;
620 }
621
622 .channel-score,
623 .channel-prob {
624   text-align: right;
625   color: var(--muted);
626 }
627
628 .channel-track {
629   height: 12px;
630   border-radius: 999px;
631   background: #ebe2d3;
632   overflow: hidden;
633 }
634
635 .channel-fill {
636   height: 100%;
637   background: linear-gradient(90deg, #2b4d3d, #75a082);

```

```

638     border-radius: 999px;
639 }
640
641 .profile-head {
642     display: flex;
643     align-items: baseline;
644     justify-content: space-between;
645     gap: 12px;
646 }
647
648 .profile-head h3 {
649     margin: 0;
650 }
651
652 .profile-meta {
653     color: var(--muted);
654     font-size: 0.92rem;
655 }
656
657 .profile-grid {
658     display: grid;
659     grid-template-columns: repeat(2, minmax(0, 1fr));
660     gap: 10px 12px;
661     margin-top: 12px;
662 }
663
664 .profile-row {
665     display: flex;
666     align-items: center;
667     justify-content: space-between;
668     gap: 12px;
669     padding: 10px 12px;
670     border-radius: 14px;
671     background: rgba(245, 176, 0, 0.07);
672 }
673
674 .profile-row span {
675     color: var(--muted);
676 }
677
678 .profile-row strong {
679     text-align: right;
680     font-size: 0.94rem;
681 }
682
683 .profile-row-total {
684     grid-column: 1 / -1;
685     background: linear-gradient(135deg, rgba(43, 77, 61, 0.13), rgba(117, 160, 130, 0.05));
686 }
687
688 .chart {
689     display: grid;
690     gap: 14px;
691 }

```

```

692
693 .bar-group {
694     display: grid;
695     grid-template-columns: 56px 1fr 90px;
696     align-items: center;
697     gap: 12px;
698 }
699
700 .bar-track {
701     height: 18px;
702     border-radius: 999px;
703     background: #ebe2d3;
704     overflow: hidden;
705 }
706
707 .bar {
708     height: 100%;
709     width: 0;
710     border-radius: 999px;
711     background: linear-gradient(90deg, #f5b000, #ffd36d);
712 }
713
714 .bar-alt {
715     background: linear-gradient(90deg, #2b4d3d, #75a082);
716 }
717
718 .status-box {
719     padding: 14px 16px;
720     border-radius: 14px;
721     background: #efe6d6;
722     color: var(--success);
723     min-width: 0;
724     width: 100%;
725     max-width: 100%;
726     overflow: hidden;
727 }
728
729 .status-box.error {
730     background: #f7dbd5;
731     color: var(--error);
732 }
733
734 .suite-btn {
735     white-space: nowrap;
736     justify-self: end;
737 }
738
739 @media (max-width: 1180px) {
740     .layout {
741         grid-template-columns: 1fr;
742         align-items: start;
743     }
744
745     .side-column {

```

```

746     gap: 14px;
747     height: auto;
748 }
749
750 .input-panel {
751     height: auto;
752 }
753
754 .suite-panel {
755     margin-top: 0;
756 }
757 }
758
759 @media (max-width: 980px) {
760     .input-top {
761         grid-template-columns: repeat(2, minmax(0, 1fr));
762         grid-template-areas:
763             "webcam preview"
764             "processed processed";
765     }
766
767     .camera-stage,
768     .preview-wrap {
769         min-height: 0;
770     }
771 }
772
773 @media (max-width: 760px) {
774     .page {
775         padding: 24px 16px 36px;
776     }
777
778     .input-top {
779         grid-template-columns: 1fr;
780         grid-template-areas:
781             "webcam"
782             "preview"
783             "processed";
784     }
785
786     .processed-wrap {
787         grid-template-columns: 1fr;
788         grid-template-areas:
789             "label"
790             "preview"
791             "actions";
792     }
793
794     .webcam-head,
795     .feedback-head,
796     .preview-head {
797         flex-direction: column;
798         align-items: flex-start;
799     }

```

```

800
801 .suite-head {
802     grid-template-columns: 1fr;
803 }
804
805 .profile-grid {
806     grid-template-columns: 1fr;
807 }
808
809 .confidence-summary {
810     grid-template-columns: 1fr;
811 }
812
813 .metrics {
814     grid-template-columns: 1fr;
815 }
816
817 .prediction-strip {
818     grid-template-columns: 1fr;
819 }
820
821 .feedback-controls {
822     grid-template-columns: 1fr;
823 }
824
825 .preprocess-card {
826     grid-template-columns: 1fr;
827 }
828
829 .webcam-actions,
830 .actions {
831     width: 100%;
832 }
833
834 .webcam-actions button,
835 .actions button,
836 .suite-btn {
837     width: 100%;
838 }
839
840 .preview-note {
841     text-align: left;
842 }
843 }

```

## A.4 HPS Board C Runtime

### A.4.1 MMIO Register ABI

*Defines the software-visible register numbers, bit fields, control flags, and default scratchpad layout.*

```

1  #ifndef CNN_MMIO_REGS_H
2  #define CNN_MMIO_REGS_H
3
4  #include <stdint.h>
5
6  /*
7   * Shared register/memory map for cnn_mmio_interface.
8   *
9   * This header is the ABI contract between the HPS C programs and the FPGA RTL.
10  * If a register index, bit field, or scratchpad address changes here, the
11  * matching decode logic in cnn_mmio_interface.v must change in the same way.
12  *
13  * The RTL uses address[12] to split the space:
14  *   0x0000..0x0FFF : 16 KiB 32-bit scratchpad memory space
15  *   0x1000..0x101F : 32-bit config/status register space
16  *
17  * Userspace/HPS code should address 32-bit words, not bytes.
18  */
19
20 #define CNN_MMIO_SCRATCHPAD_AW      12u
21 #define CNN_MMIO_SCRATCHPAD_WORDS  (1u << CNN_MMIO_SCRATCHPAD_AW)
22 #define CNN_MMIO_SCRATCHPAD_MASK  (CNN_MMIO_SCRATCHPAD_WORDS - 1u)
23
24 #define CNN_MMIO_MEM_SPACE_BIT      (0u << CNN_MMIO_SCRATCHPAD_AW)
25 #define CNN_MMIO_CFG_SPACE_BIT     (1u << CNN_MMIO_SCRATCHPAD_AW)
26
27 #define CNN_MMIO_REG_CONTROL        0u
28 #define CNN_MMIO_REG_STATUS         1u
29 #define CNN_MMIO_REG_CONV_CFG_BASE  2u
30 #define CNN_MMIO_REG_CONV_CFG_LEN   3u
31 #define CNN_MMIO_REG_CONV_WT_BASE   4u
32 #define CNN_MMIO_REG_CONV_WT_LEN   5u
33 #define CNN_MMIO_REG_FC_BIAS_BASE   6u
34 #define CNN_MMIO_REG_FC_BIAS_LEN    7u
35 #define CNN_MMIO_REG_FCW_BASE       8u
36 #define CNN_MMIO_REG_FCW_LEN        9u
37 #define CNN_MMIO_REG_IMAGE_BASE     10u
38 #define CNN_MMIO_REG_IMAGE_LEN      11u
39 #define CNN_MMIO_REG_PREDICT        12u
40 #define CNN_MMIO_REG_IF_ERROR       13u
41 #define CNN_MMIO_REG_PROFILE_L1_LO   14u
42 #define CNN_MMIO_REG_PROFILE_L1_HI   15u
43 #define CNN_MMIO_REG_PROFILE_L2_P0_LO 16u
44 #define CNN_MMIO_REG_PROFILE_L2_P0_HI 17u
45 #define CNN_MMIO_REG_PROFILE_L2_P1_LO 18u
46 #define CNN_MMIO_REG_PROFILE_L2_P1_HI 19u
47 #define CNN_MMIO_REG_PROFILE_L3_P0_LO 20u
48 #define CNN_MMIO_REG_PROFILE_L3_P0_HI 21u
49 #define CNN_MMIO_REG_PROFILE_L3_P1_LO 22u
50 #define CNN_MMIO_REG_PROFILE_L3_P1_HI 23u
51 #define CNN_MMIO_REG_PROFILE_FC_LO   24u
52 #define CNN_MMIO_REG_PROFILE_FC_HI   25u

```

```

53 #define CNN_MMIO_REG_PROFILE_ARGMAX_LO 26u
54 #define CNN_MMIO_REG_PROFILE_ARGMAX_HI 27u
55 #define CNN_MMIO_REG_PROFILE_TOTAL_LO 28u
56 #define CNN_MMIO_REG_PROFILE_TOTAL_HI 29u
57 /* Optional debug/readback helpers. Not used by the normal host flow. */
58 #define CNN_MMIO_REG_LAST_WRITE 30u
59 #define CNN_MMIO_REG_MAGIC 31u
60
61 #define CNN_MMIO_CTRL_MODEL_LOAD 0x0001u
62 #define CNN_MMIO_CTRL_INFER 0x0002u
63 #define CNN_MMIO_CTRL_CLEAR_STATUS 0x0004u
64
65 #define CNN_MMIO_STATUS_BUSY_SHIFT 1u
66 #define CNN_MMIO_STATUS_MODEL_LOADED_SHIFT 2u
67 #define CNN_MMIO_STATUS_PREDICT_DONE_SHIFT 3u
68 #define CNN_MMIO_STATUS_PREDICT_SHIFT 4u
69
70 /*
71  * Default staged memory layout for the current 64x64x1 -> 10-class model.
72  *
73  * All *_WORDS values count 32-bit host words. All *_BASE_W values are
74  * 32-bit word addresses because the MMIO scratchpad is word-indexed.
75  *
76  * Model structure behind these constants:
77  * - L1 conv: 3x3x1 -> DOT_K = 9
78  * - L2 conv: 3x3x4 -> DOT_K = 36, split into pass0/pass1
79  * - L3 conv: 3x3x8 -> DOT_K = 72, split into pass0/pass1
80  * - FC: 6x6x8 = 288 inputs -> 10 outputs
81  *
82  * Derived word counts:
83  * - CONV_CFG_WORDS = 5 passes x 9 cfg words/pass = 45
84  *   (L1, L2p0, L2p1, L3p0, L3p1)
85  * - CONV_WT_WORDS = 9 + 36 + 36 + 72 + 72 = 225
86  * - FC_BIAS_WORDS = 10 output classes = 10
87  * - FCW_WORDS = 288 FC positions x 3 host words/position = 864
88  *   Each FC position stores 10x8b = 80b of weights, packed from 3x32b.
89  * - IMAGE_WORDS = 64x64x1 int8 pixels / 4 pixels per 32-bit word = 1024
90  *
91  * Base addresses are packed contiguously in word units:
92  * conv_cfg @ 0
93  * conv_wt @ 0 + 45 = 45
94  * fc_bias @ 45 + 225 = 270
95  * fcw @ 270 + 10 = 280
96  * image @ 280 + 864 = 1144
97  */
98 #define CNN_MMIO_DEFAULT_CONV_CFG_BASE_W 0u
99 #define CNN_MMIO_DEFAULT_CONV_CFG_WORDS 45u
100 #define CNN_MMIO_DEFAULT_CONV_WT_BASE_W 45u
101 #define CNN_MMIO_DEFAULT_CONV_WT_WORDS 225u
102 #define CNN_MMIO_DEFAULT_FC_BIAS_BASE_W 270u
103 #define CNN_MMIO_DEFAULT_FC_BIAS_WORDS 10u
104 #define CNN_MMIO_DEFAULT_FCW_BASE_W 280u
105 #define CNN_MMIO_DEFAULT_FCW_WORDS 864u
106 #define CNN_MMIO_DEFAULT_IMAGE_BASE_W 1144u

```

```

107 #define CNN_MMIO_DEFAULT_IMAGE_WORDS      1024u
108
109 static inline uint32_t cnn_mmio_cfg_addr(uint32_t reg_idx) {
110     /* Config/status registers live in the address[12] == 1 half of the window. */
111     return CNN_MMIO_CFG_SPACE_BIT | (reg_idx & 0x1Fu);
112 }
113
114 static inline uint32_t cnn_mmio_mem_addr(uint32_t word_addr) {
115     /* Scratchpad addresses are 32-bit word indices written by the HPS. */
116     return CNN_MMIO_MEM_SPACE_BIT | (word_addr & CNN_MMIO_SCRATCHPAD_MASK);
117 }
118
119 static inline uint32_t cnn_mmio_pack_status_predict(uint16_t status_word) {
120     return (status_word >> CNN_MMIO_STATUS_PREDICT_SHIFT) & 0xFu;
121 }
122
123 static inline uint32_t cnn_mmio_status_busy(uint16_t status_word) {
124     return (status_word >> CNN_MMIO_STATUS_BUSY_SHIFT) & 0x1u;
125 }
126
127 static inline uint32_t cnn_mmio_status_model_loaded(uint16_t status_word) {
128     return (status_word >> CNN_MMIO_STATUS_MODEL_LOADED_SHIFT) & 0x1u;
129 }
130
131 static inline uint32_t cnn_mmio_status_predict_done(uint16_t status_word) {
132     return (status_word >> CNN_MMIO_STATUS_PREDICT_DONE_SHIFT) & 0x1u;
133 }
134
135 #endif

```

## A.4.2 HPS MMIO Helper API

*Declares the shared C structs and helper functions used by all board-side MMIO commands.*

```
code/include/cnn_mmio_host.h
```

```

1  #ifndef CNN_MMIO_HOST_H
2  #define CNN_MMIO_HOST_H
3
4  #include <stdint.h>
5
6  #include "cnn_mmio_regs.h"
7
8  #ifdef __cplusplus
9  extern "C" {
10 #endif
11
12 /*
13  * The HPS maps one small page-aligned lightweight-bridge window. The RTL uses
14  * 16 KiB of scratchpad, and the highest 32-bit config/status register ends at
15  * byte offset 0x407F. A 20 KiB span is the smallest 4 KiB page-aligned mapping

```

```

16  * that covers the implemented scratchpad and register window.
17  */
18  #define CNN_MMIO_MAP_SPAN_BYTES (20 * 1024)
19  #define CNN_MMIO_DEFAULT_TIMEOUT_MS 1000
20  #define CNN_MMIO_DEFAULT_FABRIC_MHZ 50.0
21
22  struct cnn_mmio_preload_bundle {
23      int32_t conv_cfg[CNN_MMIO_DEFAULT_CONV_CFG_WORDS];
24      uint32_t conv_wt[CNN_MMIO_DEFAULT_CONV_WT_WORDS];
25      int32_t fc_bias[CNN_MMIO_DEFAULT_FC_BIAS_WORDS];
26      int32_t fcw[CNN_MMIO_DEFAULT_FCW_WORDS];
27  };
28
29  struct cnn_mmio_inference_case {
30      uint32_t image[CNN_MMIO_DEFAULT_IMAGE_WORDS];
31      int expected_class;
32  };
33
34  struct cnn_mmio_device {
35      int fd;
36      void *map_base;
37      volatile uint32_t *mmio_base;
38      uintptr_t csr_base;
39  };
40
41  struct cnn_mmio_profile {
42      uint32_t l1_cycles;
43      uint32_t l2_p0_cycles;
44      uint32_t l2_p1_cycles;
45      uint32_t l3_p0_cycles;
46      uint32_t l3_p1_cycles;
47      uint32_t fc_cycles;
48      uint32_t argmax_cycles;
49      uint32_t total_cycles;
50  };
51
52  int cnn_mmio_load_preload_bundle(const char *preload_root, struct cnn_mmio_preload_bundle *bundle);
53  int cnn_mmio_load_inference_case(const char *case_root, struct cnn_mmio_inference_case *tc);
54
55  int cnn_mmio_open(struct cnn_mmio_device *dev, uintptr_t csr_base, const char *devmem_path);
56  void cnn_mmio_close(struct cnn_mmio_device *dev);
57
58  void cnn_mmio_program_default_registers(volatile uint32_t *mmio_base);
59  void cnn_mmio_write_preload_bundle(volatile uint32_t *mmio_base, const struct cnn_mmio_preload_bundle *
    bundle);
60  void cnn_mmio_write_inference_case(volatile uint32_t *mmio_base, const struct cnn_mmio_inference_case *tc
    );
61
62  uint16_t cnn_mmio_read_status(volatile uint32_t *mmio_base);
63  uint16_t cnn_mmio_read_error(volatile uint32_t *mmio_base);
64  uint16_t cnn_mmio_read_predict(volatile uint32_t *mmio_base);
65  void cnn_mmio_read_profile(volatile uint32_t *mmio_base, struct cnn_mmio_profile *profile);
66
67  void cnn_mmio_clear_status(volatile uint32_t *mmio_base);

```

```

68 void cnn_mmio_start_model_load(volatile uint32_t *mmio_base);
69 void cnn_mmio_start_infer(volatile uint32_t *mmio_base);
70
71 int cnn_mmio_wait_for_status_bit(
72     volatile uint32_t *mmio_base,
73     unsigned bit_idx,
74     unsigned expected_value,
75     int timeout_ms,
76     uint16_t *last_status);
77
78 #ifdef __cplusplus
79 }
80 #endif
81
82 #endif

```

### A.4.3 Common MMIO Implementation

*Maps /dev/mem, writes preload/image words, starts model load or inference, polls status, and reads profile counters.*

code/tools/cnn\_mmio\_host.c

```

1 #define _POSIX_C_SOURCE 200809L
2
3 #include "../include/cnn_mmio_host.h"
4
5 #include <errno.h>
6 #include <fcntl.h>
7 #include <inttypes.h>
8 #include <stdio.h>
9 #include <stdlib.h>
10 #include <string.h>
11 #include <sys/mman.h>
12 #include <sys/stat.h>
13 #include <sys/types.h>
14 #include <time.h>
15 #include <unistd.h>
16
17 static int read_i32_lines(const char *path, int32_t *dst, size_t count) {
18     FILE *fp = fopen(path, "r");
19     size_t i;
20     if (!fp) {
21         perror(path);
22         return -1;
23     }
24     for (i = 0; i < count; ++i) {
25         if (fscanf(fp, "%" SCNd32, &dst[i]) != 1) {
26             fprintf(stderr, "short read in %s at index %zu\n", path, i);
27             fclose(fp);

```

```

28     return -1;
29 }
30 }
31 fclose(fp);
32 return 0;
33 }
34
35 static int read_packed_i8_words(const char *path, uint32_t *dst, size_t word_count) {
36     FILE *fp = fopen(path, "r");
37     size_t i;
38     if (!fp) {
39         perror(path);
40         return -1;
41     }
42     for (i = 0; i < word_count; ++i) {
43         int32_t b0, b1, b2, b3;
44         if (fscanf(fp, "%" SCNd32, &b0) != 1 ||
45             fscanf(fp, "%" SCNd32, &b1) != 1 ||
46             fscanf(fp, "%" SCNd32, &b2) != 1 ||
47             fscanf(fp, "%" SCNd32, &b3) != 1) {
48             fprintf(stderr, "short packed read in %s at word %zu\n", path, i);
49             fclose(fp);
50             return -1;
51         }
52         dst[i] = ((uint32_t)(b0 & 0xFF)) |
53                 ((uint32_t)(b1 & 0xFF) << 8) |
54                 ((uint32_t)(b2 & 0xFF) << 16) |
55                 ((uint32_t)(b3 & 0xFF) << 24);
56     }
57     fclose(fp);
58     return 0;
59 }
60
61 static int load_manifest_expected_class(const char *path, int *expected_class) {
62     FILE *fp = fopen(path, "r");
63     char line[256];
64     if (!fp) {
65         perror(path);
66         return -1;
67     }
68     while (fgets(line, sizeof(line), fp)) {
69         if (strncmp(line, "predict_class=", 14) == 0) {
70             *expected_class = atoi(line + 14);
71             fclose(fp);
72             return 0;
73         }
74     }
75     fclose(fp);
76     fprintf(stderr, "predict_class not found in %s\n", path);
77     return -1;
78 }
79
80 static void build_path(char *dst, size_t dst_len, const char *root, const char *leaf) {
81     snprintf(dst, dst_len, "%s/%s", root, leaf);

```

```

82 }
83
84 static int try_read_i32_lines(const char *root, const char *const *leaves,
85                             int32_t *dst, size_t count) {
86     char path[1024];
87     size_t i;
88     for (i = 0; leaves[i] != NULL; ++i) {
89         build_path(path, sizeof(path), root, leaves[i]);
90         if (read_i32_lines(path, dst, count) == 0)
91             return 0;
92     }
93     return -1;
94 }
95
96 static int try_read_packed_i8_words(const char *root, const char *const *leaves,
97                                    uint32_t *dst, size_t word_count) {
98     char path[1024];
99     size_t i;
100    for (i = 0; leaves[i] != NULL; ++i) {
101        build_path(path, sizeof(path), root, leaves[i]);
102        if (read_packed_i8_words(path, dst, word_count) == 0)
103            return 0;
104    }
105    return -1;
106 }
107
108 static void mmio_write_cfg_reg(volatile uint32_t *base, uint32_t reg_idx, uint32_t value) {
109     base[cnn_mmio_cfg_addr(reg_idx)] = value;
110 }
111
112 static uint32_t mmio_read_cfg_reg(volatile uint32_t *base, uint32_t reg_idx) {
113     return base[cnn_mmio_cfg_addr(reg_idx)];
114 }
115
116 static void write_word_image(volatile uint32_t *base, uint32_t start_word,
117                             const uint32_t *words, size_t word_count) {
118     size_t i;
119
120     for (i = 0; i < word_count; ++i) {
121         base[cnn_mmio_mem_addr(start_word + (uint32_t)i)] = words[i];
122     }
123 }
124
125 int cnn_mmio_load_preload_bundle(const char *preload_root, struct cnn_mmio_preload_bundle *bundle) {
126     static const char *const conv_cfg_candidates[] = {
127         "preload_conv_cfg_45w.txt",
128         "conv_cfg_words.txt",
129         NULL};
130     static const char *const conv_wt_candidates[] = {
131         "preload_conv_wt_225w_bytes.txt",
132         "conv_wt_words.txt",
133         NULL};
134     static const char *const fc_bias_candidates[] = {
135         "preload_fc_bias_10w.txt",

```

```

136     "fc_bias_words.txt",
137     NULL};
138 static const char *const fcw_candidates[] = {
139     "preload_fcw_864w.txt",
140     "fcw_words.txt",
141     NULL};
142
143 if (try_read_i32_lines(preload_root, conv_cfg_candidates, bundle->conv_cfg,
144     CNN_MMIO_DEFAULT_CONV_CFG_WORDS) != 0)
145     return -1;
146
147 if (try_read_packed_i8_words(preload_root, conv_wt_candidates, bundle->conv_wt,
148     CNN_MMIO_DEFAULT_CONV_WT_WORDS) != 0)
149     return -1;
150
151 if (try_read_i32_lines(preload_root, fc_bias_candidates, bundle->fc_bias,
152     CNN_MMIO_DEFAULT_FC_BIAS_WORDS) != 0)
153     return -1;
154
155 if (try_read_i32_lines(preload_root, fcw_candidates, (int32_t *)bundle->fcw,
156     CNN_MMIO_DEFAULT_FCW_WORDS) != 0)
157     return -1;
158
159 return 0;
160 }
161
162 int cnn_mmio_load_inference_case(const char *case_root, struct cnn_mmio_inference_case *tc) {
163     char path[1024];
164     static const char *const image_candidates[] = {
165         "tb_conv1_in_i8_64x64x1.txt",
166         "image_words.txt",
167         NULL};
168
169 if (try_read_packed_i8_words(case_root, image_candidates, tc->image,
170     CNN_MMIO_DEFAULT_IMAGE_WORDS) != 0)
171     return -1;
172
173 build_path(path, sizeof(path), case_root, "manifest.txt");
174 if (load_manifest_expected_class(path, &tc->expected_class) != 0)
175     return -1;
176
177 return 0;
178 }
179
180 int cnn_mmio_open(struct cnn_mmio_device *dev, uintptr_t csr_base, const char *devmem_path) {
181     memset(dev, 0, sizeof(*dev));
182     dev->fd = open(devmem_path, O_RDWR | O_SYNC);
183     if (dev->fd < 0) {
184         perror(devmem_path);
185         return -1;
186     }
187
188     dev->map_base = mmap(NULL, CNN_MMIO_MAP_SPAN_BYTES, PROT_READ | PROT_WRITE,
189         MAP_SHARED, dev->fd, csr_base);

```

```

190     if (dev->map_base == MAP_FAILED) {
191         perror("mmap");
192         close(dev->fd);
193         dev->fd = -1;
194         return -1;
195     }
196
197     dev->mmio_base = (volatile uint32_t *)dev->map_base;
198     dev->csr_base = csr_base;
199     return 0;
200 }
201
202 void cnn_mmio_close(struct cnn_mmio_device *dev) {
203     if (dev->map_base && dev->map_base != MAP_FAILED) {
204         munmap(dev->map_base, CNN_MMIO_MAP_SPAN_BYTES);
205     }
206     if (dev->fd >= 0) {
207         close(dev->fd);
208     }
209     dev->fd = -1;
210     dev->map_base = NULL;
211     dev->mmio_base = NULL;
212     dev->csr_base = 0;
213 }
214
215 void cnn_mmio_program_default_registers(volatile uint32_t *mmio_base) {
216     mmio_write_cfg_reg(mmio_base, CNN_MMIO_REG_CONV_CFG_BASE, CNN_MMIO_DEFAULT_CONV_CFG_BASE_W);
217     mmio_write_cfg_reg(mmio_base, CNN_MMIO_REG_CONV_CFG_LEN, CNN_MMIO_DEFAULT_CONV_CFG_WORDS);
218     mmio_write_cfg_reg(mmio_base, CNN_MMIO_REG_CONV_WT_BASE, CNN_MMIO_DEFAULT_CONV_WT_BASE_W);
219     mmio_write_cfg_reg(mmio_base, CNN_MMIO_REG_CONV_WT_LEN, CNN_MMIO_DEFAULT_CONV_WT_WORDS);
220     mmio_write_cfg_reg(mmio_base, CNN_MMIO_REG_FC_BIAS_BASE, CNN_MMIO_DEFAULT_FC_BIAS_BASE_W);
221     mmio_write_cfg_reg(mmio_base, CNN_MMIO_REG_FC_BIAS_LEN, CNN_MMIO_DEFAULT_FC_BIAS_WORDS);
222     mmio_write_cfg_reg(mmio_base, CNN_MMIO_REG_FCW_BASE, CNN_MMIO_DEFAULT_FCW_BASE_W);
223     mmio_write_cfg_reg(mmio_base, CNN_MMIO_REG_FCW_LEN, CNN_MMIO_DEFAULT_FCW_WORDS);
224     mmio_write_cfg_reg(mmio_base, CNN_MMIO_REG_IMAGE_BASE, CNN_MMIO_DEFAULT_IMAGE_BASE_W);
225     mmio_write_cfg_reg(mmio_base, CNN_MMIO_REG_IMAGE_LEN, CNN_MMIO_DEFAULT_IMAGE_WORDS);
226 }
227
228 void cnn_mmio_write_preload_bundle(volatile uint32_t *mmio_base, const struct cnn_mmio_preload_bundle *
    bundle) {
229     uint32_t conv_cfg_base = mmio_read_cfg_reg(mmio_base, CNN_MMIO_REG_CONV_CFG_BASE);
230     uint32_t conv_wt_base = mmio_read_cfg_reg(mmio_base, CNN_MMIO_REG_CONV_WT_BASE);
231     uint32_t fc_bias_base = mmio_read_cfg_reg(mmio_base, CNN_MMIO_REG_FC_BIAS_BASE);
232     uint32_t fcw_base = mmio_read_cfg_reg(mmio_base, CNN_MMIO_REG_FCW_BASE);
233
234     write_word_image(mmio_base, conv_cfg_base, (const uint32_t *)bundle->conv_cfg,
235                     CNN_MMIO_DEFAULT_CONV_CFG_WORDS);
236     write_word_image(mmio_base, conv_wt_base, bundle->conv_wt,
237                     CNN_MMIO_DEFAULT_CONV_WT_WORDS);
238     write_word_image(mmio_base, fc_bias_base, (const uint32_t *)bundle->fc_bias,
239                     CNN_MMIO_DEFAULT_FC_BIAS_WORDS);
240     write_word_image(mmio_base, fcw_base, (const uint32_t *)bundle->fcw,
241                     CNN_MMIO_DEFAULT_FCW_WORDS);
242 }

```

```

243
244 void cnn_mmio_write_inference_case(volatile uint32_t *mmio_base, const struct cnn_mmio_inference_case *tc
    ) {
245     uint32_t image_base = mmio_read_cfg_reg(mmio_base, CNN_MMIO_REG_IMAGE_BASE);
246
247     write_word_image(mmio_base, image_base, tc->image,
248                     CNN_MMIO_DEFAULT_IMAGE_WORDS);
249 }
250
251 uint16_t cnn_mmio_read_status(volatile uint32_t *mmio_base) {
252     return (uint16_t)(mmio_read_cfg_reg(mmio_base, CNN_MMIO_REG_STATUS) & 0xFFFFu);
253 }
254
255 uint16_t cnn_mmio_read_error(volatile uint32_t *mmio_base) {
256     return (uint16_t)(mmio_read_cfg_reg(mmio_base, CNN_MMIO_REG_IF_ERROR) & 0xFFFFu);
257 }
258
259 uint16_t cnn_mmio_read_predict(volatile uint32_t *mmio_base) {
260     return (uint16_t)cnn_mmio_pack_status_predict(cnn_mmio_read_status(mmio_base));
261 }
262
263 void cnn_mmio_read_profile(volatile uint32_t *mmio_base, struct cnn_mmio_profile *profile) {
264     if (!profile)
265         return;
266
267     profile->l1_cycles = mmio_read_cfg_reg(mmio_base, CNN_MMIO_REG_PROFILE_L1_HI);
268     profile->l2_p0_cycles = mmio_read_cfg_reg(mmio_base, CNN_MMIO_REG_PROFILE_L2_P0_HI);
269     profile->l2_p1_cycles = mmio_read_cfg_reg(mmio_base, CNN_MMIO_REG_PROFILE_L2_P1_HI);
270     profile->l3_p0_cycles = mmio_read_cfg_reg(mmio_base, CNN_MMIO_REG_PROFILE_L3_P0_HI);
271     profile->l3_p1_cycles = mmio_read_cfg_reg(mmio_base, CNN_MMIO_REG_PROFILE_L3_P1_HI);
272     profile->fc_cycles = mmio_read_cfg_reg(mmio_base, CNN_MMIO_REG_PROFILE_FC_HI);
273     profile->argmax_cycles = mmio_read_cfg_reg(mmio_base, CNN_MMIO_REG_PROFILE_ARGMAX_HI);
274     profile->total_cycles = mmio_read_cfg_reg(mmio_base, CNN_MMIO_REG_PROFILE_TOTAL_HI);
275 }
276
277 void cnn_mmio_clear_status(volatile uint32_t *mmio_base) {
278     mmio_write_cfg_reg(mmio_base, CNN_MMIO_REG_CONTROL, CNN_MMIO_CTRL_CLEAR_STATUS);
279 }
280
281 void cnn_mmio_start_model_load(volatile uint32_t *mmio_base) {
282     mmio_write_cfg_reg(mmio_base, CNN_MMIO_REG_CONTROL, CNN_MMIO_CTRL_MODEL_LOAD);
283 }
284
285 void cnn_mmio_start_infer(volatile uint32_t *mmio_base) {
286     mmio_write_cfg_reg(mmio_base, CNN_MMIO_REG_CONTROL, CNN_MMIO_CTRL_INFER);
287 }
288
289 int cnn_mmio_wait_for_status_bit(
290     volatile uint32_t *mmio_base,
291     unsigned bit_idx,
292     unsigned expected_value,
293     int timeout_ms,
294     uint16_t *last_status) {
295     struct timespec req;

```

```

296 int elapsed_ms = 0;
297 if (bit_idx >= 16)
298     return -1;
299
300 req.tv_sec = 0;
301 req.tv_nsec = 1000000L;
302
303 while (elapsed_ms < timeout_ms) {
304     uint16_t status = cnn_mmio_read_status(mmio_base);
305     if (last_status)
306         *last_status = status;
307     if (((status >> bit_idx) & 0x1u) == (expected_value & 0x1u))
308         return 0;
309     nanosleep(&req, NULL);
310     elapsed_ms += 1;
311 }
312
313 return -1;
314 }

```

#### A.4.4 Status Probe Command

*Reads STATUS, PREDICT, IF\_ERROR, and profile registers for health checks before web requests.*

code/tools/hps\_mmio\_status.c

```

1 #include <stdint.h>
2 #include <stdio.h>
3 #include <stdlib.h>
4
5 #include "../include/cnn_mmio_host.h"
6
7 int main(int argc, char **argv) {
8     const char *devmem_path = "/dev/mem";
9     uintptr_t csr_base;
10    struct cnn_mmio_device dev;
11    uint16_t status;
12    uint16_t predict;
13    uint16_t error_reg;
14
15    if (argc < 2 || argc > 3) {
16        fprintf(stderr, "usage: %s <csr_base_hex> [devmem_path]\n", argv[0]);
17        return 1;
18    }
19
20    csr_base = (uintptr_t)strtoull(argv[1], NULL, 0);
21    if (argc > 2)
22        devmem_path = argv[2];
23
24    if (cnn_mmio_open(&dev, csr_base, devmem_path) != 0)
25        return 1;

```

```

26
27     status = cnn_mmio_read_status(dev.mmio_base);
28     predict = cnn_mmio_read_predict(dev.mmio_base);
29     error_reg = cnn_mmio_read_error(dev.mmio_base);
30
31     printf("status=0x%04x\n", status);
32     printf("model_loaded=%u\n", (unsigned)cnn_mmio_status_model_loaded(status));
33     printf("predict_done=%u\n", (unsigned)cnn_mmio_status_predict_done(status));
34     printf("predict_class=%u\n", (unsigned)(predict & 0xF));
35     printf("error=0x%04x\n", error_reg);
36
37     cnn_mmio_close(&dev);
38     return 0;
39 }

```

#### A.4.5 Model Load Command

*Loads the preload bundle into the scratchpad and triggers CONTROL[0] model replay.*

```
code/tools/hps_mmio_load_model.c
```

```

1  #include <stdint.h>
2  #include <stdio.h>
3  #include <stdlib.h>
4
5  #include "../include/cnn_mmio_host.h"
6
7  int main(int argc, char **argv) {
8     const char *devmem_path = "/dev/mem";
9     const char *preload_root;
10    uintptr_t csr_base;
11    uint16_t status = 0;
12    struct cnn_mmio_device dev;
13    struct cnn_mmio_preload_bundle preload;
14
15    if (argc < 3 || argc > 4) {
16        fprintf(stderr,
17            "usage: %s <csr_base_hex> <preload_root> [devmem_path]\n",
18            argv[0]);
19        return 1;
20    }
21
22    csr_base = (uintptr_t)strtoull(argv[1], NULL, 0);
23    preload_root = argv[2];
24    if (argc > 3)
25        devmem_path = argv[3];
26
27    if (cnn_mmio_load_preload_bundle(preload_root, &preload) != 0)
28        return 1;
29    if (cnn_mmio_open(&dev, csr_base, devmem_path) != 0)
30        return 1;

```

```

31
32 cnn_mmio_program_default_registers(dev.mmio_base);
33 cnn_mmio_write_preload_bundle(dev.mmio_base, &preload);
34 cnn_mmio_clear_status(dev.mmio_base);
35 cnn_mmio_start_model_load(dev.mmio_base);
36
37 if (cnn_mmio_wait_for_status_bit(
38     dev.mmio_base,
39     CNN_MMIO_STATUS_MODEL_LOADED_SHIFT,
40     1,
41     CNN_MMIO_DEFAULT_TIMEOUT_MS,
42     &status) != 0) {
43     fprintf(stderr, "timeout waiting for model_loaded, status=0x%04x\n", status);
44     cnn_mmio_close(&dev);
45     return 1;
46 }
47
48 printf("status=0x%04x\n", status);
49 printf("model_loaded=%u\n", (unsigned)((status >> CNN_MMIO_STATUS_MODEL_LOADED_SHIFT) & 0x1));
50 cnn_mmio_close(&dev);
51 return 0;
52 }

```

#### A.4.6 FPGA Prediction Command

*Writes one exported image case, triggers CONTROL[1] inference, waits for predict\_done, and prints web-readable key-value output.*

```
code/tools/hps_mmio_predict.c
```

```

1  #define _POSIX_C_SOURCE 200809L
2
3  #include <stdint.h>
4  #include <stdio.h>
5  #include <stdlib.h>
6  #include <time.h>
7
8  #include "../include/cnn_mmio_host.h"
9
10 static double monotonic_ms(void) {
11     struct timespec ts;
12     clock_gettime(CLOCK_MONOTONIC, &ts);
13     return (double)ts.tv_sec * 1000.0 + (double)ts.tv_nsec / 1000000.0;
14 }
15
16 int main(int argc, char **argv) {
17     const char *devmem_path = "/dev/mem";
18     const char *case_root;
19     uintptr_t csr_base;
20     uint16_t status = 0;

```

```

21  uint16_t error_reg = 0;
22  struct cnn_mmio_device dev;
23  struct cnn_mmio_inference_case tc;
24  struct cnn_mmio_profile profile;
25  double total_started_ms;
26  double case_load_started_ms;
27  double case_load_ended_ms;
28  double mmio_program_started_ms;
29  double mmio_program_ended_ms;
30  double infer_started_ms;
31  double infer_ended_ms;
32
33  if (argc < 3 || argc > 4) {
34      fprintf(stderr, "usage: %s <csr_base_hex> <case_root> [devmem_path]\n", argv[0]);
35      return 1;
36  }
37
38  csr_base = (uintptr_t)strtoull(argv[1], NULL, 0);
39  case_root = argv[2];
40  if (argc > 3)
41      devmem_path = argv[3];
42
43  total_started_ms = monotonic_ms();
44  case_load_started_ms = monotonic_ms();
45  if (cnn_mmio_load_inference_case(case_root, &tc) != 0)
46      return 1;
47  case_load_ended_ms = monotonic_ms();
48  if (cnn_mmio_open(&dev, csr_base, devmem_path) != 0)
49      return 1;
50
51  mmio_program_started_ms = monotonic_ms();
52  cnn_mmio_program_default_registers(dev.mmio_base);
53  cnn_mmio_write_inference_case(dev.mmio_base, &tc);
54  mmio_program_ended_ms = monotonic_ms();
55
56  infer_started_ms = monotonic_ms();
57  cnn_mmio_start_infer(dev.mmio_base);
58
59  if (cnn_mmio_wait_for_status_bit(
60      dev.mmio_base,
61      CNN_MMIO_STATUS_PREDICT_DONE_SHIFT,
62      1,
63      CNN_MMIO_DEFAULT_TIMEOUT_MS,
64      &status) != 0) {
65      fprintf(stderr, "timeout waiting for predict_done, status=0x%04x\n", status);
66      cnn_mmio_close(&dev);
67      return 1;
68  }
69  infer_ended_ms = monotonic_ms();
70
71  error_reg = cnn_mmio_read_error(dev.mmio_base);
72  cnn_mmio_read_profile(dev.mmio_base, &profile);
73  printf("predict_class=%u\n", (unsigned)cnn_mmio_pack_status_predict(status));
74  printf("status=0x%04x\n", status);

```

```

75  printf("error=0x%04x\n", error_reg);
76  printf("l1_cycles=%u\n", profile.l1_cycles);
77  printf("l2_p0_cycles=%u\n", profile.l2_p0_cycles);
78  printf("l2_p1_cycles=%u\n", profile.l2_p1_cycles);
79  printf("l3_p0_cycles=%u\n", profile.l3_p0_cycles);
80  printf("l3_p1_cycles=%u\n", profile.l3_p1_cycles);
81  printf("fc_cycles=%u\n", profile.fc_cycles);
82  printf("argmax_cycles=%u\n", profile.argmax_cycles);
83  printf("total_cycles=%u\n", profile.total_cycles);
84  printf("board_case_load_ms=%.3f\n", case_load_ended_ms - case_load_started_ms);
85  printf("board_program_ms=%.3f\n", mmio_program_ended_ms - mmio_program_started_ms);
86  printf("board_infer_wait_ms=%.3f\n", infer_ended_ms - infer_started_ms);
87  printf("board_total_ms=%.3f\n", infer_ended_ms - total_started_ms);
88
89  cnn_mmio_close(&dev);
90  return (error_reg == 0) ? 0 : 1;
91 }

```

#### A.4.7 Board CPU Baseline

*Runs the same exported case on the HPS CPU baseline so the web demo can compare CPU and FPGA results.*

code/tools/hps\_cpu\_reference.c

```

1  #define _POSIX_C_SOURCE 200809L
2
3  #include <stdint.h>
4  #include <stdio.h>
5  #include <stdlib.h>
6  #include <string.h>
7  #include <time.h>
8
9  #define IMG_H 64
10 #define IMG_W 64
11 #define C1_OUT 4
12 #define C2_OUT 8
13 #define C3_OUT 8
14 #define FC_OUT 10
15
16 #define C1_H 62
17 #define C1_W 62
18 #define P1_H 31
19 #define P1_W 31
20 #define C2_H 29
21 #define C2_W 29
22 #define P2_H 14
23 #define P2_W 14
24 #define C3_H 12
25 #define C3_W 12

```

```

26 #define P3_H 6
27 #define P3_W 6
28 #define FC_IN 288
29
30 static double monotonic_ms(void) {
31     struct timespec ts;
32     clock_gettime(CLOCK_MONOTONIC, &ts);
33     return (double)ts.tv_sec * 1000.0 + (double)ts.tv_nsec / 1000000.0;
34 }
35
36 static int read_i8_lines(const char *path, int8_t *dst, size_t count) {
37     FILE *fp = fopen(path, "r");
38     size_t i;
39     long v;
40     if (!fp) {
41         perror(path);
42         return -1;
43     }
44     for (i = 0; i < count; ++i) {
45         if (fscanf(fp, "%ld", &v) != 1) {
46             fprintf(stderr, "short read in %s at index %zu\n", path, i);
47             fclose(fp);
48             return -1;
49         }
50         if (v < -128 || v > 127) {
51             fprintf(stderr, "i8 out of range in %s at index %zu: %ld\n", path, i, v);
52             fclose(fp);
53             return -1;
54         }
55         dst[i] = (int8_t)v;
56     }
57     fclose(fp);
58     return 0;
59 }
60
61 static int read_i32_lines(const char *path, int32_t *dst, size_t count) {
62     FILE *fp = fopen(path, "r");
63     size_t i;
64     long v;
65     if (!fp) {
66         perror(path);
67         return -1;
68     }
69     for (i = 0; i < count; ++i) {
70         if (fscanf(fp, "%ld", &v) != 1) {
71             fprintf(stderr, "short read in %s at index %zu\n", path, i);
72             fclose(fp);
73             return -1;
74         }
75         dst[i] = (int32_t)v;
76     }
77     fclose(fp);
78     return 0;
79 }

```

```

80
81 static void build_path(char *dst, size_t dst_len, const char *root, const char *leaf) {
82     snprintf(dst, dst_len, "%s/%s", root, leaf);
83 }
84
85 static int8_t clamp_int8(int64_t v) {
86     if (v > 127)
87         return 127;
88     if (v < -128)
89         return -128;
90     return (int8_t)v;
91 }
92
93 static int idx_hwc(int h, int w, int c, int width, int channels) {
94     return ((h * width) + w) * channels + c;
95 }
96
97 static int idx_oihw(int oc, int ic, int kh, int kw, int cin, int ksize) {
98     return (((oc * cin + ic) * ksize + kh) * ksize + kw);
99 }
100
101 static int idx_ohw(int oc, int kh, int kw, int ksize) {
102     return ((oc * ksize + kh) * ksize + kw);
103 }
104
105 static int32_t requantize_single(int32_t mac,
106                                 int32_t eff_bias,
107                                 int32_t mult,
108                                 int32_t shift,
109                                 int32_t zp_out) {
110     int64_t acc = (int64_t)mac + (int64_t)eff_bias;
111     int64_t prod = acc * (int64_t)mult;
112     int64_t scaled;
113
114     if (shift > 0) {
115         int64_t round_term = (int64_t)1 << (shift - 1);
116         int64_t prod_adj = prod + round_term - (prod < 0 ? 1 : 0);
117         scaled = prod_adj >> shift;
118     } else if (shift == 0) {
119         scaled = prod;
120     } else {
121         scaled = prod << (-shift);
122     }
123
124     return (int32_t)scaled + zp_out;
125 }
126
127 static void conv1_forward(const int8_t *input,
128                           const int8_t *weights,
129                           const int32_t *bias,
130                           const int32_t *mult,
131                           const int32_t *shift,
132                           int8_t *pool_out) {
133     int8_t requant[C1_H * C1_W * C1_OUT];

```

```

134 int h, w, oc, kh, kw;
135
136 for (h = 0; h < C1_H; ++h) {
137     for (w = 0; w < C1_W; ++w) {
138         for (oc = 0; oc < C1_OUT; ++oc) {
139             int32_t mac = 0;
140             for (kh = 0; kh < 3; ++kh) {
141                 for (kw = 0; kw < 3; ++kw) {
142                     int32_t x = input[idx_hwc(h + kh, w + kw, 0, IMG_W, 1)];
143                     int32_t wt = weights[idx_ohw(oc, kh, kw, 3)];
144                     mac += x * wt;
145                 }
146             }
147             requant[idx_hwc(h, w, oc, C1_W, C1_OUT)] =
148                 clamp_int8(requantize_single(mac, bias[oc], mult[oc], shift[oc], -128));
149         }
150     }
151 }
152
153 for (h = 0; h < P1_H; ++h) {
154     for (w = 0; w < P1_W; ++w) {
155         for (oc = 0; oc < C1_OUT; ++oc) {
156             int base_h = h * 2;
157             int base_w = w * 2;
158             int8_t a = requant[idx_hwc(base_h + 0, base_w + 0, oc, C1_W, C1_OUT)];
159             int8_t b = requant[idx_hwc(base_h + 0, base_w + 1, oc, C1_W, C1_OUT)];
160             int8_t c = requant[idx_hwc(base_h + 1, base_w + 0, oc, C1_W, C1_OUT)];
161             int8_t d = requant[idx_hwc(base_h + 1, base_w + 1, oc, C1_W, C1_OUT)];
162             int8_t m1 = a > b ? a : b;
163             int8_t m2 = c > d ? c : d;
164             pool_out[idx_hwc(h, w, oc, P1_W, C1_OUT)] = m1 > m2 ? m1 : m2;
165         }
166     }
167 }
168 }
169
170 static void conv_generic_forward(const int8_t *input,
171                                 int in_h,
172                                 int in_w,
173                                 int in_c,
174                                 const int8_t *weights,
175                                 int out_c,
176                                 const int32_t *bias,
177                                 const int32_t *mult,
178                                 const int32_t *shift,
179                                 int out_h,
180                                 int out_w,
181                                 int pool_h,
182                                 int pool_w,
183                                 int8_t *pool_out) {
184     int8_t *requant = (int8_t *)malloc((size_t)out_h * (size_t)out_w * (size_t)out_c);
185     int h, w, oc, ic, kh, kw;
186     if (!requant) {
187         fprintf(stderr, "alloc failed for requant buffer\n");

```

```

188     exit(1);
189 }
190
191 for (h = 0; h < out_h; ++h) {
192     for (w = 0; w < out_w; ++w) {
193         for (oc = 0; oc < out_c; ++oc) {
194             int32_t mac = 0;
195             for (ic = 0; ic < in_c; ++ic) {
196                 for (kh = 0; kh < 3; ++kh) {
197                     for (kw = 0; kw < 3; ++kw) {
198                         int32_t x = input[idx_hwc(h + kh, w + kw, ic, in_w, in_c)];
199                         int32_t wt = weights[idx_oihw(oc, ic, kh, kw, in_c, 3)];
200                         mac += x * wt;
201                     }
202                 }
203             }
204             requant[idx_hwc(h, w, oc, out_w, out_c)] =
205                 clamp_int8(requantize_single(mac, bias[oc], mult[oc], shift[oc], -128));
206         }
207     }
208 }
209
210 for (h = 0; h < pool_h; ++h) {
211     for (w = 0; w < pool_w; ++w) {
212         for (oc = 0; oc < out_c; ++oc) {
213             int base_h = h * 2;
214             int base_w = w * 2;
215             int8_t a = requant[idx_hwc(base_h + 0, base_w + 0, oc, out_w, out_c)];
216             int8_t b = requant[idx_hwc(base_h + 0, base_w + 1, oc, out_w, out_c)];
217             int8_t c = requant[idx_hwc(base_h + 1, base_w + 0, oc, out_w, out_c)];
218             int8_t d = requant[idx_hwc(base_h + 1, base_w + 1, oc, out_w, out_c)];
219             int8_t m1 = a > b ? a : b;
220             int8_t m2 = c > d ? c : d;
221             pool_out[idx_hwc(h, w, oc, pool_w, out_c)] = m1 > m2 ? m1 : m2;
222         }
223     }
224 }
225
226 free(requant);
227 }
228
229 static int fc_argmax(const int8_t *input, const int8_t *weights, const int32_t *bias) {
230     int oc, k;
231     int best_idx = 0;
232     int32_t best_val = 0;
233
234     for (oc = 0; oc < FC_OUT; ++oc) {
235         int64_t acc = bias[oc];
236         for (k = 0; k < FC_IN; ++k) {
237             acc += (int64_t)input[k] * (int64_t)weights[oc * FC_IN + k];
238         }
239         if (oc == 0 || acc > best_val) {
240             best_val = (int32_t)acc;
241             best_idx = oc;

```

```

242     }
243 }
244
245 return best_idx;
246 }
247
248 int main(int argc, char **argv) {
249     char path[1024];
250     const char *image_case_root;
251     const char *reference_case_root;
252     int8_t input[IMG_H * IMG_W];
253     int8_t pool1[P1_H * P1_W * C1_OUT];
254     int8_t pool2[P2_H * P2_W * C2_OUT];
255     int8_t pool3[P3_H * P3_W * C3_OUT];
256     int8_t conv1_w[C1_OUT * 3 * 3];
257     int8_t conv2_w[C2_OUT * C1_OUT * 3 * 3];
258     int8_t conv3_w[C3_OUT * C2_OUT * 3 * 3];
259     int8_t fc_w[FC_OUT * FC_IN];
260     int32_t conv1_bias[C1_OUT], conv2_bias[C2_OUT], conv3_bias[C3_OUT], fc_bias[FC_OUT];
261     int32_t conv1_m[C1_OUT], conv2_m[C2_OUT], conv3_m[C3_OUT];
262     int32_t conv1_sh[C1_OUT], conv2_sh[C2_OUT], conv3_sh[C3_OUT];
263     int predicted_class;
264     double started_ms, ended_ms;
265
266     if (argc != 3) {
267         fprintf(stderr, "usage: %s <image_case_root> <reference_case_root>\n", argv[0]);
268         return 1;
269     }
270
271     image_case_root = argv[1];
272     reference_case_root = argv[2];
273
274     build_path(path, sizeof(path), image_case_root, "tb_conv1_in_i8_64x64x1.txt");
275     if (read_i8_lines(path, input, IMG_H * IMG_W) != 0)
276         return 1;
277
278     build_path(path, sizeof(path), reference_case_root, "tb_conv1_w_i8_3x3x4.txt");
279     if (read_i8_lines(path, conv1_w, C1_OUT * 3 * 3) != 0)
280         return 1;
281     build_path(path, sizeof(path), reference_case_root, "tb_conv2_w_i8_3x3x4x8.txt");
282     if (read_i8_lines(path, conv2_w, C2_OUT * C1_OUT * 3 * 3) != 0)
283         return 1;
284     build_path(path, sizeof(path), reference_case_root, "tb_conv3_w_i8_3x3x8x8.txt");
285     if (read_i8_lines(path, conv3_w, C3_OUT * C2_OUT * 3 * 3) != 0)
286         return 1;
287     build_path(path, sizeof(path), reference_case_root, "tb_fc_w_i8_10x288.txt");
288     if (read_i8_lines(path, fc_w, FC_OUT * FC_IN) != 0)
289         return 1;
290
291     build_path(path, sizeof(path), reference_case_root, "tb_conv1_quant_bias_eff_i32_4.txt");
292     if (read_i32_lines(path, conv1_bias, C1_OUT) != 0)
293         return 1;
294     build_path(path, sizeof(path), reference_case_root, "tb_conv2_quant_bias_eff_i32_8.txt");
295     if (read_i32_lines(path, conv2_bias, C2_OUT) != 0)

```

```

296     return 1;
297     build_path(path, sizeof(path), reference_case_root, "tb_conv3_quant_bias_eff_i32_8.txt");
298     if (read_i32_lines(path, conv3_bias, C3_OUT) != 0)
299         return 1;
300     build_path(path, sizeof(path), reference_case_root, "tb_fc_bias_eff_i32_10.txt");
301     if (read_i32_lines(path, fc_bias, FC_OUT) != 0)
302         return 1;
303
304     build_path(path, sizeof(path), reference_case_root, "tb_conv1_quant_M_i32_4.txt");
305     if (read_i32_lines(path, conv1_m, C1_OUT) != 0)
306         return 1;
307     build_path(path, sizeof(path), reference_case_root, "tb_conv2_quant_M_i32_8.txt");
308     if (read_i32_lines(path, conv2_m, C2_OUT) != 0)
309         return 1;
310     build_path(path, sizeof(path), reference_case_root, "tb_conv3_quant_M_i32_8.txt");
311     if (read_i32_lines(path, conv3_m, C3_OUT) != 0)
312         return 1;
313
314     build_path(path, sizeof(path), reference_case_root, "tb_conv1_quant_sh_i32_4.txt");
315     if (read_i32_lines(path, conv1_sh, C1_OUT) != 0)
316         return 1;
317     build_path(path, sizeof(path), reference_case_root, "tb_conv2_quant_sh_i32_8.txt");
318     if (read_i32_lines(path, conv2_sh, C2_OUT) != 0)
319         return 1;
320     build_path(path, sizeof(path), reference_case_root, "tb_conv3_quant_sh_i32_8.txt");
321     if (read_i32_lines(path, conv3_sh, C3_OUT) != 0)
322         return 1;
323
324     started_ms = monotonic_ms();
325     conv1_forward(input, conv1_w, conv1_bias, conv1_m, conv1_sh, pool1);
326     conv_generic_forward(pool1, P1_H, P1_W, C1_OUT, conv2_w, C2_OUT, conv2_bias, conv2_m, conv2_sh,
327         C2_H, C2_W, P2_H, P2_W, pool2);
328     conv_generic_forward(pool2, P2_H, P2_W, C2_OUT, conv3_w, C3_OUT, conv3_bias, conv3_m, conv3_sh,
329         C3_H, C3_W, P3_H, P3_W, pool3);
330     predicted_class = fc_argmax(pool3, fc_w, fc_bias);
331     ended_ms = monotonic_ms();
332
333     printf("predict_class=%d\n", predicted_class);
334     printf("board_cpu_infer_ms=%.3f\n", ended_ms - started_ms);
335     return 0;
336 }

```

## A.5 Golden Model And Export Code

### A.5.1 Golden Batch Entry Point

*Regenerates all hardware-aligned digit cases and SRAM preload bundles.*

```
code/Golden-Module/matlab/hardware_aligned/run_all.m
```

```
1 % RUN_ALL Multi-case golden TXT generator.
2 %
3 % Loads the model once, then for every PNG matching 'image_glob' runs:
4 %   export_case -> hardware_aligned/debug/txt_cases/<case>/
5 %
6 % By default this folder is code-only. It reads the model from
7 % <repo>/models/v1.int8.params.mat and images from <repo>/matlab/digit_*.png.
8
9 clc; clear; close all;
10
11 hw_dir = fileparts(mfilename('fullpath'));
12 matlab_dir = fileparts(hw_dir);
13 repo_root = fileparts(matlab_dir);
14 addpath(hw_dir);
15
16 % ----- USER CONFIG -----
17 params_candidates = {
18     fullfile(repo_root, 'models', 'v1.int8.params.mat')
19     fullfile(matlab_dir, 'main', 'v1.int8.params.mat')
20     fullfile(hw_dir, 'v1.int8.params.mat')
21 };
22
23 % Glob all 'digit_<n>*.png' images in matlab/. Adjust as you add cases.
24 image_glob = fullfile(matlab_dir, 'digit_*.png');
25 out_root = fullfile(hw_dir, 'debug', 'txt_cases');
26 sram_out_root = fullfile(hw_dir, 'debug', 'sram_preload');
27 % -----
28
29 params_path = "";
30 for k = 1:numel(params_candidates)
31     if exist(params_candidates{k}, 'file')
32         params_path = string(params_candidates{k});
33         break;
34     end
35 end
36
37 if ~exist(params_path, 'file')
38     error('run_all: params .mat not found. Tried repo models/, matlab/main/, and hardware_aligned/.');
39 end
40
41 P = load_params(params_path);
42 imgs = dir(image_glob);
43 if isempty(imgs)
44     error('run_all: no images match %s', image_glob);
45 end
46
47 fprintf('Loaded params from %s\n', params_path);
48 fprintf('Processing %d image(s) matching %s\n', numel(imgs), image_glob);
49
50 for i = 1:numel(imgs)
51     img_path = fullfile(imgs(i).folder, imgs(i).name);
52     [~, base, ~] = fileparts(imgs(i).name);
```

```

53     case_name = sanitize_case_name(base);
54
55     R = export_case(img_path, P, case_name, out_root);
56     gen_sram_preload(case_name, R, P, sram_out_root);
57 end
58
59 fprintf('run_all: done.\n');
60
61
62 function s = sanitize_case_name(s)
63 % Make a filesystem-safe case tag from the image basename.
64     s = lower(s);
65     s = regexp(s, '[^a-z0-9_]+', '_');
66     s = regexp(s, '^_|_+$', '');
67 end

```

## A.5.2 Golden Case Exporter

*Exports one source image into input tensors, intermediate tensors, final logits, and manifest data.*

```
code/Golden-Module/matlab/hardware_aligned/export_case.m
```

```

1 function R = export_case(image_path, P, case_name, out_root)
2 % EXPORT_CASE Run hardware-aligned forward + write all golden TXT for one image.
3 %
4 % R = export_case(image_path, P, case_name, out_root)
5 %
6 % image_path : path to source PNG (any size; converted to 64x64 grayscale)
7 % P          : params struct from load_params()
8 % case_name  : short tag for the case directory (e.g. 'digit_0')
9 % out_root   : root directory; case files go to <out_root>/<case_name>/
10 %
11 % Returns the hw_forward result struct R for further inspection.
12 %
13 % Files written to <out_root>/<case_name>/ follow GOLDEN_FORMAT.md naming
14 % 'tb_<layer>_<role>_<dtype>_<shape>.txt' plus a 'manifest.txt'.
15
16 if nargin < 4 || isempty(out_root)
17     matlab_dir = fileparts(mfilename('fullpath'));
18     out_root   = fullfile(matlab_dir, 'debug', 'txt_cases');
19 end
20
21 case_dir = fullfile(out_root, case_name);
22 if ~exist(case_dir, 'dir'), mkdir(case_dir); end
23
24 % -----
25 % Load + quantize image to INT8
26 % -----
27 img_uint8 = imread(image_path);
28 if ndims(img_uint8) == 3
29     img_uint8 = rgb2gray(img_uint8);

```

```

30 end
31 if size(img_uint8, 1) ~= 64 || size(img_uint8, 2) ~= 64
32     img_uint8 = imresize(img_uint8, [64, 64]);
33 end
34 img_float = single(img_uint8) / 255.0;
35 img_q     = int32(round(img_float / P.s_in)) + P.in_zp;
36 img_q     = max(min(img_q, int32(127)), int32(-128));
37 img_int8  = int8(img_q);
38
39 % -----
40 % Hardware-aligned forward
41 % -----
42 R = hw_forward(img_int8, P);
43
44 % -----
45 % Write all golden TXT files
46 % -----
47 cd_ = @(name) fullfile(case_dir, name);
48
49 % Conv1
50 dump_txt(cd_('tb_conv1_in_i8_64x64x1.txt'),      R.input_int8,      'i8', 'hwc');
51 dump_txt(cd_('tb_conv1_w_i8_3x3x4.txt'),        squeeze(R.conv1_w), 'i8', 'ohw'); % Cin=1
52 dump_txt(cd_('tb_conv1_out_i32_62x62x4.txt'),   R.conv1_mac,      'i32', 'hwc');
53 dump_txt(cd_('tb_conv1_quant_bias_eff_i32_4.txt'), R.conv1_eff_bias, 'i32', 'flat');
54 dump_txt(cd_('tb_conv1_quant_M_i32_4.txt'),     R.conv1_M,        'i32', 'flat');
55 dump_txt(cd_('tb_conv1_quant_sh_i32_4.txt'),    R.conv1_sh,       'i32', 'flat');
56 dump_txt(cd_('tb_conv1_requant_i8_62x62x4.txt'), R.conv1_requant,  'i8', 'hwc');
57 dump_txt(cd_('tb_conv1_pool_i8_31x31x4.txt'),   R.pool1,          'i8', 'hwc');
58
59 % Conv2
60 dump_txt(cd_('tb_conv2_in_i8_31x31x4.txt'),     R.pool1,          'i8', 'hwc');
61 dump_txt(cd_('tb_conv2_w_i8_3x3x4x8.txt'),     R.conv2_w,        'i8', 'oihw');
62 dump_txt(cd_('tb_conv2_out_i32_29x29x8.txt'),  R.conv2_mac,      'i32', 'hwc');
63 dump_txt(cd_('tb_conv2_quant_bias_eff_i32_8.txt'), R.conv2_eff_bias, 'i32', 'flat');
64 dump_txt(cd_('tb_conv2_quant_M_i32_8.txt'),    R.conv2_M,        'i32', 'flat');
65 dump_txt(cd_('tb_conv2_quant_sh_i32_8.txt'),   R.conv2_sh,       'i32', 'flat');
66 dump_txt(cd_('tb_conv2_requant_i8_29x29x8.txt'), R.conv2_requant,  'i8', 'hwc');
67 dump_txt(cd_('tb_conv2_pool_i8_14x14x8.txt'),  R.pool2,          'i8', 'hwc');
68
69 % Conv3
70 dump_txt(cd_('tb_conv3_in_i8_14x14x8.txt'),     R.pool2,          'i8', 'hwc');
71 dump_txt(cd_('tb_conv3_w_i8_3x3x8x8.txt'),     R.conv3_w,        'i8', 'oihw');
72 dump_txt(cd_('tb_conv3_out_i32_12x12x8.txt'),  R.conv3_mac,      'i32', 'hwc');
73 dump_txt(cd_('tb_conv3_quant_bias_eff_i32_8.txt'), R.conv3_eff_bias, 'i32', 'flat');
74 dump_txt(cd_('tb_conv3_quant_M_i32_8.txt'),    R.conv3_M,        'i32', 'flat');
75 dump_txt(cd_('tb_conv3_quant_sh_i32_8.txt'),   R.conv3_sh,       'i32', 'flat');
76 dump_txt(cd_('tb_conv3_requant_i8_12x12x8.txt'), R.conv3_requant,  'i8', 'hwc');
77 dump_txt(cd_('tb_conv3_pool_i8_6x6x8.txt'),    R.pool3,          'i8', 'hwc');
78
79 % FC (no quant -- raw INT32 accumulator with eff_bias, argmax-ready)
80 dump_txt(cd_('tb_fc_in_i8_288.txt'),           R.fc_in,          'i8', 'flat');
81 dump_txt(cd_('tb_fc_w_i8_10x288.txt'),        R.fc_w,           'i8', 'fc_oik');
82 dump_txt(cd_('tb_fc_w_interleaved_i8_288x10.txt'), R.fc_w,           'i8', 'fc_int_kxo');

```

```

83 dump_txt(cd_('tb_fc_bias_eff_i32_10.txt'),          R.fc_eff_bias,          'i32', 'flat');
84 dump_txt(cd_('tb_fc_out_i32_10.txt'),             R.fc_acc,               'i32', 'flat');
85
86 % -----
87 % Manifest (key=value, GOLDEN_FORMAT.md compliant)
88 % -----
89 manifest_path = fullfile(case_dir, 'manifest.txt');
90 fid = fopen(manifest_path, 'w');
91 if fid < 0, error('export_case:io', 'cannot open %s', manifest_path); end
92 cu = onCleanup(@() fclose(fid));
93
94 [~, img_name, img_ext] = fileparts(image_path);
95 fprintf(fid, 'case=%s\n',          case_name);
96 fprintf(fid, 'image_file=%s%s\n',  img_name, img_ext);
97
98 fprintf(fid, 'conv1_input=tb_conv1_in_i8_64x64x1.txt\n');
99 fprintf(fid, 'conv1_weight=tb_conv1_w_i8_3x3x4.txt\n');
100 fprintf(fid, 'conv1_output=tb_conv1_out_i32_62x62x4.txt\n');
101 fprintf(fid, 'conv1_requant=tb_conv1_requant_i8_62x62x4.txt\n');
102 fprintf(fid, 'conv1_pool=tb_conv1_pool_i8_31x31x4.txt\n');
103 fprintf(fid, 'conv2_input=tb_conv2_in_i8_31x31x4.txt\n');
104 fprintf(fid, 'conv2_weight=tb_conv2_w_i8_3x3x4x8.txt\n');
105 fprintf(fid, 'conv2_output=tb_conv2_out_i32_29x29x8.txt\n');
106 fprintf(fid, 'conv2_requant=tb_conv2_requant_i8_29x29x8.txt\n');
107 fprintf(fid, 'conv2_pool=tb_conv2_pool_i8_14x14x8.txt\n');
108 fprintf(fid, 'conv3_input=tb_conv3_in_i8_14x14x8.txt\n');
109 fprintf(fid, 'conv3_weight=tb_conv3_w_i8_3x3x8x8.txt\n');
110 fprintf(fid, 'conv3_output=tb_conv3_out_i32_12x12x8.txt\n');
111 fprintf(fid, 'conv3_requant=tb_conv3_requant_i8_12x12x8.txt\n');
112 fprintf(fid, 'conv3_pool=tb_conv3_pool_i8_6x6x8.txt\n');
113 fprintf(fid, 'fc_input=tb_fc_in_i8_288.txt\n');
114 fprintf(fid, 'fc_weight=tb_fc_w_i8_10x288.txt\n');
115 fprintf(fid, 'fc_weight_interleaved=tb_fc_w_interleaved_i8_288x10.txt\n');
116 fprintf(fid, 'fc_bias_eff=tb_fc_bias_eff_i32_10.txt\n');
117 fprintf(fid, 'fc_output=tb_fc_out_i32_10.txt\n');
118 fprintf(fid, 'predict_class=%d\n',          R.predict_class);
119
120 fprintf(fid, 'input_zero_point=%d\n',          R.zp.in);
121 fprintf(fid, 'conv1_output_zero_point=%d\n',  R.zp.conv1_out);
122 fprintf(fid, 'conv2_output_zero_point=%d\n',  R.zp.conv2_out);
123 fprintf(fid, 'conv3_output_zero_point=%d\n',  R.zp.conv3_out);
124 fprintf(fid, 'fc_input_zero_point=%d\n',      R.zp.fc_in);
125
126 fprintf(fid, ['conv_rule=layer input keeps full upstream tensor values; ' ...
127             'current-layer conv MAC uses x_zp=0,w_zp=0,bias=0; ' ...
128             'eff_bias = bias - x_zp*sum(W) added in Quant stage; ' ...
129             'no standalone ReLU (subsumed by Quant clamp when out_zp=-128); ' ...
130             'FC has no requant (argmax on raw INT32 acc)\n']);
131
132 fprintf(fid, 'conv1_input_shape=64x64x1\n');
133 fprintf(fid, 'conv2_input_shape=31x31x4\n');
134 fprintf(fid, 'conv3_input_shape=14x14x8\n');
135 fprintf(fid, 'conv1_output_shape=62x62x4\n');
136 fprintf(fid, 'conv2_output_shape=29x29x8\n');

```

```

137 fprintf(fid, 'conv3_output_shape=12x12x8\n');
138 fprintf(fid, 'fc_output_shape=10\n');
139
140 fprintf(fid, 'conv1_output_range=[%d,%d]\n', min(R.conv1_mac(:)), max(R.conv1_mac(:)));
141 fprintf(fid, 'conv2_output_range=[%d,%d]\n', min(R.conv2_mac(:)), max(R.conv2_mac(:)));
142 fprintf(fid, 'conv3_output_range=[%d,%d]\n', min(R.conv3_mac(:)), max(R.conv3_mac(:)));
143 fprintf(fid, 'fc_acc_range=[%d,%d]\n', min(R.fc_acc), max(R.fc_acc));
144
145 fprintf('export_case: wrote case=%s -> %s (predict_class=%d)\n', ...
146         case_name, case_dir, R.predict_class);
147 end

```

### A.5.3 Hardware-Aligned Forward Pass

*Runs the integer CNN reference using RTL-compatible layer ordering and quantization behavior.*

code/Golden-Module/matlab/hardware\_aligned/hw\_forward.m

```

1 function R = hw_forward(img_int8, P)
2 % HW_FORWARD Hardware-aligned forward pass for one 64x64x1 image.
3 %
4 % R = hw_forward(img_int8, P)
5 %
6 % img_int8 : int8 [64, 64] or [64, 64, 1] -- already quantized input
7 % P       : params struct (see load_params helper below)
8 %
9 % Returns struct R with all intermediates needed by the testbench.
10 %
11 % -----
12 % Hardware conventions (vs the TFLite reference rps_conv2.m):
13 %
14 % 1. **Conv MAC stage**: x_zp = 0, w_zp = 0, bias = 0.
15 %    Output is raw INT32 MAC (sum(x_q * w_q)) with no zp folding.
16 %    These goldens go to tb_convN_out_i32_*.txt.
17 %
18 % 2. **Quant stage**: receives eff_bias (NOT raw bias):
19 %    eff_bias = b - x_zp * sum(W_per_oc)
20 %    Then: y = clamp(((MAC + eff_bias) * M >> sh) + zp_out, -128, 127)
21 %    These eff_bias / M / shift goldens go to tb_convN_quant_*_C.txt.
22 %
23 % 3. **No standalone ReLU**: With out_zp = -128 the Quant clamp itself
24 %    is the ReLU; we do NOT call relu() before pooling.
25 %
26 % 4. **Pool**: pure 2x2 max, no ReLU fold (use maxpool2x2_int8.m).
27 %
28 % 5. **FC**: NO requantization. Raw INT32 accumulator with eff_bias is
29 %    what the hardware feeds straight into argmax. Golden output is
30 %    tb_fc_out_i32_<OUT_CHANNELS>.txt; predicted class = argmax(fc_acc).
31 % -----
32 %
33 % Returned fields (every tensor matches the format the TB will read):

```

```

34 % R.input_int8 int8 [64,64,1]
35 % R.conv1_w / R.conv2_w / R.conv3_w int8 weight tensors (OIHW logical)
36 % R.conv1_mac / R.conv2_mac / R.conv3_mac
37 % int32 raw MAC outputs (no bias, no zp)
38 % R.conv1_eff_bias / .._M / .._sh int32 quant params per channel
39 % R.conv1_requant / R.conv2_requant / R.conv3_requant int8 quantized
40 % R.pool1 / R.pool2 / R.pool3 int8 maxpool outputs
41 % R.fc_in int8 [288] (= flatten of pool3)
42 % R.fc_w int8 [10, 288]
43 % R.fc_eff_bias int32 [10]
44 % R.fc_acc int32 [10] <-- HW final, argmax-ready
45 % R.predict_class 0..9
46 % R.zp struct of zero points used
47
48 if ndims(img_int8) == 2
49     img_int8 = reshape(img_int8, size(img_int8,1), size(img_int8,2), 1);
50 end
51 [imgH, imgW, imgC] = size3(img_int8);
52 assert(isa(img_int8, 'int8') && imgH == 64 && imgW == 64 && imgC == 1, ...
53     'hw_forward: img must be int8 [64,64] or [64,64,1]');
54
55 R.input_int8 = img_int8;
56 R.conv1_w = P.W1; % [Cout, Cin, H, W] logical (OIHW)
57 R.conv2_w = P.W2;
58 R.conv3_w = P.W3;
59 R.fc_w = P.W_fc; % [Cout, K]
60
61 R.zp = struct( ...
62     'in', P.in_zp, ...
63     'conv1_out', P.z_conv1_out, ...
64     'conv2_out', P.z_conv2_out, ...
65     'conv3_out', P.z_conv3_out, ...
66     'fc_in', P.z_fc_in);
67
68 % --- Layer 1 -----
69 % R.conv*_sh is the HARDWARE-facing right-shift amount consumed by
70 % Quantization_PE.v ('correct >>> sh'), so we pre-convert from the
71 % TFLite exponent 'shift' to '31 - shift' = total right-shift count.
72 R.conv1_mac = conv2d_mac(img_int8, P.W1); % int32 [62,62,4]
73 R.conv1_eff_bias = compute_eff_bias(P.b1, P.W1, P.in_zp); % int32 [4]
74 R.conv1_M = int32(P.qm1);
75 R.conv1_sh = int32(int32(31) - int32(P.shift1));
76 R.conv1_requant = quant_int32_to_int8(R.conv1_mac, R.conv1_eff_bias, ...
77     R.conv1_M, R.conv1_sh, P.z_conv1_out);
78 R.pool1 = maxpool2x2_int8(R.conv1_requant); % int8 [31,31,4]
79
80 % --- Layer 2 -----
81 R.conv2_mac = conv2d_mac(R.pool1, P.W2); % int32 [29,29,8]
82 R.conv2_eff_bias = compute_eff_bias(P.b2, P.W2, P.z_conv1_out); % int32 [8]
83 R.conv2_M = int32(P.qm2);
84 R.conv2_sh = int32(int32(31) - int32(P.shift2));
85 R.conv2_requant = quant_int32_to_int8(R.conv2_mac, R.conv2_eff_bias, ...
86     R.conv2_M, R.conv2_sh, P.z_conv2_out);
87 R.pool2 = maxpool2x2_int8(R.conv2_requant); % int8 [14,14,8]

```

```

88
89 % --- Layer 3 -----
90 R.conv3_mac      = conv2d_mac(R.pool2, P.W3);           % int32 [12,12,8]
91 R.conv3_eff_bias = compute_eff_bias(P.b3, P.W3, P.z_conv2_out); % int32 [8]
92 R.conv3_M        = int32(P.qm3);
93 R.conv3_sh       = int32(int32(31) - int32(P.shift3));
94 R.conv3_requant  = quant_int32_to_int8(R.conv3_mac, R.conv3_eff_bias, ...
95                                     R.conv3_M, R.conv3_sh, P.z_conv3_out);
96 R.pool3          = maxpool2x2_int8(R.conv3_requant);   % int8 [6,6,8]
97
98 % --- Flatten + FC -----
99 R.fc_in = flatten_nhw_int8(R.pool3); % int8 [288]
100 K      = numel(R.fc_in);
101 Cout   = size(P.W_fc, 1);
102 assert(size(P.W_fc, 2) == K, 'FC weight K (%d) must match flatten (%d)', size(P.W_fc,2), K);
103 assert(numel(P.b_fc) == Cout, 'FC bias must have length %d', Cout);
104
105 R.fc_eff_bias = compute_fc_eff_bias(P.b_fc, P.W_fc, P.z_fc_in); % int32 [Cout]
106
107 % FC accumulator (hardware path): acc[oc] = eff_bias[oc] + sum(x * w[oc,:])
108 % where x is the raw INT8 stored in SRAM_B (zp already baked in).
109 x_i32 = int32(R.fc_in(:));
110 acc   = zeros(Cout, 1, 'int32');
111 for oc = 1:Cout
112     w_i32 = int32(P.W_fc(oc, :)).';
113     % Use int64 accumulation internally to avoid intermediate overflow,
114     % then saturate to int32. RTL accumulator is 32 bits; saturation here
115     % matches what the hardware would silently wrap, so flag any case
116     % where it actually exceeds 32-bit range (would indicate a model bug).
117     s64 = int64(R.fc_eff_bias(oc)) + sum(int64(x_i32) .* int64(w_i32));
118     if s64 > intmax('int32') || s64 < intmin('int32')
119         warning('hw_forward:fc_overflow', ...
120               'FC oc=%d accumulator %d out of int32; saturating', oc, s64);
121     end
122     acc(oc) = int32(max(min(s64, int64(intmax('int32'))), int64(intmin('int32'))));
123 end
124 R.fc_acc = acc;
125
126 % Argmax: strict-greater so ties go to the lower index (matches RTL ST_ARGMAX).
127 best_idx = int32(0); % 0-based class id (matches predict_class[3:0])
128 best_val = R.fc_acc(1);
129 for oc = 2:Cout
130     if R.fc_acc(oc) > best_val
131         best_val = R.fc_acc(oc);
132         best_idx = int32(oc - 1);
133     end
134 end
135 R.predict_class = best_idx;
136 end
137
138
139 % =====
140 % Helpers
141 % =====

```

```

142
143 function mac = conv2d_mac(x_int8, W_oihw)
144 % Raw MAC convolution: x_zp = 0, w_zp = 0, bias = 0; INT32 output.
145 % W_oihw: int8 [Cout, Cin, H, W] (OIHw).
146 % Stride 1, valid padding.
147
148 [Cout, Cin, kH, kW] = size_oihw(W_oihw);
149 [Hi, Wi, Ci] = size3(x_int8);
150 assert(Cin == Ci, 'conv2d_mac: weight Cin (%d) != input C (%d)', Cin, Ci);
151
152 Ho = Hi - kH + 1;
153 Wo = Wi - kW + 1;
154 mac = zeros(Ho, Wo, Cout, 'int32');
155
156 x32 = int32(x_int8);
157 W32 = int32(W_oihw);
158
159 for oc = 1:Cout
160     for r = 1:Ho
161         for c = 1:Wo
162             s = int64(0);
163             for ic = 1:Cin
164                 patch = x32(r:r+kH-1, c:c+kW-1, ic);
165                 ker = squeeze(W32(oc, ic, :, :)); % [kH, kW]
166                 s = s + sum(int64(patch(:)) .* int64(ker(:)));
167             end
168             if s > intmax('int32') || s < intmin('int32')
169                 warning('hw_forward:mac_overflow', ...
170                     'MAC overflow at oc=%d r=%d c=%d : %d', oc, r, c, s);
171             end
172             mac(r, c, oc) = int32(s);
173         end
174     end
175 end
176 end
177
178
179 function eb = compute_eff_bias(bias, W_oihw, x_zp)
180 % eff_bias[oc] = bias[oc] - x_zp * sum(W[oc, :, :])
181 % W_oihw: int8 [Cout, Cin, H, W]
182 Cout = size(W_oihw, 1);
183 eb = zeros(Cout, 1, 'int32');
184 xz = int64(x_zp);
185 for oc = 1:Cout
186     sw = sum(int64(W_oihw(oc, :, :, :)), 'all');
187     eb(oc) = int32(int64(bias(oc)) - xz * sw);
188 end
189 end
190
191
192 function eb = compute_fc_eff_bias(bias, W_fc, x_zp)
193 % eff_bias[oc] = bias[oc] - x_zp * sum(W_fc[oc, :])
194 % W_fc: int8 [Cout, K]
195 Cout = size(W_fc, 1);

```

```

196     eb = zeros(Cout, 1, 'int32');
197     xz = int64(x_zp);
198     for oc = 1:Cout
199         sw = sum(int64(W_fc(oc, :)));
200         eb(oc) = int32(int64(bias(oc)) - xz * sw);
201     end
202 end
203
204
205 function out = quant_int32_to_int8(mac_i32, eff_bias_i32, M_i32, sh_i32, zp_out)
206 % Mirror of RTL Quantization_PE.v requant:
207 % acc      = mac + eff_bias           // signed 32b
208 % mul      = acc * M                 // signed 64b
209 % correct  = mul + (1<<(sh-1)) - (mul<0?1:0)
210 % shifft   = correct >>> sh         // arithmetic right shift
211 % y        = shifft - 128           // zp_out = -128 baked in
212 % out      = clamp(y, -128, 127)
213 %
214 % sh_i32 is the HARDWARE-native shift (= 31 - TFLite_exp). Passed through
215 % verbatim. For typical conv layers this is a positive value ~39..41.
216 %
217 % mac_i32:   int32 [H, W, C]
218 % eff_bias_i32: int32 [C]
219 % M_i32:     int32 [C]
220 % sh_i32:    int32 [C] (RTL-style right-shift amount)
221 % zp_out:    int32 scalar
222 [H, W, C] = size3(mac_i32);
223 out = zeros(H, W, C, 'int8');
224 for c = 1:C
225     acc64 = int64(mac_i32(:, :, c)) + int64(eff_bias_i32(c));
226     prod  = acc64 * int64(M_i32(c));
227     sh    = int64(sh_i32(c));
228     if sh > 0
229         round_term = bitshift(int64(1), sh - 1);
230         prod_adj   = prod + round_term - int64(prod < 0);
231         scaled     = bitshift(prod_adj, -sh);
232     elseif sh == 0
233         scaled = prod;
234     else
235         scaled = bitshift(prod, -sh);
236     end
237     scaled = scaled + int64(zp_out);
238     scaled(scaled > 127) = 127;
239     scaled(scaled < -128) = -128;
240     out(:, :, c) = int8(scaled);
241 end
242 end
243
244
245 function [Cout, Cin, kH, kW] = size_oihw(W)
246 sz = size(W);
247 Cout = sz(1); Cin = sz(2); kH = sz(3); kW = sz(4);
248 end
249

```

```

250
251 function [H, W, C] = size3(x)
252     if ismatrix(x)
253         x = reshape(x, size(x,1), size(x,2), 1);
254     end
255     sz = size(x);
256     H = sz(1); W = sz(2);
257     if numel(sz) >= 3, C = sz(3); else, C = 1; end
258 end

```

#### A.5.4 SRAM Preload Packer

*Packs convolution configuration, convolution weights, FC weights, and FC bias for HPS model replay.*

```
code/Golden-Module/matlab/hardware_aligned/gen_sram_preload.m
```

```

1 function gen_sram_preload(case_name, R, P, out_root)
2 % GEN_SRAM_PRELOAD Produce the four host preload streams the TB drives in.
3 %
4 % gen_sram_preload(case_name, R, P, out_root)
5 %
6 % case_name : tag (e.g. 'digit_0')
7 % R         : forward result struct from hw_forward / export_case
8 % P         : params struct from load_params
9 % out_root  : root output dir (default: matlab/main/debug/sram_preload/)
10 %
11 % Writes to <out_root>/<case_name>/:
12 %   preload_conv_cfg_<NCFG>w.txt   conv L1/L2/L3 cfg (45 host words to SRAM_A@0x000)
13 %   preload_conv_wt_225w_bytes.txt conv L1/L2/L3 weights (900 INT8 bytes = 225 packed words)
14 %   preload_fc_bias_10w.txt        FC bias eff (10 host words to SRAM_A@0x111 via LAYER_FC+SEL_CFG)
15 %   preload_fcw_864w.txt           FC weight host stream (864 words -> packer -> 288 x 80b)
16 %
17 % Each host word is one INT32 line (signed decimal) -- the TB feeds these
18 % directly into load_data[31:0] one per beat. Host words are signed but the
19 % bit-pattern is what the wrapper writes to SRAM.
20 %
21 % -----
22 % STATUS NOTE (read before relying on this script):
23 %
24 % * FC bias (10w) : FULLY IMPLEMENTED
25 % * FC weight (864w) : FULLY IMPLEMENTED -- matches RTL fcw_preload_packer
26 %                   spec (3 host words -> 80-bit slot, ch0 in LSB)
27 % * Conv weight (900B) : IMPLEMENTED via the reorder ported from
28 %                   matlab_old/gen_sram_preload.py. TB packs this
29 %                   byte stream into 225 host words.
30 % * Conv cfg (45w) : SCAFFOLD ONLY. Each layer-pass uses 9 cfg words,
31 %                   but the bit-packing inside each word depends on
32 %                   how Quantization_Top.v unpacks them. See the TODO
33 %                   inside conv_cfg_stream() below; consult the per-

```

```

34 %                                     layer cfg writer in the standalone TBs to verify.
35 % -----
36
37 if nargin < 4 || isempty(out_root)
38     matlab_dir = fileparts(mfilename('fullpath'));
39     out_root   = fullfile(matlab_dir, 'debug', 'sram_preload');
40 end
41 case_dir = fullfile(out_root, case_name);
42 if ~exist(case_dir, 'dir'), mkdir(case_dir); end
43
44 % -----
45 % conv CFG: 45 words = 9*5 (L1, L2_p0, L2_p1, L3_p0, L3_p1)
46 % -----
47 cfg_stream = [conv_cfg_stream(R.conv1_eff_bias, R.conv1_M, R.conv1_sh, R.zp.conv1_out, 'L1');
48               conv_cfg_stream(R.conv2_eff_bias(1:4), R.conv2_M(1:4), R.conv2_sh(1:4), R.zp.conv2_out,
49                               'L2p0');
50               conv_cfg_stream(R.conv2_eff_bias(5:8), R.conv2_M(5:8), R.conv2_sh(5:8), R.zp.conv2_out,
51                               'L2p1');
52               conv_cfg_stream(R.conv3_eff_bias(1:4), R.conv3_M(1:4), R.conv3_sh(1:4), R.zp.conv3_out,
53                               'L3p0');
54               conv_cfg_stream(R.conv3_eff_bias(5:8), R.conv3_M(5:8), R.conv3_sh(5:8), R.zp.conv3_out,
55                               'L3p1')];
56 assert(numel(cfg_stream) == 45, 'conv_cfg stream length %d != 45', numel(cfg_stream));
57 dump_txt(fullfile(case_dir, 'preload_conv_cfg_45w.txt'), cfg_stream, 'i32', 'flat');
58
59 % -----
60 % conv WT: 900 bytes = 4*(9 + 36*2 + 72*2), later packed to 225 host words
61 % -----
62 wt_stream = [conv_wt_pass(P.W1, 9, 1, 1, 4); % L1: dot_k=9, Cin=1, 4 channels (single pass)
63              conv_wt_pass(P.W2, 36, 4, 1, 4); % L2 pass 0: dot_k=36, Cin=4, ch 1..4
64              conv_wt_pass(P.W2, 36, 4, 5, 4); % L2 pass 1: ch 5..8
65              conv_wt_pass(P.W3, 72, 8, 1, 4); % L3 pass 0: dot_k=72, Cin=8, ch 1..4
66              conv_wt_pass(P.W3, 72, 8, 5, 4)]; % L3 pass 1: ch 5..8
67 assert(numel(wt_stream) == 4 * (9 + 36*2 + 72*2), 'conv wt stream length %d != 900 bytes', numel(
68 wt_stream));
69 % wt_stream is byte-stream (each entry one INT8 byte); host word packs 4 bytes,
70 % so the file ends up as 225 host beats only after 4-byte packing in the TB.
71 % Until the TB's packer is hooked here, dump byte-by-byte for transparency.
72 dump_txt(fullfile(case_dir, 'preload_conv_wt_225w_bytes.txt'), wt_stream, 'i8', 'flat');
73
74 % -----
75 % FC bias: 10 INT32 words. Goes to SRAM_A @ 0x111 via LAYER_FC+SEL_CFG.
76 % -----
77 dump_txt(fullfile(case_dir, 'preload_fc_bias_10w.txt'), R.fc_eff_bias, 'i32', 'flat');
78
79 % -----
80 % FC weight: 864 host words for fcw_preload_packer (3 host words per
81 % 80-bit slot). For each k in 0..287, emit 3 host words:
82 % word0 [31:0] = {k3, k2, k1, k0}      (LSB byte = ch 0)
83 % word1 [31:0] = {k7, k6, k5, k4}
84 % word2 [15:0] = {k9, k8}             (high 16b zero-padded)
85 % This matches input/RTL/SRAM/fcw_preload_packer.v.
86 % -----
87 Cout = size(P.W_fc, 1);

```

```

83     K = size(P.W_fc, 2);
84     assert(Cout == 10 && K == 288, 'FC weight expected [10,288], got [%d,%d]', Cout, K);
85     fcw_stream = zeros(K * 3, 1, 'int32');
86     for k = 1:K
87         % bytes for 10 channels at this k position
88         bytes = int32(P.W_fc(:, k)); % [10, 1]
89         bytes = bitand(bytes, int32(255)); % treat as unsigned for bit-packing
90         word0 = bitor(bitor(bitor(bytes(1), bitshift(bytes(2), 8)), ...
91             bitshift(bytes(3), 16)), bitshift(bytes(4), 24));
92         word1 = bitor(bitor(bitor(bytes(5), bitshift(bytes(6), 8)), ...
93             bitshift(bytes(7), 16)), bitshift(bytes(8), 24));
94         word2 = bitor(bytes(9), bitshift(bytes(10), 8)); % upper 16b = 0
95         % word0/1/2 are already int32 with the correct bit pattern.
96         % Do NOT pass through uint32() - MATLAB saturates negatives to 0.
97         fcw_stream(3*(k-1)+1) = word0;
98         fcw_stream(3*(k-1)+2) = word1;
99         fcw_stream(3*(k-1)+3) = word2;
100     end
101     assert(numel(fcw_stream) == 864, 'fcw stream length %d != 864', numel(fcw_stream));
102     dump_txt(fullfile(case_dir, 'preload_fcw_864w.txt'), fcw_stream, 'i32', 'flat');
103
104     fprintf('gen_sram_preload: wrote 4 streams to %s\n', case_dir);
105 end
106
107
108 % =====
109 % Helpers
110 % =====
111
112 function out = conv_cfg_stream(eff_bias, M, sh, zp_out, tag) %#ok<INUSD>
113 % Pack 9 host words for one conv pass cfg, matching input/RTL/quant_pool/
114 % integration/quant_param_loader.v (lines 6-8 header comment):
115 %
116 % word 0..3 : bias0..3 (per-channel eff_bias, one INT32 per word)
117 % word 4..7 : MO..3 (per-channel TFLite multiplier, one INT32 per word)
118 % word 8 : sh_packed = {sh1[31:24], sh2[23:16], sh3[15:8], sh4[7:0]}
119 %             i.e. ch1 is MSB, ch4 is LSB (Verilog concat order).
120 %
121 % Verified against Quantization_Top.v:32-53 which slices sh_in[31:24]->PE1,
122 % [23:16]->PE2, [15:8]->PE3, [7:0]->PE4.
123     assert(numel(eff_bias) == 4 && numel(M) == 4 && numel(sh) == 4, ...
124         'conv_cfg_stream: expected 4 channels per pass, got [%d,%d,%d] (%s)', ...
125         numel(eff_bias), numel(M), numel(sh), tag);
126
127     sh8 = bitand(int32(sh(:)), int32(255)); % truncate each shift to 8 bits
128     sh_packed = bitor(bitor(bitor( ...
129         bitshift(sh8(1), 24), ...
130         bitshift(sh8(2), 16)), ...
131         bitshift(sh8(3), 8)), ...
132         sh8(4));
133
134     out = int32(zeros(9, 1));
135     out(1:4) = int32(eff_bias(:));
136     out(5:8) = int32(M(:));

```

```

137     out(9)    = int32(sh_packed);
138 end
139
140
141 function bytes = conv_wt_pass(W_oihw, dot_k, Cin, ch_start, cols)
142 % Reorder + interleave one pass worth of weights for the systolic array.
143 % Ported from matlab_old/gen_sram_preload.py (reorder_weights + interleave_for_sram).
144 %
145 % W_oihw : int8 [Cout_total, Cin, 3, 3]
146 % dot_k  : 9 (L1) / 36 (L2) / 72 (L3) -- inner MAC sequence length
147 % Cin    : 1 / 4 / 8
148 % ch_start, cols : take ch_start..ch_start+cols-1 output channels for this pass
149 %
150 % Returns INT8 column vector of length cols*dot_k, byte-interleaved as the
151 % SRAM stores them: for each k position, emit col0_k, col1_k, ..., col(cols-1)_k.
152
153     assert(size(W_oihw, 2) == Cin, 'conv_wt_pass: W Cin %d != expected %d', size(W_oihw,2), Cin);
154
155     % Build w_tap[ch][rd_col_idx] in the SA traversal order.
156     w_tap = zeros(cols, dot_k, 'int8');
157     for ch = 0:cols-1
158         oc = ch_start + ch;          % 1-based oc index
159         for rd_col_idx = 0:dot_k-1
160             c_in_i = mod(rd_col_idx, Cin);
161             kw     = mod(floor(rd_col_idx / Cin), 3);
162             kh     = 2 - floor(rd_col_idx / (3 * Cin));
163             % MATLAB indices are 1-based and stored as OIHW.
164             w_tap(ch+1, rd_col_idx+1) = W_oihw(oc, c_in_i+1, kh+1, kw+1);
165         end
166     end
167
168     % Interleave: for each k, emit cols bytes (col0..col_{cols-1}).
169     bytes = zeros(cols * dot_k, 1, 'int8');
170     idx = 1;
171     for k = 1:dot_k
172         for ch = 1:cols
173             bytes(idx) = w_tap(ch, k);
174             idx = idx + 1;
175         end
176     end
177 end

```

### A.5.5 Golden Parameter Loader

*Loads the exported quantized parameters consumed by the MATLAB golden flow.*

```
code/Golden-Module/matlab/hardware_aligned/load_params.m
```

```

1 function P = load_params(params_mat_path)
2 % LOAD_PARAMS Load TFLite-exported INT8 model into hw_forward-ready struct.
3 %

```

```

4 % P = load_params(params_mat_path)
5 %
6 % Reads .mat produced by pytorch/export_tflite_params_mat.py with the
7 % conventions documented in matlab/main/rps_conv2.m:
8 %
9 %     values{11}/{10} : conv1 weight / bias (weight TFLite OHWI -> we permute to OIHW)
10 %     values{9}/{8}   : conv2 weight / bias
11 %     values{7}/{6}   : conv3 weight / bias
12 %     values{5}/{4}   : fc weight / bias (already [Cout, K])
13 %
14 % Output struct P:
15 %     P.W1, P.W2, P.W3   int8 weights in OIHW [Cout, Cin, H, W]
16 %     P.b1, P.b2, P.b3   int32 raw biases
17 %     P.W_fc             int8 [Cout=10, K=288]
18 %     P.b_fc             int32 [10]
19 %     P.in_zp            int32 input zero point
20 %     P.z_conv1_out / P.z_conv2_out / P.z_conv3_out int32 quant zps
21 %     P.z_fc_in          int32
22 %     P.qm1 / P.shift1   int32 (per-channel from TFLite multiplier)
23 %     P.qm2 / P.shift2
24 %     P.qm3 / P.shift3
25 %     P.s_in             single (input scale, used for image quant)
26 %
27 % FC quantization (qm_fc / shift_fc / z_fc_out) is intentionally OMITTED:
28 % the hardware FC path has no requantization stage, argmax runs on the raw
29 % INT32 accumulator. See matlab/main/hw_forward.m header for the rule.
30
31 if ~exist(params_mat_path, 'file')
32     error('load_params:missing', 'params .mat not found: %s', params_mat_path);
33 end
34 S = load(params_mat_path);
35
36 % --- Conv weights: TFLite OHWI -> OIHW [Cout, Cin, H, W] -----
37 P.W1 = int8(permute(S.values{11}, [1, 4, 2, 3])); % OHWI -> OIHW
38 P.W2 = int8(permute(S.values{9}, [1, 4, 2, 3]));
39 P.W3 = int8(permute(S.values{7}, [1, 4, 2, 3]));
40
41 P.b1 = int32(S.values{10}(:));
42 P.b2 = int32(S.values{8}(:));
43 P.b3 = int32(S.values{6}(:));
44
45 % --- FC weight / bias -----
46 P.W_fc = int8(S.values{5}); % [Cout, K]
47 P.b_fc = int32(S.values{4}(:));
48
49 % --- Zero points / scales -----
50 P.in_zp = int32(S.input_zero_points{1}(1));
51 P.s_in = single(S.input_scales{1}(1));
52 P.z_conv1_out = int32(S.activation_zero_points{1}(1));
53 P.z_conv2_out = int32(S.activation_zero_points{3}(1));
54 P.z_conv3_out = int32(S.activation_zero_points{5}(1));
55 P.z_fc_in = int32(S.activation_zero_points{10}(1));
56
57 s_conv1_out = S.activation_scales{1};

```

```

58     s_conv2_out = S.activation_scales{3};
59     s_conv3_out = S.activation_scales{5};
60
61     sw1 = S.scales{11};
62     sw2 = S.scales{9};
63     sw3 = S.scales{7};
64
65     % Per-channel multiplier/shift via TFLite formula.
66     [~, qm1, sh1] = tflite_quantize_multiplier(P.s_in,      sw1, s_conv1_out);
67     [~, qm2, sh2] = tflite_quantize_multiplier(s_conv1_out, sw2, s_conv2_out);
68     [~, qm3, sh3] = tflite_quantize_multiplier(s_conv2_out, sw3, s_conv3_out);
69
70     P.qm1 = int32(qm1(:));   P.shift1 = int32(sh1(:));
71     P.qm2 = int32(qm2(:));   P.shift2 = int32(sh2(:));
72     P.qm3 = int32(qm3(:));   P.shift3 = int32(sh3(:));
73 end

```

## A.5.6 Golden TXT Writer

*Writes integer tensors in the one-value-per-line format consumed by RTL testbenches and HPS tools.*

```
code/Golden-Module/matlab/hardware_aligned/dump_txt.m
```

```

1 function dump_txt(filepath, tensor, dtype, layout)
2 % DUMP_TXT Write a tensor as one decimal int per line (GOLDEN_FORMAT.md).
3 %
4 % dump_txt(filepath, tensor, dtype, layout)
5 %
6 % filepath : output path (will overwrite)
7 % tensor   : numeric array (any size)
8 % dtype    : 'i8' or 'i32' (validates element range)
9 % layout   : (optional) traversal hint:
10 %    'hwc'  - 3D activation [H, W, C], outer->inner = H, W, C (default for ndims=3)
11 %    'oihw' - 4D weight [Cout, Cin, H, W] (default for ndims=4)
12 %    'ohw'  - 3D weight [Cout, H, W] (e.g. conv1 with Cin=1, after squeeze)
13 %    'flat' - just iterate column-major (default for vectors / matrices)
14 %    'fc_oik'- 2D FC weight [Cout, K], outer=Cout, inner=K (default for 2D)
15 %    'fc_int_kxo' - SRAM-interleaved FC weight: outer=K, inner=Cout
16 %                (writes Cout bytes per K position, LSB=ch0)
17 %
18 % The tensor argument is expected to already be in the *logical* layout
19 % listed above (e.g. an OIHW tensor must be passed as a [Cout, Cin, H, W]
20 % numeric array). The function reorders elements into the file's
21 % slow->fast order before writing.
22 %
23 % GOLDEN_FORMAT.md compliance:
24 % * one signed decimal integer per line, LF terminator
25 % * no header, no comments, no blank lines
26 % * file ends with a trailing newline

```

```

27 % * line count == numel(tensor)
28
29 if nargin < 4 || isempty(layout)
30     if ndims(tensor) == 4
31         layout = 'oihw';
32     elseif ndims(tensor) == 3
33         layout = 'hwc';
34     elseif ismatrix(tensor) && ~isvector(tensor)
35         layout = 'fc_oik';
36     else
37         layout = 'flat';
38     end
39 end
40
41 switch lower(dtype)
42     case 'i8'
43         lo = -128; hi = 127;
44     case 'i32'
45         lo = -2147483648; hi = 2147483647;
46     otherwise
47         error('dump_txt:dtype', 'dtype must be ''i8'' or ''i32'', got "%s"', dtype);
48 end
49
50 % Reorder into the file emission sequence.
51 switch lower(layout)
52     case 'hwc'
53         assert(ndims(tensor) <= 3, 'hwc layout expects up to 3D tensor');
54         % MATLAB stores [H, W, C] in column-major (H fastest, then W, then C).
55         % File order is H outer, W middle, C fastest -> permute then linearize.
56         t = permute(tensor, [3, 2, 1]); % now [C, W, H], column-major == file order
57         seq = t(:);
58
59     case 'oihw'
60         assert(ndims(tensor) == 4, 'oihw layout expects 4D tensor');
61         % File order: Cout outer, Cin, H, W fastest.
62         t = permute(tensor, [4, 3, 2, 1]); % [W, H, Cin, Cout] in column-major
63         seq = t(:);
64
65     case 'ohw'
66         assert(ndims(tensor) == 3, 'ohw layout expects 3D tensor');
67         % File order: Cout outer, H, W fastest.
68         t = permute(tensor, [3, 2, 1]); % [W, H, Cout] in column-major
69         seq = t(:);
70
71     case 'fc_oik'
72         assert(ismatrix(tensor), 'fc_oik layout expects 2D [Cout, K] tensor');
73         % File order: Cout outer, K fastest.
74         t = tensor.'; % [K, Cout]; column-major == [Cout][K] with K fast
75         seq = t(:);
76
77     case 'fc_int_kxo'
78         assert(ismatrix(tensor), 'fc_int_kxo layout expects 2D [Cout, K] tensor');
79         % File order: K outer, Cout fastest. SRAM_FCW has one 80-bit slot
80         % per K position; bytes within slot are ch0 (LSB) .. ch9 (MSB).

```

```

81         % Writing column-major over the original [Cout, K] = ch0,ch1,...,ch9
82         % per K position is exactly what we want.
83         seq = tensor(:);
84
85         case 'flat'
86             seq = tensor(:);
87
88         otherwise
89             error('dump_txt:layout', 'unknown layout "%s"', layout);
90     end
91
92     % Range check (after cast to a safe wide type).
93     seq_i64 = int64(seq);
94     if any(seq_i64 < lo) || any(seq_i64 > hi)
95         bad_idx = find(seq_i64 < lo | seq_i64 > hi, 1, 'first');
96         error('dump_txt:range', ...
97             'value %d at element %d out of %s range [%d, %d] (file %s)', ...
98             seq_i64(bad_idx), bad_idx, dtype, lo, hi, filepath);
99     end
100
101     % Make sure parent directory exists.
102     parent = fileparts(filepath);
103     if ~isempty(parent) && ~exist(parent, 'dir')
104         mkdir(parent);
105     end
106
107     % Write: one decimal int per line, LF, trailing newline.
108     fid = fopen(filepath, 'w');
109     if fid < 0
110         error('dump_txt:io', 'cannot open %s for write', filepath);
111     end
112     cleanup = onCleanup(@() fclose(fid));
113     fprintf(fid, '%d\n', seq_i64);
114 end

```

## A.5.7 Flatten Helper

*Flattens NHWC feature maps into the FC input order expected by the RTL datapath.*

```
code/Golden-Module/matlab/hardware_aligned/flatten_nhwc_int8.m
```

```

1 function y = flatten_nhwc_int8(x)
2 %FLATTEN_NHWC_INT8 Flatten HxWxC int8 tensor with NHWC logical order.
3 % x: int8 [H, W, C]
4 % y: int8 [1, H*W*C]
5 %
6 % TFLite Flatten on [1,H,W,C] follows NHWC contiguous order:
7 % h major, then w, and c is the fastest-changing index.
8
9     assert(isa(x, 'int8'), 'x must be int8');
10    assert(ndims(x) == 3, 'x must be HxWxC');

```

```

11
12 [H, W, C] = size(x);
13 y = zeros(1, H * W * C, 'int8');
14
15 k = 1;
16 for h = 1:H
17     for w = 1:W
18         for c = 1:C
19             y(k) = x(h, w, c);
20             k = k + 1;
21         end
22     end
23 end
24 end

```

### A.5.8 Golden Max-Pool Helper

*Implements the 2-by-2 INT8 max-pooling operation used by the golden CNN path.*

code/Golden-Module/matlab/hardware\_aligned/maxpool2x2\_int8.m

```

1 function out = maxpool2x2_int8(x)
2 % MAXPOOL2X2_INT8 Pure 2x2 stride-2 max pooling (no ReLU fold).
3 %
4 % out = maxpool2x2_int8(x)
5 %
6 % x : int8 [H, W, C]
7 % out : int8 [floor(H/2), floor(W/2), C]
8 %
9 % Hardware alignment:
10 % * The accelerator's pool stage is plain max over a 2x2 window.
11 % * ReLU is implicit in the upstream Quantization stage when out_zp = -128
12 % (clamp(.., -128, 127) == max(.., zp_neg)). We do NOT clamp again here.
13 %
14 % Use this in the hardware-aligned forward pass instead of
15 % 'relu_maxpool2x2_int8' (which fuses an explicit ReLU into the window).
16
17 assert(isa(x, 'int8'), 'maxpool2x2_int8: input must be int8');
18 [H, W, C] = size(x);
19 Ho = floor(H / 2);
20 Wo = floor(W / 2);
21 out = zeros(Ho, Wo, C, 'int8');
22
23 for c = 1:C
24     for r = 1:Ho
25         for cc = 1:Wo
26             window = x(2*r-1:2*r, 2*cc-1:2*cc, c);
27             out(r, cc, c) = max(window(:));
28         end
29     end
30 end

```

31 end

## A.5.9 TFLite Multiplier Helper

*Computes TFLite-style fixed-point quantized multipliers and shifts.*

```
code/Golden-Module/matlab/hardware_aligned/tflite_quantize_multiplier.m
```

```
1 function [M, qm, shift] = tflite_quantize_multiplier(S_in, S_w, S_out)
2 %TFLITE_QUANTIZE_MULTIPLIER Compute per-channel real multiplier and
3 %TFLite-style (quantized_multiplier, shift) for int32->int8 requant.
4 %
5 % M[c] = (S_in * S_w[c]) / S_out[c]
6 % M ~= (qm / 2^31) * 2^shift
7 %
8 % Inputs:
9 % S_in : single/double scalar
10 % S_w : single/double scalar or vector (per-channel)
11 % S_out : single/double scalar or vector (per-channel)
12 %
13 % Outputs:
14 % M : double vector of real multipliers
15 % qm : int32 vector (Q0.31 quantized multiplier, non-negative)
16 % shift : int32 vector (can be negative/positive)
17 %
18 % Notes:
19 % - This matches the common TFLite meaning used in MultiplyByQuantizedMultiplier:
20 %     total_shift = 31 - shift
21 %     y = (x * qm + round) >> total_shift
22 % - For typical conv/FC requant, M is in (0,1), shift usually <= 0,
23 %   but we support general M >= 0.
24
25 % Make everything double for stable computation
26 S_in = double(S_in);
27 S_w = double(S_w);
28 S_out = double(S_out);
29
30 % Broadcast scalars to vector length
31 n = max([numel(S_w), numel(S_out), 1]);
32 if isscalar(S_w), S_w = repmat(S_w, 1, n); end
33 if isscalar(S_out), S_out = repmat(S_out, 1, n); end
34
35 % Real multipliers
36 M = (S_in .* S_w) ./ S_out;
37
38 qm = zeros(1, n, 'int32');
39 shift = zeros(1, n, 'int32');
40
41 for i = 1:n
42     r = M(i);
```

```

44     if r == 0 || ~isfinite(r) || r < 0
45         qm(i) = int32(0);
46         shift(i) = int32(0);
47         continue;
48     end
49
50     % log2 gives: r = significand * 2^exp, significand in [0.5, 1)
51     [significand, exp] = log2(r);
52
53     % Convert significand to Q31
54     q = round(significand * 2^31);
55
56     % Edge case: rounding to exactly 2^31
57     if q == 2^31
58         q = q / 2;
59         exp = exp + 1;
60     end
61
62     % Clamp into [0, 2^31-1] (TFLite expects non-negative here)
63     if q < 0, q = 0; end
64     if q > (2^31 - 1), q = (2^31 - 1); end
65
66     qm(i) = int32(q);
67     shift(i) = int32(exp);
68 end
69 end

```

### A.5.10 TFLite Parameter Exporter

*Extracts TFLite weights, biases, scales, zero-points, multipliers, and shifts into a MAT file.*

```
code/Golden-Module/pytorch/export_tflite_params_mat.py
```

```

1  #!/usr/bin/env python3
2  from __future__ import annotations
3
4  import argparse
5  from pathlib import Path
6
7  import numpy as np
8  import scipy.io as sio
9  import tensorflow as tf
10
11
12  def _to_cell(items: list[object]) -> np.ndarray:
13      arr = np.empty((len(items), 1), dtype=object)
14      for i, v in enumerate(items):
15          arr[i, 0] = v
16      return arr
17
18

```

```

19 def collect_param_tensor_indices(interpreter: tf.lite.Interpreter) -> list[int]:
20     details = interpreter.get_tensor_details()
21     all_indices = {d["index"] for d in details}
22     model_inputs = {d["index"] for d in interpreter.get_input_details()}
23
24     op_inputs: set[int] = set()
25     op_outputs: set[int] = set()
26     for op in interpreter._get_ops_details():
27         op_inputs.update(i for i in op["inputs"] if i >= 0)
28         op_outputs.update(i for i in op["outputs"] if i >= 0)
29
30     candidates = (op_inputs - op_outputs - model_inputs) & all_indices
31     return sorted(candidates)
32
33
34 def collect_graph_tensor_indices(interpreter: tf.lite.Interpreter) -> list[int]:
35     details = interpreter.get_tensor_details()
36     all_indices = {d["index"] for d in details}
37     model_inputs = {d["index"] for d in interpreter.get_input_details()}
38     model_outputs = {d["index"] for d in interpreter.get_output_details()}
39
40     op_inputs: set[int] = set()
41     op_outputs: set[int] = set()
42     for op in interpreter._get_ops_details():
43         op_inputs.update(i for i in op["inputs"] if i >= 0)
44         op_outputs.update(i for i in op["outputs"] if i >= 0)
45
46     graph_indices = (op_inputs | op_outputs | model_inputs | model_outputs) & all_indices
47     return sorted(graph_indices)
48
49
50 def collect_activation_tensor_indices(interpreter: tf.lite.Interpreter) -> list[int]:
51     details = interpreter.get_tensor_details()
52     all_indices = {d["index"] for d in details}
53
54     op_outputs: set[int] = set()
55     for op in interpreter._get_ops_details():
56         op_outputs.update(i for i in op["outputs"] if i >= 0)
57
58     return sorted(op_outputs & all_indices)
59
60
61 def export_params_to_mat(model_path: Path, out_path: Path) -> None:
62     interpreter = tf.lite.Interpreter(model_path=str(model_path))
63     interpreter.allocate_tensors()
64
65     tensor_details = {d["index"]: d for d in interpreter.get_tensor_details()}
66     op_details = interpreter._get_ops_details()
67     input_details = interpreter.get_input_details()
68     output_details = interpreter.get_output_details()
69
70     param_indices = collect_param_tensor_indices(interpreter)
71     graph_indices = collect_graph_tensor_indices(interpreter)
72     activation_indices = collect_activation_tensor_indices(interpreter)

```

```

73
74 # Parameter tensors (weights/bias/const tensors consumed by ops)
75 names: list[str] = []
76 tensor_indices: list[int] = []
77 dtypes: list[str] = []
78 shapes: list[np.ndarray] = []
79 shape_signatures: list[np.ndarray] = []
80 quantized_dims: list[int] = []
81 values: list[np.ndarray] = []
82 scales: list[np.ndarray] = []
83 zero_points: list[np.ndarray] = []
84
85 for idx in param_indices:
86     d = tensor_details[idx]
87     q = d["quantization_parameters"]
88     v = interpreter.get_tensor(idx)
89
90     names.append(d["name"])
91     tensor_indices.append(int(idx))
92     dtypes.append(str(v.dtype))
93     shapes.append(np.asarray(d["shape"], dtype=np.int64))
94     shape_signatures.append(np.asarray(d["shape_signature"], dtype=np.int64))
95     quantized_dims.append(int(q["quantized_dimension"]))
96     values.append(v)
97     scales.append(np.asarray(q["scales"], dtype=np.float32))
98     zero_points.append(np.asarray(q["zero_points"], dtype=np.int32))
99
100 # Activation tensors (op outputs)
101 act_names: list[str] = []
102 act_indices: list[int] = []
103 act_dtypes: list[str] = []
104 act_shapes: list[np.ndarray] = []
105 act_shape_signatures: list[np.ndarray] = []
106 act_quantized_dims: list[int] = []
107 act_scales: list[np.ndarray] = []
108 act_zero_points: list[np.ndarray] = []
109
110 for idx in activation_indices:
111     d = tensor_details[idx]
112     q = d["quantization_parameters"]
113
114     act_names.append(d["name"])
115     act_indices.append(int(idx))
116     act_dtypes.append(str(np.dtype(d["dtype"])))
117     act_shapes.append(np.asarray(d["shape"], dtype=np.int64))
118     act_shape_signatures.append(np.asarray(d["shape_signature"], dtype=np.int64))
119     act_quantized_dims.append(int(q["quantized_dimension"]))
120     act_scales.append(np.asarray(q["scales"], dtype=np.float32))
121     act_zero_points.append(np.asarray(q["zero_points"], dtype=np.int32))
122
123 # Model IO tensor specs
124 input_names: list[str] = []
125 input_indices: list[int] = []
126 input_dtypes: list[str] = []

```

```

127 input_shapes: list[np.ndarray] = []
128 input_shape_signatures: list[np.ndarray] = []
129 input_quantized_dims: list[int] = []
130 input_scales: list[np.ndarray] = []
131 input_zero_points: list[np.ndarray] = []
132
133 for d in input_details:
134     q = d["quantization_parameters"]
135     input_names.append(d["name"])
136     input_indices.append(int(d["index"]))
137     input_dtypes.append(str(np.dtype(d["dtype"])))
138     input_shapes.append(np.asarray(d["shape"], dtype=np.int64))
139     input_shape_signatures.append(np.asarray(d["shape_signature"], dtype=np.int64))
140     input_quantized_dims.append(int(q["quantized_dimension"]))
141     input_scales.append(np.asarray(q["scales"], dtype=np.float32))
142     input_zero_points.append(np.asarray(q["zero_points"], dtype=np.int32))
143
144 output_names: list[str] = []
145 output_indices: list[int] = []
146 output_dtypes: list[str] = []
147 output_shapes: list[np.ndarray] = []
148 output_shape_signatures: list[np.ndarray] = []
149 output_quantized_dims: list[int] = []
150 output_scales: list[np.ndarray] = []
151 output_zero_points: list[np.ndarray] = []
152
153 for d in output_details:
154     q = d["quantization_parameters"]
155     output_names.append(d["name"])
156     output_indices.append(int(d["index"]))
157     output_dtypes.append(str(np.dtype(d["dtype"])))
158     output_shapes.append(np.asarray(d["shape"], dtype=np.int64))
159     output_shape_signatures.append(np.asarray(d["shape_signature"], dtype=np.int64))
160     output_quantized_dims.append(int(q["quantized_dimension"]))
161     output_scales.append(np.asarray(q["scales"], dtype=np.float32))
162     output_zero_points.append(np.asarray(q["zero_points"], dtype=np.int32))
163
164 # Operator graph connectivity
165 op_indices: list[int] = []
166 op_names: list[str] = []
167 op_inputs: list[np.ndarray] = []
168 op_outputs: list[np.ndarray] = []
169
170 for op in op_details:
171     op_indices.append(int(op["index"]))
172     op_names.append(op["op_name"])
173     op_inputs.append(np.asarray(op["inputs"], dtype=np.int32))
174     op_outputs.append(np.asarray(op["outputs"], dtype=np.int32))
175
176 out_path.parent.mkdir(parents=True, exist_ok=True)
177 sio.savemat(
178     str(out_path),
179     {
180         "model_path": str(model_path),

```

```

181     "num_graph_tensors": np.int32(len(graph_indices)),
182     "graph_tensor_indices": np.asarray(graph_indices, dtype=np.int32).reshape(-1, 1),
183     "num_params": np.int32(len(param_indices)),
184     "names": _to_cell(names),
185     "tensor_indices": np.asarray(tensor_indices, dtype=np.int32).reshape(-1, 1),
186     "dtypes": _to_cell(dtypes),
187     "shapes": _to_cell(shapes),
188     "shape_signatures": _to_cell(shape_signatures),
189     "quantized_dimensions": np.asarray(quantized_dims, dtype=np.int32).reshape(-1, 1),
190     "values": _to_cell(values),
191     "scales": _to_cell(scales),
192     "zero_points": _to_cell(zero_points),
193     "num_inputs": np.int32(len(input_details)),
194     "input_names": _to_cell(input_names),
195     "input_indices": np.asarray(input_indices, dtype=np.int32).reshape(-1, 1),
196     "input_dtypes": _to_cell(input_dtypes),
197     "input_shapes": _to_cell(input_shapes),
198     "input_shape_signatures": _to_cell(input_shape_signatures),
199     "input_quantized_dimensions": np.asarray(input_quantized_dims, dtype=np.int32).reshape(-1, 1)
200 ,
201     "input_scales": _to_cell(input_scales),
202     "input_zero_points": _to_cell(input_zero_points),
203     "num_outputs": np.int32(len(output_details)),
204     "output_names": _to_cell(output_names),
205     "output_indices": np.asarray(output_indices, dtype=np.int32).reshape(-1, 1),
206     "output_dtypes": _to_cell(output_dtypes),
207     "output_shapes": _to_cell(output_shapes),
208     "output_shape_signatures": _to_cell(output_shape_signatures),
209     "output_quantized_dimensions": np.asarray(output_quantized_dims, dtype=np.int32).reshape(-1,
210 1),
211     "output_scales": _to_cell(output_scales),
212     "output_zero_points": _to_cell(output_zero_points),
213     "num_activations": np.int32(len(activation_indices)),
214     "activation_names": _to_cell(act_names),
215     "activation_indices": np.asarray(act_indices, dtype=np.int32).reshape(-1, 1),
216     "activation_dtypes": _to_cell(act_dtypes),
217     "activation_shapes": _to_cell(act_shapes),
218     "activation_shape_signatures": _to_cell(act_shape_signatures),
219     "activation_quantized_dimensions": np.asarray(act_quantized_dims, dtype=np.int32).reshape(-1,
220 1),
221     "activation_scales": _to_cell(act_scales),
222     "activation_zero_points": _to_cell(act_zero_points),
223     "num_ops": np.int32(len(op_details)),
224     "op_indices": np.asarray(op_indices, dtype=np.int32).reshape(-1, 1),
225     "op_names": _to_cell(op_names),
226     "op_inputs": _to_cell(op_inputs),
227     "op_outputs": _to_cell(op_outputs),
228 },
229 do_compression=True,
230 )
231
232 print(f"[OK] Exported {len(param_indices)} parameter tensors")
233 print(f"[OK] Exported {len(input_details)} input tensor specs")
234 print(f"[OK] Exported {len(output_details)} output tensor specs")

```

```

232     print(f"[OK] Exported {len(activation_indices)} activation tensor specs")
233     print(f"[OK] Exported {len(op_details)} operator nodes")
234     print(f"[OK] Saved: {out_path}")
235
236
237 def main() -> None:
238     parser = argparse.ArgumentParser(
239         description="Export int8 TFLite parameter tensors to a MATLAB .mat file."
240     )
241     parser.add_argument(
242         "--model",
243         type=Path,
244         default=Path("models/v2.int8.tflite"),
245         help="Path to .tflite model",
246     )
247     parser.add_argument(
248         "--out",
249         type=Path,
250         default=Path("models/v2.int8.params.mat"),
251         help="Output .mat path",
252     )
253     args = parser.parse_args()
254
255     if not args.model.exists():
256         raise FileNotFoundError(f"Model not found: {args.model}")
257
258     export_params_to_mat(args.model.resolve(), args.out.resolve())
259
260
261 if __name__ == "__main__":
262     main()

```

### A.5.11 Digit CNN Retraining Script

*Trains and exports the compact digit CNN model family used by the accelerator flow.*

```
code/Golden-Module/pytorch/retrain_digits_cnn.py
```

```

1  #!/usr/bin/env python3
2  from __future__ import annotations
3
4  import argparse
5  import math
6  import subprocess
7  import sys
8  import textwrap
9  from pathlib import Path
10
11 from tensorflow.keras import callbacks
12 from tensorflow.keras.layers import Conv2D, Dense, Flatten, InputLayer, MaxPooling2D
13 from tensorflow.keras.models import Sequential

```

```

14 from tensorflow.keras.preprocessing.image import ImageDataGenerator
15
16
17 def build_model(num_classes: int = 10) -> Sequential:
18     model = Sequential(
19         [
20             InputLayer(input_shape=(64, 64, 1)),
21             Conv2D(filters=4, kernel_size=(3, 3), strides=(1, 1), padding="valid", activation="relu"),
22             MaxPooling2D(pool_size=(2, 2)),
23             Conv2D(filters=8, kernel_size=(3, 3), strides=(1, 1), padding="valid", activation="relu"),
24             MaxPooling2D(pool_size=(2, 2)),
25             Conv2D(filters=8, kernel_size=(3, 3), strides=(1, 1), padding="valid", activation="relu"),
26             MaxPooling2D(pool_size=(2, 2)),
27             Flatten(),
28             Dense(units=num_classes, activation="softmax"),
29         ]
30     )
31     model.compile(loss="categorical_crossentropy", optimizer="adam", metrics=["accuracy"])
32     return model
33
34
35 def main() -> None:
36     parser = argparse.ArgumentParser(description="Retrain the 64x64 grayscale 0-9 sign-language CNN and
37     export INT8 TFLite.")
38     parser.add_argument("--dataset", type=Path, required=True, help="Merged dataset root with train/
39     validation/test splits.")
40     parser.add_argument("--models-dir", type=Path, default=Path("../models"), help="Model output
41     directory (default: ../models).")
42     parser.add_argument("--name", type=str, default="v3_merged", help="Base model name (default:
43     v3_merged).")
44     parser.add_argument("--epochs", type=int, default=60, help="Training epochs (default: 60).")
45     parser.add_argument("--batch-size", type=int, default=32, help="Batch size (default: 32).")
46     parser.add_argument("--rep-per-class", type=int, default=120, help="Representative samples per class
47     for INT8 export.")
48     args = parser.parse_args()
49
50     dataset_root = args.dataset.resolve()
51     train_folder = dataset_root / "train"
52     validation_folder = dataset_root / "validation"
53     test_folder = dataset_root / "test"
54     if not train_folder.exists():
55         raise FileNotFoundError(f"Missing train split: {train_folder}")
56     if not validation_folder.exists():
57         raise FileNotFoundError(f"Missing validation split: {validation_folder}")
58     if not test_folder.exists():
59         raise FileNotFoundError(f"Missing test split: {test_folder}")
60
61     train_gen_cfg = ImageDataGenerator(
62         rescale=1.0 / 255.0,
63         rotation_range=12,
64         width_shift_range=0.10,
65         height_shift_range=0.10,
66         shear_range=0.10,
67         zoom_range=0.10,

```

```

63     horizontal_flip=False,
64     vertical_flip=False,
65     fill_mode="nearest",
66 )
67 eval_gen_cfg = ImageDataGenerator(rescale=1.0 / 255.0)
68
69 train_gen = train_gen_cfg.flow_from_directory(
70     train_folder,
71     target_size=(64, 64),
72     color_mode="grayscale",
73     batch_size=args.batch_size,
74     class_mode="categorical",
75     shuffle=True,
76 )
77 validation_gen = eval_gen_cfg.flow_from_directory(
78     validation_folder,
79     target_size=(64, 64),
80     color_mode="grayscale",
81     batch_size=args.batch_size,
82     class_mode="categorical",
83     shuffle=False,
84 )
85 test_gen = eval_gen_cfg.flow_from_directory(
86     test_folder,
87     target_size=(64, 64),
88     color_mode="grayscale",
89     batch_size=args.batch_size,
90     class_mode="categorical",
91     shuffle=False,
92 )
93
94 expected_mapping = {str(i): i for i in range(10)}
95 if train_gen.class_indices != expected_mapping:
96     raise ValueError(f"Unexpected class mapping: {train_gen.class_indices}")
97
98 models_dir = args.models_dir.resolve()
99 models_dir.mkdir(parents=True, exist_ok=True)
100 h5_path = models_dir / f"{args.name}.h5"
101 tflite_path = models_dir / f"{args.name}.int8.tflite"
102
103 model = build_model(num_classes=10)
104 checkpoint = callbacks.ModelCheckpoint(str(h5_path), monitor="val_accuracy", save_best_only=True,
105     verbose=1)
106 early_stop = callbacks.EarlyStopping(monitor="val_accuracy", patience=10, restore_best_weights=True)
107
108 history = model.fit(
109     train_gen,
110     epochs=args.epochs,
111     steps_per_epoch=math.ceil(train_gen.samples / train_gen.batch_size),
112     validation_data=validation_gen,
113     validation_steps=math.ceil(validation_gen.samples / validation_gen.batch_size),
114     callbacks=[checkpoint, early_stop],
115     verbose=2,
116 )

```

```

116
117 test_loss, test_acc = model.evaluate(test_gen, verbose=0)
118 print(f"[TEST] loss={test_loss:.4f} acc={test_acc:.4f}")
119
120 export_script = textwrap.dedent(
121     """
122     import sys
123     from pathlib import Path
124
125     import numpy as np
126     import tensorflow as tf
127
128     model_path = Path(sys.argv[1])
129     rep_dir = Path(sys.argv[2])
130     out_path = Path(sys.argv[3])
131     per_class_samples = int(sys.argv[4])
132
133     model = tf.keras.models.load_model(str(model_path), compile=False)
134     model(np.zeros((1, 64, 64, 1), dtype=np.float32))
135
136     def representative_dataset_gen():
137         exts = {'.jpg', '.jpeg', '.png', '.bmp', '.webp'}
138         class_dirs = sorted([d for d in rep_dir.iterdir() if d.is_dir()], key=lambda p: int(p.name))
139         picked = []
140         for d in class_dirs:
141             files = [p for p in sorted(d.rglob('*')) if p.suffix.lower() in exts]
142             picked.extend(files[: min(per_class_samples, len(files))])
143         if not picked:
144             raise FileNotFoundError(f'No representative images found in {rep_dir}')
145         for p in picked:
146             img = tf.keras.utils.load_img(p, color_mode='grayscale', target_size=(64, 64))
147             arr = tf.keras.utils.img_to_array(img).astype(np.float32) / 255.0
148             arr = np.expand_dims(arr, axis=0)
149             yield [arr]
150
151     serving_fn = tf.function(lambda x: model(x))
152     concrete = serving_fn.get_concrete_function(tf.TensorSpec([1, 64, 64, 1], tf.float32))
153     converter = tf.lite.TFLiteConverter.from_concrete_functions([concrete], model)
154     converter.optimizations = [tf.lite.Optimize.DEFAULT]
155     converter.representative_dataset = representative_dataset_gen
156     converter.target_spec.supported_ops = [tf.lite.OpsSet.TFLITE_BUILTINS_INT8]
157     converter.inference_input_type = tf.int8
158     converter.inference_output_type = tf.int8
159     out_path.write_bytes(converter.convert())
160
161     print(f'[OK] wrote {out_path}')
162     """
163 )
164 cmd = [sys.executable, "-c", export_script, str(h5_path), str(train_folder), str(tflite_path), str(
    args.rep_per_class)]
165 subprocess.run(cmd, check=True)
166
167 summary = {
168     "name": args.name,

```

```

169     "dataset": str(dataset_root),
170     "train_samples": int(train_gen.samples),
171     "validation_samples": int(validation_gen.samples),
172     "test_samples": int(test_gen.samples),
173     "best_val_accuracy": max(history.history.get("val_accuracy", [0.0])),
174     "test_accuracy": float(test_acc),
175     "h5_path": str(h5_path),
176     "tflite_path": str(tflite_path),
177 }
178 print(summary)
179
180
181 if __name__ == "__main__":
182     main()

```

## A.6 Quartus And Platform Designer Setup

The Quartus and Platform Designer GUI was used only as an editing and inspection front end. The reproducible project state is captured by the saved source configuration files listed below: the .qsys Platform Designer system, the custom component \*\_hw.tcl, the Quartus project-generation Tcl script, and the timing constraint .sdc. Generated HDL, QIP files, reports, programming files, and build databases are intentionally excluded from the submission archive because they can be regenerated from these sources.

### A.6.1 Platform Designer Component Descriptor

*Registers the custom CNN MMIO interface as a Platform Designer component.*

```
code/platform_designer/cnn_mmio_interface_hw.tcl
```

```

1  # TCL File adapted for Component Editor / Platform Designer 21.1
2  #
3  # cnn_mmio_interface "cnn_mmio_interface" v1.0
4
5  package require -exact qsys 16.1
6
7  # -----
8  # module cnn_mmio_interface
9  # -----
10 set_module_property DESCRIPTION "Memory-mapped wrapper for the CNN_ACC gesture accelerator"
11 set_module_property NAME cnn_mmio_interface
12 set_module_property VERSION 1.0
13 set_module_property GROUP "Project"
14 set_module_property INTERNAL false
15 set_module_property HIDE_FROM_SOPC false
16 set_module_property HIDE_FROM_QUARTUS false
17 set_module_property OPAQUE_ADDRESS_MAP true

```

```

18 set_module_property AUTHOR ""
19 set_module_property DISPLAY_NAME cnn_mmio_interface
20 set_module_property INSTANTIATE_IN_SYSTEM_MODULE true
21 set_module_property EDITABLE true
22 set_module_property REPORT_TO_TALKBACK false
23 set_module_property ALLOW_GREYBOX_GENERATION false
24 set_module_property REPORT_HIERARCHY false
25
26 # -----
27 # file sets
28 # -----
29 add_fileset QUARTUS_SYNTH QUARTUS_SYNTH "" ""
30 set_fileset_property QUARTUS_SYNTH TOP_LEVEL cnn_mmio_interface
31 set_fileset_property QUARTUS_SYNTH ENABLE_RELATIVE_INCLUDE_PATHS false
32 set_fileset_property QUARTUS_SYNTH ENABLE_FILE_OVERWRITE_MODE false
33 add_fileset_file ../input/RTL/interface/cnn_mmio_interface.v VERILOG PATH ../input/RTL/interface/
   cnn_mmio_interface.v TOP_LEVEL_FILE
34 add_fileset_file ../input/RTL/system_top.v VERILOG PATH ../input/RTL/system_top.v
35 add_fileset_file ../input/RTL/fsm/top_fsm.v VERILOG PATH ../input/RTL/fsm/top_fsm.v
36 add_fileset_file ../input/RTL/fsm/layer_runner_fsm.v VERILOG PATH ../input/RTL/fsm/layer_runner_fsm.v
37 add_fileset_file ../input/RTL/fsm/conv_data_adapter.v VERILOG PATH ../input/RTL/fsm/conv_data_adapter.v
38 add_fileset_file ../input/RTL/fsm/wt_prepad_inserter.v VERILOG PATH ../input/RTL/fsm/wt_prepad_inserter.v
39 add_fileset_file ../input/RTL/conv_core/conv_quant_pool.v VERILOG PATH ../input/RTL/conv_core/
   conv_quant_pool.v
40 add_fileset_file ../input/RTL/conv_core/conv_top.v VERILOG PATH ../input/RTL/conv_core/conv_top.v
41 add_fileset_file ../input/RTL/conv_core/conv_engine_ctrl.v VERILOG PATH ../input/RTL/conv_core/
   conv_engine_ctrl.v
42 add_fileset_file ../input/RTL/conv_core/weight_buffer.v VERILOG PATH ../input/RTL/conv_core/
   weight_buffer.v
43 add_fileset_file ../input/RTL/conv_core/Line_Buffer.v VERILOG PATH ../input/RTL/conv_core/Line_Buffer.v
44 add_fileset_file ../input/RTL/conv_core/input_row_aligner.v VERILOG PATH ../input/RTL/conv_core/
   input_row_aligner.v
45 add_fileset_file ../input/RTL/conv_core/Conv_Buffer.v VERILOG PATH ../input/RTL/conv_core/Conv_Buffer.v
46 add_fileset_file ../input/RTL/conv_core/sa_skew_feeder.v VERILOG PATH ../input/RTL/conv_core/
   sa_skew_feeder.v
47 add_fileset_file ../input/RTL/conv_core/systolic_array_top.v VERILOG PATH ../input/RTL/conv_core/
   systolic_array_top.v
48 add_fileset_file ../input/RTL/quant_pool/integration/quant_param_loader.v VERILOG PATH ../input/RTL/
   quant_pool/integration/quant_param_loader.v
49 add_fileset_file ../input/RTL/quant_pool/integration/conv_quant_adapter.v VERILOG PATH ../input/RTL/
   quant_pool/integration/conv_quant_adapter.v
50 add_fileset_file ../input/RTL/quant_pool/integration/quant_pool_adapter.v VERILOG PATH ../input/RTL/
   quant_pool/integration/quant_pool_adapter.v
51 add_fileset_file ../input/RTL/quant_pool/quant/Quantization_PE.v VERILOG PATH ../input/RTL/quant_pool/
   quant/Quantization_PE.v
52 add_fileset_file ../input/RTL/quant_pool/quant/Quantization_Top.v VERILOG PATH ../input/RTL/quant_pool/
   quant/Quantization_Top.v
53 add_fileset_file ../input/RTL/quant_pool/pool/pool_core.v VERILOG PATH ../input/RTL/quant_pool/pool/
   pool_core.v
54 add_fileset_file ../input/RTL/quant_pool/pool/pool_stream_top.v VERILOG PATH ../input/RTL/quant_pool/pool
   /pool_stream_top.v
55 add_fileset_file ../input/RTL/fc/mac.v VERILOG PATH ../input/RTL/fc/mac.v
56 add_fileset_file ../input/RTL/fc/fc_bias_loader.v VERILOG PATH ../input/RTL/fc/fc_bias_loader.v
57 add_fileset_file ../input/RTL/fc/fc_data_adapter.v VERILOG PATH ../input/RTL/fc/fc_data_adapter.v

```

```

58 add_fileset_file ../input/RTL/fc/FC.v VERILOG PATH ../input/RTL/fc/FC.v
59 add_fileset_file ../input/RTL/SRAM/Addr_Gen.v VERILOG PATH ../input/RTL/SRAM/Addr_Gen.v
60 add_fileset_file ../input/RTL/SRAM/sram_A_controller.v VERILOG PATH ../input/RTL/SRAM/sram_A_controller.v
61 add_fileset_file ../input/RTL/SRAM/sram_B_controller.v VERILOG PATH ../input/RTL/SRAM/sram_B_controller.v
62 add_fileset_file ../input/RTL/SRAM/sram_A_wrapper.v VERILOG PATH ../input/RTL/SRAM/sram_A_wrapper.v
63 add_fileset_file ../input/RTL/SRAM/sram_B_wrapper.v VERILOG PATH ../input/RTL/SRAM/sram_B_wrapper.v
64 add_fileset_file ../input/RTL/SRAM/sram_FCW_wrapper.v VERILOG PATH ../input/RTL/SRAM/sram_FCW_wrapper.v
65 add_fileset_file ../input/RTL/SRAM/fcw_preload_packer.v VERILOG PATH ../input/RTL/SRAM/
    fcw_preload_packer.v
66 add_fileset_file ../input/RTL/SRAM/top_sram_A.v VERILOG PATH ../input/RTL/SRAM/top_sram_A.v
67 add_fileset_file ../input/RTL/SRAM/top_sram_B.v VERILOG PATH ../input/RTL/SRAM/top_sram_B.v
68
69 # -----
70 # module assignments
71 # -----
72 set_module_assignment embeddedsw.dts.group cnn_mmio_interface
73 set_module_assignment embeddedsw.dts.name cnn_mmio_interface
74 set_module_assignment embeddedsw.dts.vendor csee4840
75
76 # -----
77 # connection point clock
78 # -----
79 add_interface clock clock end
80 set_interface_property clock clockRate 0
81 set_interface_property clock ENABLED true
82 set_interface_property clock EXPORT_OF ""
83 set_interface_property clock PORT_NAME_MAP ""
84 set_interface_property clock CMSIS_SVD_VARIABLES ""
85 set_interface_property clock SVD_ADDRESS_GROUP ""
86 add_interface_port clock clk clk Input 1
87
88 # -----
89 # connection point reset
90 # -----
91 add_interface reset reset end
92 set_interface_property reset associatedClock clock
93 set_interface_property reset synchronousEdges DEASSERT
94 set_interface_property reset ENABLED true
95 set_interface_property reset EXPORT_OF ""
96 set_interface_property reset PORT_NAME_MAP ""
97 set_interface_property reset CMSIS_SVD_VARIABLES ""
98 set_interface_property reset SVD_ADDRESS_GROUP ""
99 add_interface_port reset reset reset Input 1
100
101 # -----
102 # connection point avalon_slave_0
103 # -----
104 add_interface avalon_slave_0 avalon end
105 set_interface_property avalon_slave_0 addressUnits WORDS
106 set_interface_property avalon_slave_0 associatedClock clock
107 set_interface_property avalon_slave_0 associatedReset reset
108 set_interface_property avalon_slave_0 bitsPerSymbol 8
109 set_interface_property avalon_slave_0 burstOnBurstBoundariesOnly false
110 set_interface_property avalon_slave_0 burstcountUnits WORDS

```

```

111 set_interface_property avalon_slave_0 explicitAddressSpan 0
112 set_interface_property avalon_slave_0 holdTime 0
113 set_interface_property avalon_slave_0 linewidthBursts false
114 set_interface_property avalon_slave_0 maximumPendingReadTransactions 0
115 set_interface_property avalon_slave_0 maximumPendingWriteTransactions 0
116 set_interface_property avalon_slave_0 readLatency 2
117 set_interface_property avalon_slave_0 readWaitTime 0
118 set_interface_property avalon_slave_0 setupTime 0
119 set_interface_property avalon_slave_0 timingUnits Cycles
120 set_interface_property avalon_slave_0 writeWaitTime 0
121 set_interface_property avalon_slave_0 ENABLED true
122 set_interface_property avalon_slave_0 EXPORT_OF ""
123 set_interface_property avalon_slave_0 PORT_NAME_MAP ""
124 set_interface_property avalon_slave_0 CMSIS_SVD_VARIABLES ""
125 set_interface_property avalon_slave_0 SVD_ADDRESS_GROUP ""
126 add_interface_port avalon_slave_0 writedata writedata Input 32
127 add_interface_port avalon_slave_0 byteenable byteenable Input 4
128 add_interface_port avalon_slave_0 write write Input 1
129 add_interface_port avalon_slave_0 read read Input 1
130 add_interface_port avalon_slave_0 chipselect chipselect Input 1
131 add_interface_port avalon_slave_0 address address Input 13
132 add_interface_port avalon_slave_0 readdata readdata Output 32
133 set_interface_assignment avalon_slave_0 embeddedsw.configuration.isFlash 0
134 set_interface_assignment avalon_slave_0 embeddedsw.configuration.isMemoryDevice 0
135 set_interface_assignment avalon_slave_0 embeddedsw.configuration.isNonVolatileStorage 0
136 set_interface_assignment avalon_slave_0 embeddedsw.configuration.isPrintableDevice 0

```

## A.6.2 Quartus Build Wrapper

*Discovers Quartus tools, regenerates Qsys HDL, patches generated debug ports, and runs the compile flow.*

```
code/de1_soc/build_soc_system.py
```

```

1  #!/usr/bin/env python3
2  """Cross-platform build wrapper for the DE1-SoC Quartus project.
3
4  Usage:
5      python de1_soc/build_soc_system.py
6  """
7
8  from __future__ import annotations
9
10 import argparse
11 import os
12 import platform
13 import shutil
14 import subprocess
15 import sys
16 from pathlib import Path

```

```

17
18
19 PROJECT_DIR = Path(__file__).resolve().parent
20 PROJECT_NAME = "soc_system"
21 QSYS_FILE = PROJECT_DIR / f"{PROJECT_NAME}.qsys"
22 SOF_FILE = PROJECT_DIR / "output_files" / f"{PROJECT_NAME}.sof"
23 RBF_FILE = PROJECT_DIR / "output_files" / f"{PROJECT_NAME}.rbf"
24 GENERATED_SOC_SYSTEM_V = PROJECT_DIR / "soc_system" / "synthesis" / "soc_system.v"
25
26
27 def candidate_roots() -> list[Path]:
28     env_root = os.environ.get("QUARTUS_ROOTDIR")
29     roots: list[Path] = []
30     if env_root:
31         roots.append(Path(env_root))
32
33     system = platform.system()
34     if system == "Windows":
35         roots.extend(
36             [
37                 Path(r"C:\intelFPGA_lite"),
38                 Path(r"C:\intelFPGA"),
39                 Path(r"C:\Program Files\intelFPGA_lite"),
40                 Path(r"C:\Program Files\intelFPGA"),
41             ]
42         )
43     elif system == "Darwin":
44         roots.extend(
45             [
46                 Path("/Applications/intelFPGA_lite"),
47                 Path("/Applications/intelFPGA"),
48                 Path.home() / "intelFPGA_lite",
49                 Path.home() / "intelFPGA",
50             ]
51         )
52     else:
53         roots.extend(
54             [
55                 Path("/opt/intelFPGA_lite"),
56                 Path("/opt/intelFPGA"),
57                 Path.home() / "intelFPGA_lite",
58                 Path.home() / "intelFPGA",
59             ]
60         )
61     return roots
62
63
64 def tool_names(base: str) -> list[str]:
65     system = platform.system()
66     if system == "Windows":
67         return [f"{base}.exe", base]
68     return [base]
69
70

```

```

71 def env_var_name(base: str) -> str:
72     return base.replace("-", "_").upper()
73
74
75 def find_tool(base: str, verbose: bool = False) -> tuple[str | None, list[str]]:
76     searched: list[str] = []
77     override_var = env_var_name(base)
78     override = os.environ.get(override_var)
79     if override:
80         searched.append(f"{override_var}={override}")
81         if Path(override).exists():
82             return override, searched
83
84     for name in tool_names(base):
85         searched.append(f"PATH:{name}")
86         found = shutil.which(name)
87         if found:
88             return found, searched
89
90     for root in candidate_roots():
91         searched.append(f"ROOT:{root}")
92         if not root.exists():
93             continue
94
95         if root.name == "quartus":
96             search_bins = [
97                 root / "bin64",
98                 root / "bin",
99                 root / "sopc_builder" / "bin",
100             ]
101         else:
102             search_bins = []
103             for version_dir in sorted(root.glob("*")):
104                 if not version_dir.is_dir():
105                     continue
106                 quartus_dir = version_dir / "quartus"
107                 if quartus_dir.exists():
108                     search_bins.extend(
109                         [
110                             quartus_dir / "bin64",
111                             quartus_dir / "bin",
112                             quartus_dir / "sopc_builder" / "bin",
113                         ]
114                     )
115
116         for bin_dir in search_bins:
117             for name in tool_names(base):
118                 candidate = bin_dir / name
119                 searched.append(str(candidate))
120                 if candidate.exists():
121                     return str(candidate), searched
122
123     if verbose:
124         for item in searched:

```

```

125         print(f" searched {item}")
126
127     return None, searched
128
129
130 def build_env_for_subprocesses(quartus_sh: str | None) -> dict[str, str]:
131     env = os.environ.copy()
132
133     if platform.system() != "Windows" or not quartus_sh:
134         return env
135
136     quartus_path = Path(quartus_sh).resolve()
137     acds_root = quartus_path.parent.parent
138     nios_gnu_bin = acds_root / "nios2eds" / "bin" / "gnu" / "H-x86_64-mingw32" / "bin"
139
140     if nios_gnu_bin.exists():
141         env["PATH"] = str(nios_gnu_bin) + os.pathsep + env.get("PATH", "")
142
143     return env
144
145
146 def run(cmd: list[str], cwd: Path, env: dict[str, str] | None = None) -> None:
147     print("+", " ".join(f'"{part}"' if " " in part else part for part in cmd))
148     subprocess.run(cmd, cwd=str(cwd), env=env, check=True)
149
150
151 def run_capture(cmd: list[str], cwd: Path, env: dict[str, str] | None = None) -> subprocess.
    CompletedProcess[str]:
152     print("+", " ".join(f'"{part}"' if " " in part else part for part in cmd))
153     return subprocess.run(
154         cmd,
155         cwd=str(cwd),
156         env=env,
157         check=False,
158         text=True,
159         stdout=subprocess.PIPE,
160         stderr=subprocess.STDOUT,
161     )
162
163
164 def is_recoverable_windows_qsys_failure(output: str) -> bool:
165     markers = [
166         "generate_hps_sdram.tcl",
167         "Nios II Command Shell.bat",
168         "child process exited abnormally",
169         "Error: border: Error during execution of script generate_hps_sdram.tcl: seq:",
170     ]
171     return platform.system() == "Windows" and all(marker in output for marker in markers)
172
173
174 def generated_hdl_looks_usable() -> bool:
175     required = [
176         PROJECT_DIR / "soc_system" / "synthesis" / "soc_system.v",
177         PROJECT_DIR / "soc_system" / "synthesis" / "submodules" / "hps_sdram.v",

```

```

178     PROJECT_DIR / "soc_system" / "synthesis" / "submodules" / "hps_sdrnm_p0.sv",
179 ]
180 return all(path.exists() for path in required)
181
182
183 def patch_generated_soc_system_debug_ports() -> None:
184     if not GENERATED_SOC_SYSTEM_V.exists():
185         raise FileNotFoundError(f"generated top not found: {GENERATED_SOC_SYSTEM_V}")
186
187     text = GENERATED_SOC_SYSTEM_V.read_text(encoding="utf-8")
188
189     if "cnn_debug_predict_class" in text:
190         return
191
192     module_old = "module soc_system (\n\t\tinput wire      clk_clk,"
193     module_new = (
194         "module soc_system (\n"
195         "\t\tinput wire      clk_clk,\n"
196         "\t\toutput wire      cnn_debug_model_loaded,\n"
197         "\t\toutput wire      cnn_debug_predict_done,\n"
198         "\t\toutput wire [3:0] cnn_debug_predict_class,\n"
199         "\t\toutput wire [15:0] cnn_debug_interface_error,"
200     )
201
202     inst_old = (
203         "\t\t.address      (mm_interconnect_0_cnn_mmio_interface_0_avalon_slave_0_address), //
204         .address\n"
205         "\t\t.readdata      (mm_interconnect_0_cnn_mmio_interface_0_avalon_slave_0_readdata) //
206         .readdata\n"
207         "\t);;"
208     )
209
210     inst_new = (
211         "\t\t.address      (mm_interconnect_0_cnn_mmio_interface_0_avalon_slave_0_address), //
212         .address\n"
213         "\t\t.readdata      (mm_interconnect_0_cnn_mmio_interface_0_avalon_slave_0_readdata), //
214         .readdata\n"
215         "\t\t.debug_model_loaded      (cnn_debug_model_loaded), //
216         .debug_model_loaded\n"
217         "\t\t.debug_predict_done      (cnn_debug_predict_done), //
218         .debug_predict_done\n"
219         "\t\t.debug_predict_class      (cnn_debug_predict_class), //
220         .debug_predict_class\n"
221         "\t\t.debug_interface_error(cnn_debug_interface_error) //
222         .debug_interface_error\n"
223         "\t);;"
224     )
225
226     if module_old not in text:
227         raise RuntimeError("could not find soc_system module header to patch debug ports")
228     if inst_old not in text:
229         raise RuntimeError("could not find cnn_mmio_interface instantiation to patch debug ports")
230
231     text = text.replace(module_old, module_new, 1)
232     text = text.replace(inst_old, inst_new, 1)

```

```

224     GENERATED_SOC_SYSTEM_V.write_text(text, encoding="utf-8")
225
226
227 def check_windows_qsys_prereqs(quartus_sh: str | None) -> tuple[bool, list[str]]:
228     if platform.system() != "Windows" or not quartus_sh:
229         return True, []
230
231     messages: list[str] = []
232     quartus_path = Path(quartus_sh).resolve()
233     quartus_root = quartus_path.parent.parent
234     nios2_shell = quartus_root.parent / "nios2eds" / "Nios II Command Shell.bat"
235     mingw_bin = quartus_root.parent / "nios2eds" / "bin" / "gnu" / "H-x86_64-mingw32" / "bin"
236
237     ok = True
238     if not nios2_shell.exists():
239         messages.append(f"missing Nios II Command Shell: {nios2_shell}")
240         ok = False
241     if not mingw_bin.exists():
242         messages.append(f"missing Nios II GCC toolchain directory: {mingw_bin}")
243         ok = False
244
245     try:
246         subprocess.run(
247             ["wsl.exe", "bash", "-lc", "command -v dos2unix"],
248             stdout=subprocess.DEVNULL,
249             stderr=subprocess.DEVNULL,
250             check=True,
251         )
252     except subprocess.CalledProcessError:
253         messages.append("WSL is missing dos2unix (install with: wsl -u root bash -lc \"apt-get install -y dos2unix\")")
254         ok = False
255     except FileNotFoundError:
256         messages.append("WSL is not available on this machine")
257         ok = False
258
259     return ok, messages
260
261
262 def parse_args() -> argparse.Namespace:
263     parser = argparse.ArgumentParser(description="Build the DE1-SoC Quartus project.")
264     parser.add_argument(
265         "--check",
266         action="store_true",
267         help="Only check tool discovery and print what would be used.",
268     )
269     parser.add_argument(
270         "--verbose",
271         action="store_true",
272         help="Print searched locations while resolving Quartus tools.",
273     )
274     parser.add_argument(
275         "--skip-qsys",
276         action="store_true",

```

```

277     help="Skip qsys-generate and compile using the committed generated HDL.",
278 )
279 return parser.parse_args()
280
281
282 def main() -> int:
283     args = parse_args()
284     qsys_generate, qsys_search = find_tool("qsys-generate", verbose=args.verbose)
285     quartus_sh, quartus_sh_search = find_tool("quartus_sh", verbose=args.verbose)
286     quartus_cpf, quartus_cpf_search = find_tool("quartus_cpf", verbose=args.verbose)
287
288     missing = []
289     if not qsys_generate and not args.skip_qsys:
290         missing.append("qsys-generate")
291     if not quartus_sh:
292         missing.append("quartus_sh")
293
294     if missing:
295         print("Required Quartus tools were not found:", ", ".join(missing), file=sys.stderr)
296         print(
297             "Tried PATH, tool-specific env vars, QUARTUS_ROOTDIR, and common install folders.",
298             file=sys.stderr,
299         )
300         print("Optional overrides:", file=sys.stderr)
301         print("  QSYS_GENERATE=/path/to/qsys-generate", file=sys.stderr)
302         print("  QUARTUS_SH=/path/to/quartus_sh", file=sys.stderr)
303         print("  QUARTUS_CPF=/path/to/quartus_cpf", file=sys.stderr)
304         return 1
305
306     prereq_ok, prereq_messages = check_windows_qsys_prereqs(quartus_sh)
307     build_env = build_env_for_subprocesses(quartus_sh)
308     if args.check:
309         print("Quartus tool check passed.")
310         if args.skip_qsys:
311             print("  qsys-generate: skipped by request")
312         else:
313             print(f"  qsys-generate: {qsys_generate}")
314             print(f"  quartus_sh: {quartus_sh}")
315             if quartus_cpf:
316                 print(f"  quartus_cpf: {quartus_cpf}")
317             else:
318                 print("  quartus_cpf: not found (RBF conversion will be skipped)")
319     if platform.system() == "Windows":
320         if prereq_ok:
321             print("  windows qsys prereqs: OK")
322             quartus_path = Path(quartus_sh).resolve()
323             acds_root = quartus_path.parent.parent.parent
324             mingw_bin = acds_root / "nios2eds" / "bin" / "gnu" / "H-x86_64-mingw32" / "bin"
325             print(f"  injected Nios GCC path: {mingw_bin}")
326         else:
327             print("  windows qsys prereqs: MISSING")
328             for message in prereq_messages:
329                 print(f"    - {message}")
330         return 1

```

```

331     return 0
332
333 if not prereq_ok:
334     print("Windows Qsys prerequisites are incomplete:", file=sys.stderr)
335     for message in prereq_messages:
336         print(f" - {message}", file=sys.stderr)
337     return 1
338
339 try:
340     if args.skip_qsys:
341         print("[1/3] Skipping Qsys generation and using committed generated HDL...")
342     else:
343         print("[1/3] Generating Qsys HDL...")
344         qsys_result = run_capture(
345             [qsys_generate, str(QSYS_FILE), "--synthesis=VERILOG"],
346             cwd=PROJECT_DIR,
347             env=build_env,
348         )
349         sys.stdout.write(qsys_result.stdout)
350         if qsys_result.returncode != 0:
351             if is_recoverable_windows_qsys_failure(qsys_result.stdout) and generated_hdl_looks_usable
352             (:):
353                 print()
354                 print(
355                     "Warning: qsys-generate hit the known Windows/WSL HPS SDRAM "
356                     "sequencer-build failure, but the HDL outputs were still emitted."
357                 )
358                 print("Continuing to Quartus compile with the freshly generated HDL.")
359             else:
360                 print(f"Build failed with exit code {qsys_result.returncode}.", file=sys.stderr)
361                 return qsys_result.returncode
362
363         patch_generated_soc_system_debug_ports()
364
365         print("[2/3] Compiling Quartus project...")
366         run([quartus_sh, "--flow", "compile", PROJECT_NAME], cwd=PROJECT_DIR, env=build_env)
367
368         if not SOF_FILE.exists():
369             print(f"Expected SOF not found: {SOF_FILE}", file=sys.stderr)
370             return 1
371
372         if quartus_cpf:
373             print("[3/3] Converting SOF to RBF...")
374             run([quartus_cpf, "-c", str(SOF_FILE), str(RBF_FILE)], cwd=PROJECT_DIR, env=build_env)
375         else:
376             print("[3/3] quartus_cpf not found, skipping RBF conversion.")
377     except subprocess.CalledProcessError as exc:
378         print(f"Build failed with exit code {exc.returncode}.", file=sys.stderr)
379         return exc.returncode
380
381     print()
382     print(f"Build complete.")
383     print(f"SOF: {SOF_FILE}")
384     if RBF_FILE.exists():

```

```

384     print(f"RBF: {RBF_FILE}")
385     return 0
386
387
388 if __name__ == "__main__":
389     raise SystemExit(main())

```

### A.6.3 Quartus Project Script

*Creates the DE1-SoC Quartus project and assigns board-level files, SDC constraints, and QIP references.*

```

code/de1_soc/soc_system_project.tcl

1 # Generate Quartus project files for the DE1-SoC board.
2 #
3 # Adapted for CNN_ACC:
4 # - expects a generated Platform Designer system under soc_system/
5 # - expects a board-level top named soc_system_top in soc_system_top.sv
6 # - pulls in soc_system/synthesis/soc_system.qip
7 #
8 # Invoke as:
9 # quartus_sh -t de1_soc/soc_system_project.tcl
10
11 set script_dir [file dirname [file normalize [info script]]]
12 set project_name "soc_system"
13 set project_path [file join $script_dir $project_name]
14
15 set systemVerilogSource [file join $script_dir "${project_name}_top.sv"]
16 set qip [file join $script_dir $project_name "synthesis" "${project_name}.qip"]
17
18 project_new $project_path -overwrite
19
20 # Clean out any stale assignments from previous iterations of this project.
21 foreach stale_file {
22     "input/SRAM_macro/sram_A/sram_A.v"
23     "input/SRAM_macro/sram_B/sram_B.v"
24     "/homes/user/stud/fall25/lw3227/CNN_ACC/input/SRAM_macro/sram_A/sram_A.v"
25     "/homes/user/stud/fall25/lw3227/CNN_ACC/input/SRAM_macro/sram_B/sram_B.v"
26 } {
27     catch { remove_global_assignment -name VERILOG_FILE $stale_file }
28 }
29
30 foreach {name value} {
31     FAMILY "Cyclone V"
32     DEVICE 5CSEMA5F31C6
33     PROJECT_OUTPUT_DIRECTORY output_files
34     CYCLONEII_RESERVE_NCEO_AFTER_CONFIGURATION "USE AS REGULAR IO"
35     NUM_PARALLEL_PROCESSORS 4
36 } {

```

```

37     set_global_assignment -name $name $value
38 }
39
40 set_global_assignment -name TOP_LEVEL_ENTITY "${project_name}_top"
41 set_global_assignment -name SYSTEMVERILOG_FILE $systemVerilogSource
42 set_global_assignment -name QIP_FILE $qip
43 # Match the original Quartus project: use synthesizable behavioral SRAM
44 # models instead of the foundry macro models, which contain sim-only code.
45 set_global_assignment -name VERILOG_FILE [file join [file dirname $script_dir] "input" "TB" "sram_behav.v
    "]
46
47 # Minimal board-level assignment set for the generated HPS/MMIO system.
48 # Follow the lab3-hw pattern: export the HPS DDR3 interface using the board-
49 # level HPS_DDR3_* port names and let the generated HPS SDRAM assignment Tcl
50 # handle the dedicated DDR3 pin standards/termination.
51 set_location_assignment PIN_AF14 -to CLOCK_50
52 set_instance_assignment -name IO_STANDARD "3.3-V LVTTTL" -to CLOCK_50
53
54 set sdcFilename [file join $script_dir "${project_name}.sdc"]
55 set_global_assignment -name SDC_FILE $sdcFilename
56
57 set sdcf [open $sdcFilename "w"]
58 puts $sdcf {
59     create_clock -name clock_50 -period 20ns [get_ports CLOCK_50]
60     derive_pll_clocks -create_base_clocks
61     derive_clock_uncertainty
62 }
63 close $sdcf
64
65 project_close

```

#### A.6.4 Platform Designer System

*Stores the HPS, clock, reset, interconnect, and CNN MMIO component connections used by qsys-generate.*

code/de1\_soc/soc\_system.qsys

```

1 <?xml version="1.0" encoding="UTF-8"?>
2 <system name="${FILENAME}">
3   <component
4     name="${FILENAME}"
5     displayName="${FILENAME}"
6     version="1.0"
7     description=""
8     tags=""
9     categories="System" />
10  <parameter name="bonusData"><![CDATA[bonusData
11  {
12    element clk_0

```

```

13  {
14      datum _sortIndex
15      {
16          value = "0";
17          type = "int";
18      }
19  }
20  element cnn_mmio_interface_0
21  {
22      datum _sortIndex
23      {
24          value = "1";
25          type = "int";
26      }
27  }
28  element hps_0
29  {
30      datum _sortIndex
31      {
32          value = "2";
33          type = "int";
34      }
35  }
36  }
37  ]]></parameter>
38  <parameter name="clockCrossingAdapter" value="HANDSHAKE" />
39  <parameter name="device" value="5CSEMA5F31C6" />
40  <parameter name="deviceFamily" value="Cyclone V" />
41  <parameter name="deviceSpeedGrade" value="6" />
42  <parameter name="fabricMode" value="QSYS" />
43  <parameter name="generateLegacySim" value="false" />
44  <parameter name="generationId" value="0" />
45  <parameter name="globalResetBus" value="false" />
46  <parameter name="hdlLanguage" value="VERILOG" />
47  <parameter name="hideFromIPCatalog" value="false" />
48  <parameter name="lockedInterfaceDefinition" value="" />
49  <parameter name="maxAdditionalLatency" value="1" />
50  <parameter name="projectName" value="soc_system.qpf" />
51  <parameter name="sopcBorderPoints" value="false" />
52  <parameter name="systemHash" value="0" />
53  <parameter name="testBenchDutName" value="" />
54  <parameter name="timeStamp" value="0" />
55  <parameter name="useTestBenchNamingPattern" value="false" />
56  <instanceScript></instanceScript>
57  <interface name="clk" internal="clk_0.clk_in" type="clock" dir="end" />
58  <interface name="hps" internal="hps_0.hps_io" type="conduit" dir="end" />
59  <interface name="memory" internal="hps_0.memory" type="conduit" dir="end" />
60  <interface name="reset" internal="clk_0.clk_in_reset" type="reset" dir="end" />
61  <module name="clk_0" kind="clock_source" version="21.1" enabled="1">
62      <parameter name="clockFrequency" value="50000000" />
63      <parameter name="clockFrequencyKnown" value="true" />
64      <parameter name="inputClockFrequency" value="0" />
65      <parameter name="resetSynchronousEdges" value="NONE" />
66  </module>

```

```

67 <module
68     name="cnn_mmio_interface_0"
69     kind="cnn_mmio_interface"
70     version="1.0"
71     enabled="1" />
72 <module name="hps_0" kind="altera_hps" version="21.1" enabled="1">
73     <parameter name="ABSTRACT_REAL_COMPARE_TEST" value="false" />
74     <parameter name="ABS_RAM_MEM_INIT_FILENAME" value="meminit" />
75     <parameter name="ACV_PHY_CLK_ADD_FR_PHASE" value="0.0" />
76     <parameter name="AC_PACKAGE_DESKEW" value="false" />
77     <parameter name="AC_ROM_USER_ADD_0" value="0_0000_0000_0000" />
78     <parameter name="AC_ROM_USER_ADD_1" value="0_0000_0000_1000" />
79     <parameter name="ADDR_ORDER" value="0" />
80     <parameter name="ADD_EFFICIENCY_MONITOR" value="false" />
81     <parameter name="ADD_EXTERNAL_SEQ_DEBUG_NIOS" value="false" />
82     <parameter name="ADVANCED_CK_PHASES" value="false" />
83     <parameter name="ADVERTISE_SEQUENCER_SW_BUILD_FILES" value="false" />
84     <parameter name="AFI_DEBUG_INFO_WIDTH" value="32" />
85     <parameter name="ALTMEMPHY_COMPATIBLE_MODE" value="false" />
86     <parameter name="AP_MODE" value="false" />
87     <parameter name="AP_MODE_EN" value="0" />
88     <parameter name="AUTO_DEVICE_SPEEDGRADE" value="6" />
89     <parameter name="AUTO_PD_CYCLES" value="0" />
90     <parameter name="AUTO_POWERDN_EN" value="false" />
91     <parameter name="AVL_DATA_WIDTH_PORT" value="32,32,32,32,32,32" />
92     <parameter name="AVL_MAX_SIZE" value="4" />
93     <parameter name="BONDING_OUT_ENABLED" value="false" />
94     <parameter name="BOOTFROMFPGA_Enable" value="false" />
95     <parameter name="BSEL" value="1" />
96     <parameter name="BSEL_EN" value="false" />
97     <parameter name="BYTE_ENABLE" value="true" />
98     <parameter name="C2P_WRITE_CLOCK_ADD_PHASE" value="0.0" />
99     <parameter name="CALIBRATION_MODE" value="Skip" />
100    <parameter name="CALIB_REG_WIDTH" value="8" />
101    <parameter name="CAN0_Mode" value="N/A" />
102    <parameter name="CAN0_PinMuxing" value="Unused" />
103    <parameter name="CAN1_Mode" value="N/A" />
104    <parameter name="CAN1_PinMuxing" value="Unused" />
105    <parameter name="CFG_DATA_REORDERING_TYPE" value="INTER_BANK" />
106    <parameter name="CFG_REORDER_DATA" value="true" />
107    <parameter name="CFG_TCCD_NS" value="2.5" />
108    <parameter name="COMMAND_PHASE" value="0.0" />
109    <parameter name="CONTROLLER_LATENCY" value="5" />
110    <parameter name="CORE_DEBUG_CONNECTION" value="EXPORT" />
111    <parameter name="CPORT_TYPE_PORT">Bidirectional,Bidirectional,Bidirectional,Bidirectional,Bidirectional
        ,Bidirectional</parameter>
112    <parameter name="CSEL" value="0" />
113    <parameter name="CSEL_EN" value="false" />
114    <parameter name="CTI_Enable" value="false" />
115    <parameter name="CTL_AUTOPCH_EN" value="false" />
116    <parameter name="CTL_CMD_QUEUE_DEPTH" value="8" />
117    <parameter name="CTL_CSR_CONNECTION" value="INTERNAL_JTAG" />
118    <parameter name="CTL_CSR_ENABLED" value="false" />
119    <parameter name="CTL_CSR_READ_ONLY" value="1" />

```

```

120 <parameter name="CTL_DEEP_POWERDN_EN" value="false" />
121 <parameter name="CTL_DYNAMIC_BANK_ALLOCATION" value="false" />
122 <parameter name="CTL_DYNAMIC_BANK_NUM" value="4" />
123 <parameter name="CTL_ECC_AUTO_CORRECTION_ENABLED" value="false" />
124 <parameter name="CTL_ECC_ENABLED" value="false" />
125 <parameter name="CTL_ENABLE_BURST_INTERRUPT" value="false" />
126 <parameter name="CTL_ENABLE_BURST_TERMINATE" value="false" />
127 <parameter name="CTL_HRB_ENABLED" value="false" />
128 <parameter name="CTL_LOOK_AHEAD_DEPTH" value="4" />
129 <parameter name="CTL_SELF_REFRESH_EN" value="false" />
130 <parameter name="CTL_USR_REFRESH_EN" value="false" />
131 <parameter name="CTL_ZQCAL_EN" value="false" />
132 <parameter name="CUT_NEW_FAMILY_TIMING" value="true" />
133 <parameter name="DAT_DATA_WIDTH" value="32" />
134 <parameter name="DEBUGAPB_Enable" value="false" />
135 <parameter name="DEBUG_MODE" value="false" />
136 <parameter name="DEVICE_DEPTH" value="1" />
137 <parameter name="DEVICE_FAMILY_PARAM" value="" />
138 <parameter name="DISABLE_CHILD_MESSAGING" value="false" />
139 <parameter name="DISCRETE_FLY_BY" value="true" />
140 <parameter name="DLL_SHARING_MODE" value="None" />
141 <parameter name="DMA_Enable">No,No,No,No,No,No,No,No</parameter>
142 <parameter name="DQS_DQSN_MODE" value="DIFFERENTIAL" />
143 <parameter name="DQ_INPUT_REG_USE_CLKN" value="false" />
144 <parameter name="DUPLICATE_AC" value="false" />
145 <parameter name="ED_EXPORT_SEQ_DEBUG" value="false" />
146 <parameter name="EMACO_Mode" value="N/A" />
147 <parameter name="EMACO_PTP" value="false" />
148 <parameter name="EMACO_PinMuxing" value="Unused" />
149 <parameter name="EMAC1_Mode" value="RGMI" />
150 <parameter name="EMAC1_PTP" value="false" />
151 <parameter name="EMAC1_PinMuxing" value="HPS I/O Set 0" />
152 <parameter name="ENABLE_ABS_RAM_MEM_INIT" value="false" />
153 <parameter name="ENABLE_BONDING" value="false" />
154 <parameter name="ENABLE_BURST_MERGE" value="false" />
155 <parameter name="ENABLE_CTRL_AVALON_INTERFACE" value="true" />
156 <parameter name="ENABLE_DELAY_CHAIN_WRITE" value="false" />
157 <parameter name="ENABLE_EMIT_BFM_MASTER" value="false" />
158 <parameter name="ENABLE_EXPORT_SEQ_DEBUG_BRIDGE" value="false" />
159 <parameter name="ENABLE_EXTRA_REPORTING" value="false" />
160 <parameter name="ENABLE_ISS_PROBES" value="false" />
161 <parameter name="ENABLE_NON_DESTRUCTIVE_CALIB" value="false" />
162 <parameter name="ENABLE_NON_DES_CAL" value="false" />
163 <parameter name="ENABLE_NON_DES_CAL_TEST" value="false" />
164 <parameter name="ENABLE_SEQUENCER_MARGINING_ON_BY_DEFAULT" value="false" />
165 <parameter name="ENABLE_USER_ECC" value="false" />
166 <parameter name="EXPORT_AFI_HALF_CLK" value="false" />
167 <parameter name="EXTRA_SETTINGS" value="" />
168 <parameter name="F2H_AXI_CLOCK_FREQ" value="50000000" />
169 <parameter name="F2H_SDRAM0_CLOCK_FREQ" value="50000000" />
170 <parameter name="F2H_SDRAM1_CLOCK_FREQ" value="100" />
171 <parameter name="F2H_SDRAM2_CLOCK_FREQ" value="100" />
172 <parameter name="F2H_SDRAM3_CLOCK_FREQ" value="100" />
173 <parameter name="F2H_SDRAM4_CLOCK_FREQ" value="100" />

```





```

276 <parameter name="MEM_IF_SIM_VALID_WINDOW" value="0" />
277 <parameter name="MEM_INIT_EN" value="false" />
278 <parameter name="MEM_INIT_FILE" value="" />
279 <parameter name="MEM_MIRROR_ADDRESSING" value="0" />
280 <parameter name="MEM_NUMBER_OF_DIMMS" value="1" />
281 <parameter name="MEM_NUMBER_OF_RANKS_PER_DEVICE" value="1" />
282 <parameter name="MEM_NUMBER_OF_RANKS_PER_DIMM" value="1" />
283 <parameter name="MEM_PD" value="DLL off" />
284 <parameter name="MEM_RANK_MULTIPLICATION_FACTOR" value="1" />
285 <parameter name="MEM_ROW_ADDR_WIDTH" value="15" />
286 <parameter name="MEM_RTT_NOM" value="ODT Disabled" />
287 <parameter name="MEM_RTT_WR" value="Dynamic ODT off" />
288 <parameter name="MEM_SRT" value="Normal" />
289 <parameter name="MEM_TCL" value="7" />
290 <parameter name="MEM_TFAW_NS" value="40.0" />
291 <parameter name="MEM_TINIT_US" value="500" />
292 <parameter name="MEM_TMRD_CK" value="4" />
293 <parameter name="MEM_TRAS_NS" value="37.5" />
294 <parameter name="MEM_TRCD_NS" value="12.5" />
295 <parameter name="MEM_TREFI_US" value="7.8" />
296 <parameter name="MEM_TRFC_NS" value="90.0" />
297 <parameter name="MEM_TRP_NS" value="12.5" />
298 <parameter name="MEM_TRRD_NS" value="13.2" />
299 <parameter name="MEM_TRTP_NS" value="13.2" />
300 <parameter name="MEM_TWR_NS" value="15.0" />
301 <parameter name="MEM_TWTR" value="4" />
302 <parameter name="MEM_USER_LEVELING_MODE" value="Leveling" />
303 <parameter name="MEM_VENDOR" value="JEDEC" />
304 <parameter name="MEM_VERBOSE" value="true" />
305 <parameter name="MEM_VOLTAGE" value="1.5V DDR3" />
306 <parameter name="MEM_WTCL" value="6" />
307 <parameter name="MPU_EVENTS_Enable" value="false" />
308 <parameter name="MRS_MIRROR_PING_PONG_ATSO" value="false" />
309 <parameter name="MULTICAST_EN" value="false" />
310 <parameter name="NAND_Mode" value="N/A" />
311 <parameter name="NAND_PinMuxing" value="Unused" />
312 <parameter name="NEXTGEN" value="true" />
313 <parameter name="NIOS_ROM_DATA_WIDTH" value="32" />
314 <parameter name="NUM_DLL_SHARING_INTERFACES" value="1" />
315 <parameter name="NUM_EXTRA_REPORT_PATH" value="10" />
316 <parameter name="NUM_OCT_SHARING_INTERFACES" value="1" />
317 <parameter name="NUM_OF_PORTS" value="1" />
318 <parameter name="NUM_PLL_SHARING_INTERFACES" value="1" />
319 <parameter name="OCT_SHARING_MODE" value="None" />
320 <parameter name="P2C_READ_CLOCK_ADD_PHASE" value="0.0" />
321 <parameter name="PACKAGE_DESKEW" value="false" />
322 <parameter name="PARSE_FRIENDLY_DEVICE_FAMILY_PARAM" value="" />
323 <parameter name="PARSE_FRIENDLY_DEVICE_FAMILY_PARAM_VALID" value="false" />
324 <parameter name="PHY_CSR_CONNECTION" value="INTERNAL_JTAG" />
325 <parameter name="PHY_CSR_ENABLED" value="false" />
326 <parameter name="PHY_ONLY" value="false" />
327 <parameter name="PINGPONGPHY_EN" value="false" />
328 <parameter name="PLL_ADDR_CMD_CLK_DIV_PARAM" value="0" />
329 <parameter name="PLL_ADDR_CMD_CLK_FREQ_PARAM" value="0.0" />

```

```

330 <parameter name="PLL_ADDR_CMD_CLK_FREQ_SIM_STR_PARAM" value="" />
331 <parameter name="PLL_ADDR_CMD_CLK_MULT_PARAM" value="0" />
332 <parameter name="PLL_ADDR_CMD_CLK_PHASE_PS_PARAM" value="0" />
333 <parameter name="PLL_ADDR_CMD_CLK_PHASE_PS_SIM_STR_PARAM" value="" />
334 <parameter name="PLL_AFI_CLK_DIV_PARAM" value="0" />
335 <parameter name="PLL_AFI_CLK_FREQ_PARAM" value="0.0" />
336 <parameter name="PLL_AFI_CLK_FREQ_SIM_STR_PARAM" value="" />
337 <parameter name="PLL_AFI_CLK_MULT_PARAM" value="0" />
338 <parameter name="PLL_AFI_CLK_PHASE_PS_PARAM" value="0" />
339 <parameter name="PLL_AFI_CLK_PHASE_PS_SIM_STR_PARAM" value="" />
340 <parameter name="PLL_AFI_HALF_CLK_DIV_PARAM" value="0" />
341 <parameter name="PLL_AFI_HALF_CLK_FREQ_PARAM" value="0.0" />
342 <parameter name="PLL_AFI_HALF_CLK_FREQ_SIM_STR_PARAM" value="" />
343 <parameter name="PLL_AFI_HALF_CLK_MULT_PARAM" value="0" />
344 <parameter name="PLL_AFI_HALF_CLK_PHASE_PS_PARAM" value="0" />
345 <parameter name="PLL_AFI_HALF_CLK_PHASE_PS_SIM_STR_PARAM" value="" />
346 <parameter name="PLL_AFI_PHY_CLK_DIV_PARAM" value="0" />
347 <parameter name="PLL_AFI_PHY_CLK_FREQ_PARAM" value="0.0" />
348 <parameter name="PLL_AFI_PHY_CLK_FREQ_SIM_STR_PARAM" value="" />
349 <parameter name="PLL_AFI_PHY_CLK_MULT_PARAM" value="0" />
350 <parameter name="PLL_AFI_PHY_CLK_PHASE_PS_PARAM" value="0" />
351 <parameter name="PLL_AFI_PHY_CLK_PHASE_PS_SIM_STR_PARAM" value="" />
352 <parameter name="PLL_C2P_WRITE_CLK_DIV_PARAM" value="0" />
353 <parameter name="PLL_C2P_WRITE_CLK_FREQ_PARAM" value="0.0" />
354 <parameter name="PLL_C2P_WRITE_CLK_FREQ_SIM_STR_PARAM" value="" />
355 <parameter name="PLL_C2P_WRITE_CLK_MULT_PARAM" value="0" />
356 <parameter name="PLL_C2P_WRITE_CLK_PHASE_PS_PARAM" value="0" />
357 <parameter name="PLL_C2P_WRITE_CLK_PHASE_PS_SIM_STR_PARAM" value="" />
358 <parameter name="PLL_CLK_PARAM_VALID" value="false" />
359 <parameter name="PLL_CONFIG_CLK_DIV_PARAM" value="0" />
360 <parameter name="PLL_CONFIG_CLK_FREQ_PARAM" value="0.0" />
361 <parameter name="PLL_CONFIG_CLK_FREQ_SIM_STR_PARAM" value="" />
362 <parameter name="PLL_CONFIG_CLK_MULT_PARAM" value="0" />
363 <parameter name="PLL_CONFIG_CLK_PHASE_PS_PARAM" value="0" />
364 <parameter name="PLL_CONFIG_CLK_PHASE_PS_SIM_STR_PARAM" value="" />
365 <parameter name="PLL_DR_CLK_DIV_PARAM" value="0" />
366 <parameter name="PLL_DR_CLK_FREQ_PARAM" value="0.0" />
367 <parameter name="PLL_DR_CLK_FREQ_SIM_STR_PARAM" value="" />
368 <parameter name="PLL_DR_CLK_MULT_PARAM" value="0" />
369 <parameter name="PLL_DR_CLK_PHASE_PS_PARAM" value="0" />
370 <parameter name="PLL_DR_CLK_PHASE_PS_SIM_STR_PARAM" value="" />
371 <parameter name="PLL_HR_CLK_DIV_PARAM" value="0" />
372 <parameter name="PLL_HR_CLK_FREQ_PARAM" value="0.0" />
373 <parameter name="PLL_HR_CLK_FREQ_SIM_STR_PARAM" value="" />
374 <parameter name="PLL_HR_CLK_MULT_PARAM" value="0" />
375 <parameter name="PLL_HR_CLK_PHASE_PS_PARAM" value="0" />
376 <parameter name="PLL_HR_CLK_PHASE_PS_SIM_STR_PARAM" value="" />
377 <parameter name="PLL_LOCATION" value="Top_Bottom" />
378 <parameter name="PLL_MEM_CLK_DIV_PARAM" value="0" />
379 <parameter name="PLL_MEM_CLK_FREQ_PARAM" value="0.0" />
380 <parameter name="PLL_MEM_CLK_FREQ_SIM_STR_PARAM" value="" />
381 <parameter name="PLL_MEM_CLK_MULT_PARAM" value="0" />
382 <parameter name="PLL_MEM_CLK_PHASE_PS_PARAM" value="0" />
383 <parameter name="PLL_MEM_CLK_PHASE_PS_SIM_STR_PARAM" value="" />

```

```

384 <parameter name="PLL_NIOS_CLK_DIV_PARAM" value="0" />
385 <parameter name="PLL_NIOS_CLK_FREQ_PARAM" value="0.0" />
386 <parameter name="PLL_NIOS_CLK_FREQ_SIM_STR_PARAM" value="" />
387 <parameter name="PLL_NIOS_CLK_MULT_PARAM" value="0" />
388 <parameter name="PLL_NIOS_CLK_PHASE_PS_PARAM" value="0" />
389 <parameter name="PLL_NIOS_CLK_PHASE_PS_SIM_STR_PARAM" value="" />
390 <parameter name="PLL_P2C_READ_CLK_DIV_PARAM" value="0" />
391 <parameter name="PLL_P2C_READ_CLK_FREQ_PARAM" value="0.0" />
392 <parameter name="PLL_P2C_READ_CLK_FREQ_SIM_STR_PARAM" value="" />
393 <parameter name="PLL_P2C_READ_CLK_MULT_PARAM" value="0" />
394 <parameter name="PLL_P2C_READ_CLK_PHASE_PS_PARAM" value="0" />
395 <parameter name="PLL_P2C_READ_CLK_PHASE_PS_SIM_STR_PARAM" value="" />
396 <parameter name="PLL_SHARING_MODE" value="None" />
397 <parameter name="PLL_WRITE_CLK_DIV_PARAM" value="0" />
398 <parameter name="PLL_WRITE_CLK_FREQ_PARAM" value="0.0" />
399 <parameter name="PLL_WRITE_CLK_FREQ_SIM_STR_PARAM" value="" />
400 <parameter name="PLL_WRITE_CLK_MULT_PARAM" value="0" />
401 <parameter name="PLL_WRITE_CLK_PHASE_PS_PARAM" value="0" />
402 <parameter name="PLL_WRITE_CLK_PHASE_PS_SIM_STR_PARAM" value="" />
403 <parameter name="POWER_OF_TWO_BUS" value="false" />
404 <parameter name="PRIORITY_PORT" value="1,1,1,1,1,1" />
405 <parameter name="QSPI_Mode" value="N/A" />
406 <parameter name="QSPI_PinMuxing" value="Unused" />
407 <parameter name="RATE" value="Full" />
408 <parameter name="RDIMM_CONFIG" value="0" />
409 <parameter name="READ_DQ_DQS_CLOCK_SOURCE" value="INVERTED_DQS_BUS" />
410 <parameter name="READ_FIFO_SIZE" value="8" />
411 <parameter name="REFRESH_BURST_VALIDATION" value="false" />
412 <parameter name="REFRESH_INTERVAL" value="15000" />
413 <parameter name="REF_CLK_FREQ" value="125.0" />
414 <parameter name="REF_CLK_FREQ_MAX_PARAM" value="0.0" />
415 <parameter name="REF_CLK_FREQ_MIN_PARAM" value="0.0" />
416 <parameter name="REF_CLK_FREQ_PARAM_VALID" value="false" />
417 <parameter name="S2FCLK_COLDRST_Enable" value="false" />
418 <parameter name="S2FCLK_PENDINGRST_Enable" value="false" />
419 <parameter name="S2FCLK_USEROCLK_Enable" value="false" />
420 <parameter name="S2FCLK_USER1CLK_Enable" value="false" />
421 <parameter name="S2FCLK_USER1CLK_FREQ" value="100.0" />
422 <parameter name="S2FCLK_USER2CLK" value="5" />
423 <parameter name="S2FCLK_USER2CLK_Enable" value="false" />
424 <parameter name="S2FCLK_USER2CLK_FREQ" value="100.0" />
425 <parameter name="S2FINTERRUPT_CAN_Enable" value="false" />
426 <parameter name="S2FINTERRUPT_CLOCKPERIPHERAL_Enable" value="false" />
427 <parameter name="S2FINTERRUPT_CTI_Enable" value="false" />
428 <parameter name="S2FINTERRUPT_DMA_Enable" value="false" />
429 <parameter name="S2FINTERRUPT_EMAC_Enable" value="false" />
430 <parameter name="S2FINTERRUPT_FPGAMANAGER_Enable" value="false" />
431 <parameter name="S2FINTERRUPT_GPIO_Enable" value="false" />
432 <parameter name="S2FINTERRUPT_I2CEMAC_Enable" value="false" />
433 <parameter name="S2FINTERRUPT_I2CPERIPHERAL_Enable" value="false" />
434 <parameter name="S2FINTERRUPT_L4TIMER_Enable" value="false" />
435 <parameter name="S2FINTERRUPT_NAND_Enable" value="false" />
436 <parameter name="S2FINTERRUPT_OSCTIMER_Enable" value="false" />
437 <parameter name="S2FINTERRUPT_QSPI_Enable" value="false" />

```

```

438 <parameter name="S2FINTERRUPT_SDMMC_Enable" value="false" />
439 <parameter name="S2FINTERRUPT_SPIMASTER_Enable" value="false" />
440 <parameter name="S2FINTERRUPT_SPI_SLAVE_Enable" value="false" />
441 <parameter name="S2FINTERRUPT_UART_Enable" value="false" />
442 <parameter name="S2FINTERRUPT_USB_Enable" value="false" />
443 <parameter name="S2FINTERRUPT_WATCHDOG_Enable" value="false" />
444 <parameter name="S2F_Width" value="2" />
445 <parameter name="SDIO_Mode" value="4-bit Data" />
446 <parameter name="SDIO_PinMuxing" value="HPS I/O Set 0" />
447 <parameter name="SEQUENCER_TYPE" value="NIOS" />
448 <parameter name="SEQ_MODE" value="0" />
449 <parameter name="SKIP_MEM_INIT" value="true" />
450 <parameter name="SOPC_COMPAT_RESET" value="false" />
451 <parameter name="SPEED_GRADE" value="7" />
452 <parameter name="SPIM0_Mode" value="N/A" />
453 <parameter name="SPIM0_PinMuxing" value="Unused" />
454 <parameter name="SPIM1_Mode" value="Single Slave Select" />
455 <parameter name="SPIM1_PinMuxing" value="HPS I/O Set 0" />
456 <parameter name="SPIS0_Mode" value="N/A" />
457 <parameter name="SPIS0_PinMuxing" value="Unused" />
458 <parameter name="SPIS1_Mode" value="N/A" />
459 <parameter name="SPIS1_PinMuxing" value="Unused" />
460 <parameter name="STARVE_LIMIT" value="10" />
461 <parameter name="STM_Enable" value="false" />
462 <parameter name="SYS_INFO_DEVICE_FAMILY" value="Cyclone V" />
463 <parameter name="TEST_Enable" value="false" />
464 <parameter name="TIMING_BOARD_AC_EYE_REDUCTION_H" value="0.0" />
465 <parameter name="TIMING_BOARD_AC_EYE_REDUCTION_SU" value="0.0" />
466 <parameter name="TIMING_BOARD_AC_SKEW" value="0.02" />
467 <parameter name="TIMING_BOARD_AC_SLEW_RATE" value="1.0" />
468 <parameter name="TIMING_BOARD_AC_TO_CK_SKEW" value="0.0" />
469 <parameter name="TIMING_BOARD_CK_CKN_SLEW_RATE" value="2.0" />
470 <parameter name="TIMING_BOARD_DELTA_DQS_ARRIVAL_TIME" value="0.0" />
471 <parameter name="TIMING_BOARD_DELTA_READ_DQS_ARRIVAL_TIME" value="0.0" />
472 <parameter name="TIMING_BOARD_DERATE_METHOD" value="AUTO" />
473 <parameter name="TIMING_BOARD_DQS_DQSN_SLEW_RATE" value="2.0" />
474 <parameter name="TIMING_BOARD_DQ_EYE_REDUCTION" value="0.0" />
475 <parameter name="TIMING_BOARD_DQ_SLEW_RATE" value="1.0" />
476 <parameter name="TIMING_BOARD_DQ_TO_DQS_SKEW" value="0.0" />
477 <parameter name="TIMING_BOARD_ISI_METHOD" value="AUTO" />
478 <parameter name="TIMING_BOARD_MAX_CK_DELAY" value="0.6" />
479 <parameter name="TIMING_BOARD_MAX_DQS_DELAY" value="0.6" />
480 <parameter name="TIMING_BOARD_READ_DQ_EYE_REDUCTION" value="0.0" />
481 <parameter name="TIMING_BOARD_SKEW_BETWEEN_DIMMS" value="0.05" />
482 <parameter name="TIMING_BOARD_SKEW_BETWEEN_DQS" value="0.02" />
483 <parameter name="TIMING_BOARD_SKEW_CKDQS_DIMM_MAX" value="0.01" />
484 <parameter name="TIMING_BOARD_SKEW_CKDQS_DIMM_MIN" value="-0.01" />
485 <parameter name="TIMING_BOARD_SKEW_WITHIN_DQS" value="0.02" />
486 <parameter name="TIMING_BOARD_TDH" value="0.0" />
487 <parameter name="TIMING_BOARD_TDS" value="0.0" />
488 <parameter name="TIMING_BOARD_TIH" value="0.0" />
489 <parameter name="TIMING_BOARD_TIS" value="0.0" />
490 <parameter name="TIMING_TDH" value="150" />
491 <parameter name="TIMING_TDQ_SCK" value="400" />

```

```

492 <parameter name="TIMING_TDQSCKDL" value="1200" />
493 <parameter name="TIMING_TDQSCKDM" value="900" />
494 <parameter name="TIMING_TDQSCKDS" value="450" />
495 <parameter name="TIMING_TDQSH" value="0.35" />
496 <parameter name="TIMING_TDQSQ" value="200" />
497 <parameter name="TIMING_TDQSS" value="0.25" />
498 <parameter name="TIMING_TDS" value="75" />
499 <parameter name="TIMING_TDSH" value="0.2" />
500 <parameter name="TIMING_TDSS" value="0.2" />
501 <parameter name="TIMING_TIH" value="275" />
502 <parameter name="TIMING_TIS" value="200" />
503 <parameter name="TIMING_TQH" value="0.38" />
504 <parameter name="TIMING_TQHS" value="300" />
505 <parameter name="TIMING_TQSH" value="0.38" />
506 <parameter name="TPIUFPGA_Enable" value="false" />
507 <parameter name="TPIUFPGA_alt" value="false" />
508 <parameter name="TRACE_Mode" value="N/A" />
509 <parameter name="TRACE_PinMuxing" value="Unused" />
510 <parameter name="TRACKING_ERROR_TEST" value="false" />
511 <parameter name="TRACKING_WATCH_TEST" value="false" />
512 <parameter name="TREFI" value="35100" />
513 <parameter name="TRFC" value="350" />
514 <parameter name="UART0_Mode" value="No Flow Control" />
515 <parameter name="UART0_PinMuxing" value="HPS I/O Set 0" />
516 <parameter name="UART1_Mode" value="N/A" />
517 <parameter name="UART1_PinMuxing" value="Unused" />
518 <parameter name="USBO_Mode" value="N/A" />
519 <parameter name="USBO_PinMuxing" value="Unused" />
520 <parameter name="USB1_Mode" value="SDR" />
521 <parameter name="USB1_PinMuxing" value="HPS I/O Set 0" />
522 <parameter name="USER_DEBUG_LEVEL" value="1" />
523 <parameter name="USE_AXI_ADAPTOR" value="false" />
524 <parameter name="USE_FAKE_PHY" value="false" />
525 <parameter name="USE_MEM_CLK_FREQ" value="false" />
526 <parameter name="USE_MM_ADAPTOR" value="true" />
527 <parameter name="USE_SEQUENCER_BFM" value="false" />
528 <parameter name="WEIGHT_PORT" value="0,0,0,0,0,0" />
529 <parameter name="WRBUFFER_ADDR_WIDTH" value="6" />
530 <parameter name="can0_clk_div" value="1" />
531 <parameter name="can1_clk_div" value="1" />
532 <parameter name="configure_advanced_parameters" value="false" />
533 <parameter name="customize_device_pll_info" value="false" />
534 <parameter name="dbctrl_stayosci" value="true" />
535 <parameter name="dbg_at_clk_div" value="0" />
536 <parameter name="dbg_clk_div" value="1" />
537 <parameter name="dbg_trace_clk_div" value="0" />
538 <parameter name="desired_can0_clk_mhz" value="100.0" />
539 <parameter name="desired_can1_clk_mhz" value="100.0" />
540 <parameter name="desired_cfg_clk_mhz" value="100.0" />
541 <parameter name="desired_emac0_clk_mhz" value="250.0" />
542 <parameter name="desired_emac1_clk_mhz" value="250.0" />
543 <parameter name="desired_gpio_db_clk_hz" value="32000" />
544 <parameter name="desired_l4_mp_clk_mhz" value="100.0" />
545 <parameter name="desired_l4_sp_clk_mhz" value="100.0" />

```

```

546 <parameter name="desired_mpu_clk_mhz" value="800.0" />
547 <parameter name="desired_nand_clk_mhz" value="12.5" />
548 <parameter name="desired_qspi_clk_mhz" value="400.0" />
549 <parameter name="desired_sdmmc_clk_mhz" value="200.0" />
550 <parameter name="desired_spi_m_clk_mhz" value="200.0" />
551 <parameter name="desired_usb_mp_clk_mhz" value="200.0" />
552 <parameter name="device_name" value="5CSEMA5F31C6" />
553 <parameter name="device_pll_info_manual">{320000000 1600000000} {320000000 1000000000} {800000000
    400000000 400000000}</parameter>
554 <parameter name="eosc1_clk_mhz" value="25.0" />
555 <parameter name="eosc2_clk_mhz" value="25.0" />
556 <parameter name="gpio_db_clk_div" value="6249" />
557 <parameter name="l3_mp_clk_div" value="1" />
558 <parameter name="l3_sp_clk_div" value="1" />
559 <parameter name="l4_mp_clk_div" value="1" />
560 <parameter name="l4_mp_clk_source" value="1" />
561 <parameter name="l4_sp_clk_div" value="1" />
562 <parameter name="l4_sp_clk_source" value="1" />
563 <parameter name="main_pll_c3" value="3" />
564 <parameter name="main_pll_c4" value="3" />
565 <parameter name="main_pll_c5" value="15" />
566 <parameter name="main_pll_m" value="63" />
567 <parameter name="main_pll_n" value="0" />
568 <parameter name="nand_clk_source" value="2" />
569 <parameter name="periph_pll_c0" value="3" />
570 <parameter name="periph_pll_c1" value="3" />
571 <parameter name="periph_pll_c2" value="1" />
572 <parameter name="periph_pll_c3" value="19" />
573 <parameter name="periph_pll_c4" value="4" />
574 <parameter name="periph_pll_c5" value="9" />
575 <parameter name="periph_pll_m" value="79" />
576 <parameter name="periph_pll_n" value="1" />
577 <parameter name="periph_pll_source" value="0" />
578 <parameter name="qspi_clk_source" value="1" />
579 <parameter name="quartus_ini_hps_emif_pll" value="false" />
580 <parameter
581     name="quartus_ini_hps_ip_enable_all_peripheral_fpga_interfaces"
582     value="false" />
583 <parameter name="quartus_ini_hps_ip_enable_bsel_csel" value="false" />
584 <parameter
585     name="quartus_ini_hps_ip_enable_emac0_peripheral_fpga_interface"
586     value="false" />
587 <parameter
588     name="quartus_ini_hps_ip_enable_low_speed_serial_fpga_interfaces"
589     value="false" />
590 <parameter name="quartus_ini_hps_ip_enable_test_interface" value="false" />
591 <parameter name="quartus_ini_hps_ip_f2sdram_bonding_out" value="false" />
592 <parameter name="quartus_ini_hps_ip_fast_f2sdram_sim_model" value="false" />
593 <parameter name="quartus_ini_hps_ip_suppress_sdram_synth" value="false" />
594 <parameter name="sdmmc_clk_source" value="2" />
595 <parameter name="show_advanced_parameters" value="false" />
596 <parameter name="show_debug_info_as_warning_msg" value="false" />
597 <parameter name="show_warning_as_error_msg" value="false" />
598 <parameter name="spi_m_clk_div" value="0" />

```

```

599 <parameter name="usb_mp_clk_div" value="0" />
600 <parameter name="use_default_mpu_clk" value="true" />
601 </module>
602 <connection
603     kind="avalon"
604     version="21.1"
605     start="hps_0.h2f_axi_master"
606     end="cnn_mmio_interface_0.avalon_slave_0">
607 <parameter name="arbitrationPriority" value="1" />
608 <parameter name="baseAddress" value="0x0000" />
609 <parameter name="defaultConnection" value="false" />
610 </connection>
611 <connection
612     kind="avalon"
613     version="21.1"
614     start="hps_0.h2f_lw_axi_master"
615     end="cnn_mmio_interface_0.avalon_slave_0">
616 <parameter name="arbitrationPriority" value="1" />
617 <parameter name="baseAddress" value="0x0000" />
618 <parameter name="defaultConnection" value="false" />
619 </connection>
620 <connection
621     kind="clock"
622     version="21.1"
623     start="clk_0.clk"
624     end="cnn_mmio_interface_0.clock" />
625 <connection
626     kind="clock"
627     version="21.1"
628     start="clk_0.clk"
629     end="hps_0.f2h_axi_clock" />
630 <connection
631     kind="clock"
632     version="21.1"
633     start="clk_0.clk"
634     end="hps_0.f2h_sdram0_clock" />
635 <connection
636     kind="clock"
637     version="21.1"
638     start="clk_0.clk"
639     end="hps_0.h2f_axi_clock" />
640 <connection
641     kind="clock"
642     version="21.1"
643     start="clk_0.clk"
644     end="hps_0.h2f_lw_axi_clock" />
645 <connection
646     kind="reset"
647     version="21.1"
648     start="clk_0.clk_reset"
649     end="cnn_mmio_interface_0.reset" />
650 <connection
651     kind="reset"
652     version="21.1"

```

```

653     start="hps_0.h2f_reset"
654     end="clk_0.clk_in_reset" />
655 <interconnectRequirement for="$system" name="qsys_mm.clockCrossingAdapter" value="HANDSHAKE" />
656 <interconnectRequirement for="$system" name="qsys_mm.enableEccProtection" value="FALSE" />
657 <interconnectRequirement for="$system" name="qsys_mm.insertDefaultSlave" value="FALSE" />
658 <interconnectRequirement for="$system" name="qsys_mm.maxAdditionalLatency" value="1" />
659 </system>

```

## A.6.5 Timing Constraint File

*Defines the board clock, derived PLL clock handling, and clock uncertainty constraints used by Quartus timing analysis.*

```
code/de1_soc/soc_system.sdc
```

```

1
2     create_clock -name clock_50 -period 20ns [get_ports CLOCK_50]
3     derive_pll_clocks -create_base_clocks
4     derive_clock_uncertainty

```

## A.7 FPGA Board Boundary And Accelerator Top

### A.7.1 DE1-SoC Board Wrapper

*Connects the generated HPS system, board pins, clocks/resets, and seven-segment debug display.*

```
code/de1_soc/soc_system_top.sv
```

```

1 module soc_system_top (
2     input wire      CLOCK_50,
3     output wire [6:0] HEX0,
4     output wire [6:0] HEX1,
5     output wire [6:0] HEX2,
6     output wire [6:0] HEX3,
7     output wire [6:0] HEX4,
8     output wire [6:0] HEX5,
9     input wire     HPS_UART_RX,
10    output wire    HPS_UART_TX,
11    output wire    HPS_SD_CLK,
12    inout wire     HPS_SD_CMD,
13    inout wire [3:0] HPS_SD_DATA,
14    inout wire     HPS_I2C1_SCLK,
15    inout wire     HPS_I2C1_SDAT,
16    inout wire     HPS_I2C2_SCLK,
17    inout wire     HPS_I2C2_SDAT,
18    inout wire     HPS_I2C_CONTROL,

```

```

19     output wire      HPS_ENET_GTX_CLK,
20     inout  wire      HPS_ENET_INT_N,
21     output wire      HPS_ENET_MDC,
22     inout  wire      HPS_ENET_MDIO,
23     input  wire      HPS_ENET_RX_CLK,
24     input  wire [3:0] HPS_ENET_RX_DATA,
25     input  wire      HPS_ENET_RX_DV,
26     output wire [3:0] HPS_ENET_TX_DATA,
27     output wire      HPS_ENET_TX_EN,
28     output wire      HPS_SPIM_CLK,
29     input  wire      HPS_SPIM_MISO,
30     output wire      HPS_SPIM_MOSI,
31     inout  wire      HPS_SPIM_SS,
32     input  wire      HPS_USB_CLKOUT,
33     inout  wire [7:0] HPS_USB_DATA,
34     input  wire      HPS_USB_DIR,
35     input  wire      HPS_USB_NXT,
36     output wire      HPS_USB_STP,
37     inout  wire      HPS_CONV_USB_N,
38     inout  wire      HPS_GSENSOR_INT,
39     inout  wire      HPS_KEY,
40     inout  wire      HPS_LED,
41     inout  wire      HPS_LTC_GPIO,
42
43     output wire [14:0] HPS_DDR3_ADDR,
44     output wire [2:0]  HPS_DDR3_BA,
45     output wire      HPS_DDR3_CAS_N,
46     output wire      HPS_DDR3_CKE,
47     output wire      HPS_DDR3_CK_N,
48     output wire      HPS_DDR3_CK_P,
49     output wire      HPS_DDR3_CS_N,
50     output wire [3:0] HPS_DDR3_DM,
51     inout  wire [31:0] HPS_DDR3_DQ,
52     inout  wire [3:0] HPS_DDR3_DQS_N,
53     inout  wire [3:0] HPS_DDR3_DQS_P,
54     output wire      HPS_DDR3_ODT,
55     output wire      HPS_DDR3_RAS_N,
56     output wire      HPS_DDR3_RESET_N,
57     input  wire      HPS_DDR3_RZQ,
58     output wire      HPS_DDR3_WE_N
59 );
60 wire [3:0] predict_class_dbg;
61 wire      predict_done_dbg;
62 wire      model_loaded_dbg;
63 wire [15:0] interface_error_dbg;
64
65 localparam [6:0] SEVEN_SEG_BLANK = 7'b1111111;
66
67 function automatic [6:0] seven_seg_decode;
68     input [3:0] value;
69     begin
70         case (value)
71             4'h0: seven_seg_decode = 7'b1000000;
72             4'h1: seven_seg_decode = 7'b1111001;

```

```

73         4'h2: seven_seg_decode = 7'b0100100;
74         4'h3: seven_seg_decode = 7'b0110000;
75         4'h4: seven_seg_decode = 7'b0011001;
76         4'h5: seven_seg_decode = 7'b0010010;
77         4'h6: seven_seg_decode = 7'b0000010;
78         4'h7: seven_seg_decode = 7'b1111000;
79         4'h8: seven_seg_decode = 7'b0000000;
80         4'h9: seven_seg_decode = 7'b0010000;
81         4'hA: seven_seg_decode = 7'b0001000; // A
82         4'hB: seven_seg_decode = 7'b0000011; // b
83         4'hC: seven_seg_decode = 7'b1000110; // C
84         4'hD: seven_seg_decode = 7'b0100001; // d
85         4'hE: seven_seg_decode = 7'b0000110; // E
86         default: seven_seg_decode = SEVEN_SEG_BLANK;
87     endcase
88 end
89 endfunction
90
91 soc_system u_soc_system (
92     .clk_clk          (CLOCK_50),
93     .cnn_debug_model_loaded(model_loaded_dbg),
94     .cnn_debug_predict_done(predict_done_dbg),
95     .cnn_debug_predict_class(predict_class_dbg),
96     .cnn_debug_interface_error(interface_error_dbg),
97     .reset_reset_n    (1'b1),
98     .memory_mem_a     (HPS_DDR3_ADDR),
99     .memory_mem_ba    (HPS_DDR3_BA),
100    .memory_mem_ck     (HPS_DDR3_CK_P),
101    .memory_mem_ck_n   (HPS_DDR3_CK_N),
102    .memory_mem_cke    (HPS_DDR3_CKE),
103    .memory_mem_cs_n   (HPS_DDR3_CS_N),
104    .memory_mem_ras_n  (HPS_DDR3_RAS_N),
105    .memory_mem_cas_n  (HPS_DDR3_CAS_N),
106    .memory_mem_we_n   (HPS_DDR3_WE_N),
107    .memory_mem_reset_n (HPS_DDR3_RESET_N),
108    .memory_mem_dq     (HPS_DDR3_DQ),
109    .memory_mem_dqs    (HPS_DDR3_DQS_P),
110    .memory_mem_dqs_n  (HPS_DDR3_DQS_N),
111    .memory_mem_odt    (HPS_DDR3_ODT),
112    .memory_mem_dm     (HPS_DDR3_DM),
113    .memory_oct_rzqin  (HPS_DDR3_RZQ),
114    .hps_hps_io_emac1_inst_TX_CLK (HPS_ENET_GTX_CLK),
115    .hps_hps_io_emac1_inst_TXD0  (HPS_ENET_TX_DATA[0]),
116    .hps_hps_io_emac1_inst_TXD1  (HPS_ENET_TX_DATA[1]),
117    .hps_hps_io_emac1_inst_TXD2  (HPS_ENET_TX_DATA[2]),
118    .hps_hps_io_emac1_inst_TXD3  (HPS_ENET_TX_DATA[3]),
119    .hps_hps_io_emac1_inst_RXD0  (HPS_ENET_RX_DATA[0]),
120    .hps_hps_io_emac1_inst_MDIO  (HPS_ENET_MDIO),
121    .hps_hps_io_emac1_inst_MDC   (HPS_ENET_MDC),
122    .hps_hps_io_emac1_inst_RX_CTL (HPS_ENET_RX_DV),
123    .hps_hps_io_emac1_inst_TX_CTL (HPS_ENET_TX_EN),
124    .hps_hps_io_emac1_inst_RX_CLK (HPS_ENET_RX_CLK),
125    .hps_hps_io_emac1_inst_RXD1  (HPS_ENET_RX_DATA[1]),
126    .hps_hps_io_emac1_inst_RXD2  (HPS_ENET_RX_DATA[2]),

```

```

127     .hps_hps_io_emac1_inst_RXD3    (HPS_ENET_RX_DATA[3]),
128     .hps_hps_io_sdio_inst_CMD      (HPS_SD_CMD),
129     .hps_hps_io_sdio_inst_D0       (HPS_SD_DATA[0]),
130     .hps_hps_io_sdio_inst_D1       (HPS_SD_DATA[1]),
131     .hps_hps_io_sdio_inst_CLK      (HPS_SD_CLK),
132     .hps_hps_io_sdio_inst_D2       (HPS_SD_DATA[2]),
133     .hps_hps_io_sdio_inst_D3       (HPS_SD_DATA[3]),
134     .hps_hps_io_usb1_inst_D0       (HPS_USB_DATA[0]),
135     .hps_hps_io_usb1_inst_D1       (HPS_USB_DATA[1]),
136     .hps_hps_io_usb1_inst_D2       (HPS_USB_DATA[2]),
137     .hps_hps_io_usb1_inst_D3       (HPS_USB_DATA[3]),
138     .hps_hps_io_usb1_inst_D4       (HPS_USB_DATA[4]),
139     .hps_hps_io_usb1_inst_D5       (HPS_USB_DATA[5]),
140     .hps_hps_io_usb1_inst_D6       (HPS_USB_DATA[6]),
141     .hps_hps_io_usb1_inst_D7       (HPS_USB_DATA[7]),
142     .hps_hps_io_usb1_inst_CLK      (HPS_USB_CLKOUT),
143     .hps_hps_io_usb1_inst_STP      (HPS_USB_STP),
144     .hps_hps_io_usb1_inst_DIR      (HPS_USB_DIR),
145     .hps_hps_io_usb1_inst_NXT      (HPS_USB_NXT),
146     .hps_hps_io_spim1_inst_CLK     (HPS_SPIM_CLK),
147     .hps_hps_io_spim1_inst_MOSI    (HPS_SPIM_MOSI),
148     .hps_hps_io_spim1_inst_MISO    (HPS_SPIM_MISO),
149     .hps_hps_io_spim1_inst_SS0     (HPS_SPIM_SS),
150     .hps_hps_io_uart0_inst_RX      (HPS_UART_RX),
151     .hps_hps_io_uart0_inst_TX      (HPS_UART_TX),
152     .hps_hps_io_i2c0_inst_SDA      (HPS_I2C1_SDAT),
153     .hps_hps_io_i2c0_inst_SCL      (HPS_I2C1_SCLK),
154     .hps_hps_io_i2c1_inst_SDA      (HPS_I2C2_SDAT),
155     .hps_hps_io_i2c1_inst_SCL      (HPS_I2C2_SCLK),
156     .hps_hps_io_gpio_inst_GPI009   (HPS_CONV_USB_N),
157     .hps_hps_io_gpio_inst_GPI035   (HPS_ENET_INT_N),
158     .hps_hps_io_gpio_inst_GPI040   (HPS_LTC_GPIO),
159     .hps_hps_io_gpio_inst_GPI048   (HPS_I2C_CONTROL),
160     .hps_hps_io_gpio_inst_GPI053   (HPS_LED),
161     .hps_hps_io_gpio_inst_GPI054   (HPS_KEY),
162     .hps_hps_io_gpio_inst_GPI061   (HPS_GSENSOR_INT)
163 );
164
165 // Current seven-segment policy:
166 // - HEX0 shows the predicted class once predict_done is asserted
167 // - HEX1 shows model_loaded / predict_done as two debug bits
168 // - HEX2 shows the low nibble of interface_error when non-zero
169 // - HEX3/HEX4 stay blank
170 // - HEX5 shows "E" only when an interface error is latched
171 assign HEX0 = predict_done_dbg ? seven_seg_decode(predict_class_dbg) : SEVEN_SEG_BLANK;
172 assign HEX1 = seven_seg_decode({2'b00, model_loaded_dbg, predict_done_dbg});
173 assign HEX2 = (interface_error_dbg != 16'h0000) ? seven_seg_decode(interface_error_dbg[3:0]) :
    SEVEN_SEG_BLANK;
174 assign HEX3 = SEVEN_SEG_BLANK;
175 assign HEX4 = SEVEN_SEG_BLANK;
176 assign HEX5 = (interface_error_dbg != 16'h0000) ? seven_seg_decode(4'hE) : SEVEN_SEG_BLANK;
177
178 endmodule

```

## A.7.2 Avalon-MM MMIO Wrapper

*Implements the scratchpad, register file, replay FSM, status/profile readback, and stream bridge into the accelerator.*

```
code/input/RTL/interface/cnn_mmio_interface.v

1  `timescale 1ns / 1ps
2
3  // cnn_mmio_interface
4  // -----
5  // A thin MMIO/BRAM wrapper around system_top.
6  //
7  // address[12] = 0 : memory space   (16 KiB 32-bit word BRAM)
8  // address[12] = 1 : config space   (32-bit MMIO registers, low 16 bits keep
9  //                                     the original control/status ABI)
10 //
11 // The wrapper keeps the existing accelerator datapath intact and only adds:
12 //   1) a stable memory image interface for model/image staging
13 //   2) a small register block for base/length programming
14 //   3) a sequencer that replays staged 32-bit words into system_top
15 //
16 // This keeps the current CNN_ACC architecture and borrows only the transaction
17 // pattern from the reference cnn_interface design.
18 //
19 // Notes on the slightly "mixed-width" looking register map:
20 //   - Control/status bits stay in the low 16 bits so the existing HPS helpers
21 //     and docs keep working after the move to a 32-bit Avalon-MM slave.
22 //   - Profile counters are already 32-bit values, so the legacy *_HI slots are
23 //     treated as the canonical full-width readback locations. The matching
24 //     *_LO slots are kept as zero/reserved padding for compatibility.
25
26 module cnn_mmio_interface #(
27     parameter integer MEM_AW      = 12,
28     parameter integer MEM_WORDS   = (1 << MEM_AW)
29 ) (
30     input  wire      clk,
31     input  wire      reset,
32     input  wire [31:0] writedata,
33     input  wire [3:0]  byteenable,
34     input  wire      write,
35     input  wire      read,
36     input  wire      chipselect,
37     input  wire [MEM_AW:0] address,
38     output reg [31:0] readdata,
39     output wire      debug_model_loaded,
40     output wire      debug_predict_done,
41     output wire [3:0] debug_predict_class,
42     output wire [15:0] debug_interface_error
43 );
44
45 localparam [4:0] REG_CONTROL      = 5'd0;
46 localparam [4:0] REG_STATUS      = 5'd1;
```

```

47 localparam [4:0] REG_CONV_CFG_BASE = 5'd2;
48 localparam [4:0] REG_CONV_CFG_LEN = 5'd3;
49 localparam [4:0] REG_CONV_WT_BASE = 5'd4;
50 localparam [4:0] REG_CONV_WT_LEN = 5'd5;
51 localparam [4:0] REG_FC_BIAS_BASE = 5'd6;
52 localparam [4:0] REG_FC_BIAS_LEN = 5'd7;
53 localparam [4:0] REG_FCW_BASE = 5'd8;
54 localparam [4:0] REG_FCW_LEN = 5'd9;
55 localparam [4:0] REG_IMAGE_BASE = 5'd10;
56 localparam [4:0] REG_IMAGE_LEN = 5'd11;
57 localparam [4:0] REG_PREDICT = 5'd12;
58 localparam [4:0] REG_IF_ERROR = 5'd13;
59 localparam [4:0] REG_PROFILE_L1_LO = 5'd14;
60 localparam [4:0] REG_PROFILE_L1_HI = 5'd15;
61 localparam [4:0] REG_PROFILE_L2_PO_LO = 5'd16;
62 localparam [4:0] REG_PROFILE_L2_PO_HI = 5'd17;
63 localparam [4:0] REG_PROFILE_L2_P1_LO = 5'd18;
64 localparam [4:0] REG_PROFILE_L2_P1_HI = 5'd19;
65 localparam [4:0] REG_PROFILE_L3_PO_LO = 5'd20;
66 localparam [4:0] REG_PROFILE_L3_PO_HI = 5'd21;
67 localparam [4:0] REG_PROFILE_L3_P1_LO = 5'd22;
68 localparam [4:0] REG_PROFILE_L3_P1_HI = 5'd23;
69 localparam [4:0] REG_PROFILE_FC_LO = 5'd24;
70 localparam [4:0] REG_PROFILE_FC_HI = 5'd25;
71 localparam [4:0] REG_PROFILE_ARGMAX_LO = 5'd26;
72 localparam [4:0] REG_PROFILE_ARGMAX_HI = 5'd27;
73 localparam [4:0] REG_PROFILE_TOTAL_LO = 5'd28;
74 localparam [4:0] REG_PROFILE_TOTAL_HI = 5'd29;
75 localparam [4:0] REG_LAST_WRITE = 5'd30;
76 localparam [4:0] REG_MAGIC = 5'd31;
77
78 localparam [31:0] CFG_MAGIC_WORD = 32'h434E4E32; // "CNN2"
79
80 localparam [1:0] ENG_IDLE = 2'd0;
81 localparam [1:0] ENG_MODEL = 2'd1;
82 localparam [1:0] ENG_INFER = 2'd2;
83
84 localparam [2:0] ST_IDLE = 3'd0;
85 localparam [2:0] ST_RD = 3'd1;
86 localparam [2:0] ST_WAIT_RD = 3'd2;
87 localparam [2:0] ST_SEND = 3'd3;
88 localparam [2:0] ST_GAP = 3'd4;
89
90 localparam [1:0] SEG_CONV_CFG = 2'd0;
91 localparam [1:0] SEG_CONV_WT = 2'd1;
92 localparam [1:0] SEG_FC_BIAS = 2'd2;
93 localparam [1:0] SEG_FCW = 2'd3;
94
95 // -----
96 // 32-bit staging BRAM
97 // -----
98 (* ramstyle = "M10K" *)
99 reg [31:0] mem [0:MEM_WORDS-1];
100

```

```

101 reg [31:0] mem_rdata_q;
102 reg [31:0] host_rdata_q;
103 reg [31:0] last_writedata_q;
104 reg [3:0] last_byteenable_q;
105 reg host_mem_rd_q;
106 reg [MEM_AW-1:0] mem_rd_addr_q;
107 reg mem_rd_pending_q;
108 reg mem_rd_host_q;
109
110 wire mem_sel = chipselect && !address[MEM_AW];
111 wire cfg_sel = chipselect && address[MEM_AW];
112 wire mem_wr_req = mem_sel && write;
113 wire mem_rd_req = mem_sel && read;
114 wire cfg_wr_req = cfg_sel && write;
115 wire cfg_rd_req = cfg_sel && read;
116 wire mem_addr_ok = (address[MEM_AW-1:0] < MEM_WORDS);
117 wire mem_any_byte_en = |byteenable;
118 wire cfg_lo_byte_en = |byteenable[1:0];
119
120 // -----
121 // Config registers
122 // -----
123 reg [15:0] conv_cfg_base_w;
124 reg [15:0] conv_cfg_len_w;
125 reg [15:0] conv_wt_base_w;
126 reg [15:0] conv_wt_len_w;
127 reg [15:0] fc_bias_base_w;
128 reg [15:0] fc_bias_len_w;
129 reg [15:0] fcw_base_w;
130 reg [15:0] fcw_len_w;
131 reg [15:0] image_base_w;
132 reg [15:0] image_len_w;
133
134 reg model_loaded;
135 reg predict_done;
136 reg [3:0] predict_class_latched;
137 reg [15:0] interface_error;
138 reg [1:0] eng_mode;
139 wire eng_busy = (eng_mode != ENG_IDLE);
140
141 wire [15:0] cfg_status_word = {
142     8'd0,
143     predict_class_latched,
144     predict_done,
145     model_loaded,
146     eng_busy,
147     1'b0
148 };
149
150 reg [31:0] cfg_read_data;
151 always @(*) begin
152     cfg_read_data = 32'd0;
153     case (address[4:0])
154         REG_STATUS:     cfg_read_data = cfg_status_word;

```

```

155     REG_CONV_CFG_BASE: cfg_read_data = conv_cfg_base_w;
156     REG_CONV_CFG_LEN:  cfg_read_data = conv_cfg_len_w;
157     REG_CONV_WT_BASE:  cfg_read_data = conv_wt_base_w;
158     REG_CONV_WT_LEN:   cfg_read_data = conv_wt_len_w;
159     REG_FC_BIAS_BASE:  cfg_read_data = fc_bias_base_w;
160     REG_FC_BIAS_LEN:   cfg_read_data = fc_bias_len_w;
161     REG_FCW_BASE:      cfg_read_data = fcw_base_w;
162     REG_FCW_LEN:       cfg_read_data = fcw_len_w;
163     REG_IMAGE_BASE:    cfg_read_data = image_base_w;
164     REG_IMAGE_LEN:     cfg_read_data = image_len_w;
165     REG_PREDICT:        cfg_read_data = {12'd0, predict_class_latched};
166     REG_IF_ERROR:      cfg_read_data = interface_error;
167     // Keep the legacy low-half register numbers reserved. Software should
168     // read the *_HI aliases for the full 32-bit counters.
169     REG_PROFILE_L1_LO:  cfg_read_data = 32'd0;
170     REG_PROFILE_L1_HI:  cfg_read_data = accel_profile_l1_cycles;
171     REG_PROFILE_L2_P0_LO: cfg_read_data = 32'd0;
172     REG_PROFILE_L2_P0_HI: cfg_read_data = accel_profile_l2_p0_cycles;
173     REG_PROFILE_L2_P1_LO: cfg_read_data = 32'd0;
174     REG_PROFILE_L2_P1_HI: cfg_read_data = accel_profile_l2_p1_cycles;
175     REG_PROFILE_L3_P0_LO: cfg_read_data = 32'd0;
176     REG_PROFILE_L3_P0_HI: cfg_read_data = accel_profile_l3_p0_cycles;
177     REG_PROFILE_L3_P1_LO: cfg_read_data = 32'd0;
178     REG_PROFILE_L3_P1_HI: cfg_read_data = accel_profile_l3_p1_cycles;
179     REG_PROFILE_FC_LO:   cfg_read_data = 32'd0;
180     REG_PROFILE_FC_HI:   cfg_read_data = accel_profile_fc_cycles;
181     REG_PROFILE_ARGMAX_LO: cfg_read_data = 32'd0;
182     REG_PROFILE_ARGMAX_HI: cfg_read_data = accel_profile_argmax_cycles;
183     REG_PROFILE_TOTAL_LO:  cfg_read_data = 32'd0;
184     REG_PROFILE_TOTAL_HI:  cfg_read_data = accel_profile_total_cycles;
185     REG_LAST_WRITE:        cfg_read_data = { last_byteenable_q, last_writedata_q[27:0] };
186     REG_MAGIC:             cfg_read_data = CFG_MAGIC_WORD;
187     default:               cfg_read_data = 32'd0;
188 endcase
189 end
190
191 // -----
192 // system_top host-stream interface
193 // -----
194 reg      accel_load_sel;
195 reg      accel_load_valid;
196 reg [31:0] accel_load_data;
197 reg      accel_load_last;
198 wire     accel_load_ready;
199 wire     accel_busy;
200 wire     accel_predict_valid;
201 wire [3:0] accel_predict_class;
202 wire [31:0] accel_profile_l1_cycles;
203 wire [31:0] accel_profile_l2_p0_cycles;
204 wire [31:0] accel_profile_l2_p1_cycles;
205 wire [31:0] accel_profile_l3_p0_cycles;
206 wire [31:0] accel_profile_l3_p1_cycles;
207 wire [31:0] accel_profile_fc_cycles;
208 wire [31:0] accel_profile_argmax_cycles;

```

```

209 wire [31:0] accel_profile_total_cycles;
210
211 system_top u_system_top (
212     .clk(clk),
213     .rst_n(~reset),
214     .load_sel(accel_load_sel),
215     .load_valid(accel_load_valid),
216     .load_data(accel_load_data),
217     .load_last(accel_load_last),
218     .load_ready(accel_load_ready),
219     .busy(accel_busy),
220     .predict_valid(accel_predict_valid),
221     .predict_class(accel_predict_class),
222     .profile_l1_cycles(accel_profile_l1_cycles),
223     .profile_l2_p0_cycles(accel_profile_l2_p0_cycles),
224     .profile_l2_p1_cycles(accel_profile_l2_p1_cycles),
225     .profile_l3_p0_cycles(accel_profile_l3_p0_cycles),
226     .profile_l3_p1_cycles(accel_profile_l3_p1_cycles),
227     .profile_fc_cycles(accel_profile_fc_cycles),
228     .profile_argmax_cycles(accel_profile_argmax_cycles),
229     .profile_total_cycles(accel_profile_total_cycles)
230 );
231
232 // -----
233 // Replay engine
234 // -----
235 reg [2:0] eng_state;
236 reg [1:0] model_seg_idx;
237 reg [15:0] seg_word_idx;
238 reg [15:0] seg_word_count;
239
240 function [15:0] model_seg_base;
241     input [1:0] seg;
242     begin
243         case (seg)
244             SEG_CONV_CFG: model_seg_base = conv_cfg_base_w;
245             SEG_CONV_WT:  model_seg_base = conv_wt_base_w;
246             SEG_FC_BIAS:  model_seg_base = fc_bias_base_w;
247             default:     model_seg_base = fcw_base_w;
248         endcase
249     end
250 endfunction
251
252 function [15:0] model_seg_len;
253     input [1:0] seg;
254     begin
255         case (seg)
256             SEG_CONV_CFG: model_seg_len = conv_cfg_len_w;
257             SEG_CONV_WT:  model_seg_len = conv_wt_len_w;
258             SEG_FC_BIAS:  model_seg_len = fc_bias_len_w;
259             default:     model_seg_len = fcw_len_w;
260         endcase
261     end
262 endfunction

```

```

263
264 // Memory read arbitration: the engine owns BRAM while active.
265 wire issue_eng_read = (eng_state == ST_RD);
266 wire [15:0] current_seg_base = (eng_mode == ENG_INFER) ? image_base_w
267                                     : model_seg_base(model_seg_idx);
268 wire [15:0] current_seg_len  = (eng_mode == ENG_INFER) ? image_len_w
269                                     : model_seg_len(model_seg_idx);
270 wire [15:0] current_word_addr = current_seg_base + seg_word_idx;
271 wire eng_read_in_range = (current_word_addr < MEM_WORDS);
272
273 always @(posedge clk) begin
274     if (reset) begin
275         mem_rdata_q    <= 32'd0;
276         host_rdata_q   <= 32'd0;
277         last_writedata_q <= 32'd0;
278         last_byteenable_q <= 4'd0;
279         host_mem_rd_q   <= 1'b0;
280         mem_rd_addr_q   <= {MEM_AW{1'b0}};
281         mem_rd_pending_q <= 1'b0;
282         mem_rd_host_q   <= 1'b0;
283     end else begin
284         host_mem_rd_q <= 1'b0;
285
286         if (cfg_rd_req) begin
287             host_rdata_q <= cfg_read_data;
288         end
289
290         if (host_mem_rd_q)
291             host_rdata_q <= mem_rdata_q;
292
293         if (mem_wr_req || cfg_wr_req) begin
294             last_writedata_q <= writedata;
295             last_byteenable_q <= byteenable;
296         end
297
298         // The staged scratchpad expects full-word host writes. We only use
299         // byteenable as a write-present qualifier here to keep the memory inferable
300         // as block RAM on Cyclone V.
301         if (mem_wr_req && !eng_busy && mem_addr_ok && mem_any_byte_en)
302             mem[address[MEM_AW-1:0]] <= writedata;
303
304         if (mem_rd_pending_q) begin
305             mem_rdata_q <= mem[mem_rd_addr_q];
306             host_mem_rd_q <= mem_rd_host_q;
307         end
308
309         mem_rd_pending_q <= 1'b0;
310
311         if (issue_eng_read && eng_read_in_range) begin
312             mem_rd_addr_q    <= current_word_addr[MEM_AW-1:0];
313             mem_rd_pending_q <= 1'b1;
314             mem_rd_host_q    <= 1'b0;
315         end else if (mem_rd_req && !eng_busy && mem_addr_ok) begin
316             mem_rd_addr_q    <= address[MEM_AW-1:0];

```

```

317     mem_rd_pending_q <= 1'b1;
318     mem_rd_host_q    <= 1'b1;
319     end
320 end
321 end
322
323 // -----
324 // Register block, status latches, and replay FSM
325 // -----
326 always @(posedge clk) begin
327     if (reset) begin
328         conv_cfg_base_w    <= 16'd0;
329         conv_cfg_len_w     <= 16'd45;
330         conv_wt_base_w     <= 16'd45;
331         conv_wt_len_w      <= 16'd225;
332         fc_bias_base_w     <= 16'd270;
333         fc_bias_len_w      <= 16'd10;
334         fcw_base_w         <= 16'd280;
335         fcw_len_w          <= 16'd864;
336         image_base_w       <= 16'd1144;
337         image_len_w        <= 16'd1024;
338         model_loaded       <= 1'b0;
339         predict_done       <= 1'b0;
340         predict_class_latched <= 4'd0;
341         interface_error    <= 16'd0;
342         eng_mode           <= ENG_IDLE;
343         eng_state          <= ST_IDLE;
344         model_seg_idx      <= SEG_CONV_CFG;
345         seg_word_idx       <= 16'd0;
346         seg_word_count     <= 16'd0;
347     end else begin
348         if (accel_predict_valid) begin
349             predict_done <= 1'b1;
350             predict_class_latched <= accel_predict_class;
351         end
352
353         if (cfg_wr_req && cfg_lo_byte_en) begin
354             case (address[4:0])
355                 REG_CONTROL: begin
356                     if (writedata[2]) begin
357                         predict_done <= 1'b0;
358                         interface_error <= 16'd0;
359                     end
360
361                     if (writedata[0]) begin
362                         if (eng_busy || accel_busy) begin
363                             interface_error[0] <= 1'b1;
364                         end else begin
365                             eng_mode <= ENG_MODEL;
366                             eng_state <= ST_RD;
367                             model_seg_idx <= SEG_CONV_CFG;
368                             seg_word_idx <= 16'd0;
369                             seg_word_count <= conv_cfg_len_w;
370                             predict_done <= 1'b0;

```

```

371         model_loaded    <= 1'b0;
372     end
373 end
374
375     if (writedata[1]) begin
376         if (eng_busy || accel_busy) begin
377             interface_error[1] <= 1'b1;
378         end else if (!model_loaded) begin
379             interface_error[2] <= 1'b1;
380         end else begin
381             eng_mode          <= ENG_INFER;
382             eng_state         <= ST_RD;
383             seg_word_idx      <= 16'd0;
384             seg_word_count    <= image_len_w;
385             predict_done      <= 1'b0;
386         end
387     end
388 end
389
390     REG_CONV_CFG_BASE: conv_cfg_base_w <= writedata[15:0];
391     REG_CONV_CFG_LEN:  conv_cfg_len_w  <= writedata[15:0];
392     REG_CONV_WT_BASE:  conv_wt_base_w  <= writedata[15:0];
393     REG_CONV_WT_LEN:   conv_wt_len_w   <= writedata[15:0];
394     REG_FC_BIAS_BASE:  fc_bias_base_w  <= writedata[15:0];
395     REG_FC_BIAS_LEN:   fc_bias_len_w   <= writedata[15:0];
396     REG_FCW_BASE:      fcw_base_w      <= writedata[15:0];
397     REG_FCW_LEN:       fcw_len_w       <= writedata[15:0];
398     REG_IMAGE_BASE:    image_base_w     <= writedata[15:0];
399     REG_IMAGE_LEN:     image_len_w     <= writedata[15:0];
400     default: begin
401         end
402     endcase
403 end
404
405 case (eng_state)
406     ST_IDLE: begin
407         end
408
409     ST_RD: begin
410         if (!eng_read_in_range) begin
411             interface_error[3] <= 1'b1;
412             eng_mode           <= ENG_IDLE;
413             eng_state          <= ST_IDLE;
414         end else begin
415             eng_state <= ST_WAIT_RD;
416         end
417     end
418
419     ST_WAIT_RD: begin
420         eng_state <= ST_SEND;
421     end
422
423     ST_SEND: begin
424         if (accel_load_ready) begin

```

```

425     if (seg_word_idx == (seg_word_count - 16'd1)) begin
426         if (eng_mode == ENG_MODEL) begin
427             if (model_seg_idx == SEG_FCW) begin
428                 model_loaded <= 1'b1;
429                 eng_mode     <= ENG_IDLE;
430                 eng_state    <= ST_IDLE;
431             end else begin
432                 model_seg_idx <= model_seg_idx + 2'd1;
433                 seg_word_idx  <= 16'd0;
434                 eng_state    <= ST_GAP;
435             end
436         end else begin
437             eng_mode     <= ENG_IDLE;
438             eng_state    <= ST_IDLE;
439             seg_word_idx <= 16'd0;
440         end
441     end else begin
442         seg_word_idx <= seg_word_idx + 16'd1;
443         eng_state    <= ST_RD;
444     end
445 end
446 end
447
448 ST_GAP: begin
449     seg_word_count <= model_seg_len(model_seg_idx);
450     eng_state     <= ST_RD;
451 end
452
453 default: begin
454     eng_state <= ST_IDLE;
455     eng_mode  <= ENG_IDLE;
456 end
457 endcase
458 end
459 end
460
461 always @(*) begin
462     accel_load_valid = (eng_state == ST_SEND);
463     accel_load_data  = mem_rdata_q;
464     accel_load_sel   = (eng_mode == ENG_INFER);
465     accel_load_last  = (eng_state == ST_SEND) && (seg_word_idx == (seg_word_count - 16'd1));
466 end
467
468 // -----
469 // Read response
470 // -----
471 always @(*) begin
472     if (host_mem_rd_q)
473         readdata = mem_rdata_q;
474     else
475         readdata = host_rdata_q;
476 end
477
478 assign debug_model_loaded = model_loaded;

```

```

479 assign debug_predict_done = predict_done;
480 assign debug_predict_class = predict_class_latched;
481 assign debug_interface_error = interface_error;
482
483 endmodule

```

### A.7.3 Accelerator Top

*Integrates model-load, CNN layer sequencing, SRAM datapath, FC classifier, argmax, and profiling outputs.*

code/input/RTL/system\_top.v

```

1 // system_top: ASIC top-level module
2 //
3 // Instantiates:
4 // TopFSM + LayerRunnerFSM (control)
5 // top_sram_A + top_sram_B (storage, from teammate)
6 // conv_quant_pool (Conv-Quant-Pool)
7 // FC + fc_data_adapter (FC classifier)
8 // wt_prepad_inserter (3-zero-beat inserter for conv weights)
9 //
10 // Chip pins: clk, rst_n, load_data[31:0], load_valid, load_sel,
11 // load_last, load_ready, busy, predict_valid, predict_class[1:0]
12
13 module system_top #(
14     parameter integer OUT_CHANNELS = 10
15 ) (
16     input wire clk,
17     input wire rst_n,
18
19     // --- External data port ---
20     input wire load_sel,
21     input wire load_valid,
22     input wire [31:0] load_data,
23     input wire load_last,
24     output wire load_ready,
25
26     // --- Status ---
27     output wire busy,
28     output wire predict_valid,
29     output wire [3:0] predict_class,
30     output wire [31:0] profile_l1_cycles,
31     output wire [31:0] profile_l2_p0_cycles,
32     output wire [31:0] profile_l2_p1_cycles,
33     output wire [31:0] profile_l3_p0_cycles,
34     output wire [31:0] profile_l3_p1_cycles,
35     output wire [31:0] profile_fc_cycles,
36     output wire [31:0] profile_argmax_cycles,
37     output wire [31:0] profile_total_cycles

```

```

38 );
39
40 // =====
41 // Internal wires
42 // =====
43
44 // --- TopFSM <-> LayerRunnerFSM ---
45 wire runner_start, runner_done;
46 wire [1:0] runner_layer_sel;
47 wire runner_pass_id, runner_is_fc;
48
49 // --- TopFSM -> SRAM_A preload control ---
50 wire top_sram_a_start;
51 wire [2:0] top_sram_a_layer_sel;
52 wire [1:0] top_sram_a_data_sel;
53 wire top_sram_a_done;
54 wire preload_wr_valid;
55 wire [31:0] preload_wr_data;
56
57 // --- LayerRunnerFSM -> SRAM_A inference control ---
58 wire run_sram_a_start;
59 wire [2:0] run_sram_a_layer_sel;
60 wire [1:0] run_sram_a_data_sel;
61 wire run_sram_a_pass_id;
62
63 // --- LayerRunnerFSM -> SRAM_B control ---
64 wire run_sram_b_start;
65 wire [2:0] run_sram_b_layer_sel;
66 wire [1:0] run_sram_b_data_sel;
67 wire run_sram_b_pass_id;
68
69 // --- SRAM_A muxed control ---
70 wire sram_a_start;
71 wire [2:0] sram_a_layer_sel;
72 wire [1:0] sram_a_data_sel;
73 wire sram_a_pass_id;
74 wire sram_a_done;
75
76 // --- SRAM_A streams ---
77 wire sram_a_data_valid, sram_a_data_last;
78 wire [31:0] sram_a_read_data;
79 wire sram_a_data_ready;
80 wire sram_a_pool_ready;
81
82 // --- SRAM_B streams ---
83 wire sram_b_data_valid, sram_b_data_last;
84 wire [31:0] sram_b_read_data;
85 wire sram_b_data_ready;
86 wire sram_b_pool_ready;
87 wire sram_b_done;
88
89 // --- conv_quant_pool signals ---
90 wire conv_cfg_ready, conv_wt_ready, conv_in_ready;
91 wire signed [31:0] conv_pool_data;

```

```

92 wire      conv_pool_valid, conv_pool_last;
93 wire      conv_cfg_load_done, conv_wt_load_done, conv_pool_frame_done;
94 wire      conv_frame_rearm_done;
95
96 // --- wt_prepad_inserter ---
97 wire      wt_prepad_dn_valid, wt_prepad_dn_last;
98 wire [31:0] wt_prepad_dn_data;
99 wire      wt_prepad_dn_ready;
100 wire     wt_prepad_up_ready;
101
102 // --- conv_data_adapter ---
103 wire     conv_in_adapter_up_ready;
104
105 // --- L1 pixel bypass control ---
106 wire     pixel_stream_active;
107
108 // --- FC signals ---
109 wire     fc_cfg_ready, fc_param_load_done;
110 wire [OUT_CHANNELS*32-1:0] fc_acc_vec;
111 wire [OUT_CHANNELS-1:0] fc_mul_done_vec;
112 wire     fc_all_done;
113
114 // --- fc_data_adapter signals ---
115 wire     fc_mul_en;
116 wire [OUT_CHANNELS*8-1:0] fc_pixel_vec;
117 wire [OUT_CHANNELS*8-1:0] fc_kernel_vec;
118 wire     fc_adapter_wt_ready, fc_adapter_data_ready;
119 wire     fc_done;
120
121 // --- FC weight dedicated 80-bit stream (from SRAM_FCW via top_sram_A) ---
122 wire [OUT_CHANNELS*8-1:0] fc_wt_read_data;
123 wire     fc_wt_stream_valid;
124 wire     fc_wt_stream_last;
125
126 // =====
127 // MUX control: who owns SRAM_A?
128 // =====
129 // TopFSM drives SRAM_A during preload, LayerRunnerFSM during inference.
130 // preload_mode latches on start pulse and HOLDS until the next start,
131 // so layer_sel/data_sel/pass_id remain stable throughout the transaction
132 // (wrapper requires stability until done).
133
134 // preload_mode: latches on start pulse, determines MUX routing.
135 // sram_a_start is delayed 1 cycle so that preload_mode is stable
136 // when the SRAM wrapper samples layer_sel / data_sel / pass_id.
137 reg preload_mode;
138 always @(posedge clk or negedge rst_n) begin
139     if (!rst_n)          preload_mode <= 1'b1;
140     else if (top_sram_a_start) preload_mode <= 1'b1;
141     else if (run_sram_a_start) preload_mode <= 1'b0;
142 end
143
144 reg sram_a_start_d;
145 always @(posedge clk or negedge rst_n) begin

```

```

146     if (!rst_n) sram_a_start_d <= 1'b0;
147     else       sram_a_start_d <= top_sram_a_start | run_sram_a_start;
148 end
149
150 assign sram_a_start      = sram_a_start_d;
151 assign sram_a_layer_sel = preload_mode ? top_sram_a_layer_sel : run_sram_a_layer_sel;
152 assign sram_a_data_sel  = preload_mode ? top_sram_a_data_sel  : run_sram_a_data_sel;
153 assign sram_a_pass_id   = preload_mode ? 1'b0                  : run_sram_a_pass_id;
154
155 // =====
156 // SRAM_A write port MUX
157 // =====
158 // Preload: load_data from external host
159 // Inference: pool_data from Conv (L2 writeback to SRAM_A)
160 wire [31:0] sram_a_pool_data = preload_mode ? load_data      : conv_pool_data;
161 wire       sram_a_pool_valid = preload_mode ? (preload_wr_valid && load_valid)
162                                     : conv_pool_valid;
163 wire       sram_a_pool_last  = preload_mode ? load_last      : conv_pool_last;
164
165 // SRAM_A preload_wr_ready = sram_a_pool_ready (when in preload mode)
166 wire       preload_wr_ready = preload_mode ? sram_a_pool_ready : 1'b0;
167
168 // =====
169 // SRAM_A read stream routing
170 // =====
171 // Latched data_sel and is_fc to know where to route read data.
172 reg [1:0] active_data_sel;
173 reg       active_is_fc;
174 always @(posedge clk or negedge rst_n) begin
175     if (!rst_n) begin
176         active_data_sel <= 2'd0;
177         active_is_fc   <= 1'b0;
178     end else if (run_sram_a_start) begin
179         active_data_sel <= run_sram_a_data_sel;
180         active_is_fc   <= runner_is_fc;
181     end
182 end
183
184 localparam [1:0] SEL_CFG = 2'd0, SEL_WT = 2'd1, SEL_DATA = 2'd2;
185
186 // Route SRAM_A read to correct consumer:
187 //   CFG + conv -> conv_quant_pool.cfg
188 //   CFG + FC   -> FC.cfg
189 //   WT  + conv -> wt_prepad_inserter -> conv_quant_pool.wt
190 //   WT  + FC   -> fc_data_adapter.wt (no prepad for FC)
191 //   DATA + conv (L1/L3) -> conv_quant_pool.in
192 //   DATA + FC   -> never (FC data comes from SRAM_B)
193
194 wire route_conv_cfg = !active_is_fc && (active_data_sel == SEL_CFG);
195 wire route_conv_wt  = !active_is_fc && (active_data_sel == SEL_WT);
196 wire route_conv_in  = !active_is_fc && (active_data_sel == SEL_DATA);
197 wire route_fc_cfg   = active_is_fc && (active_data_sel == SEL_CFG);
198 // FC weights no longer flow through the 32-bit SRAM_A read bus; they travel
199 // via the dedicated 80-bit fc_wt_read_data port out of top_sram_A.

```

```

200
201 // Conv cfg port
202 wire      conv_cfg_valid = route_conv_cfg ? sram_a_data_valid : 1'b0;
203 wire [31:0] conv_cfg_data = sram_a_read_data;
204 wire      conv_cfg_last  = route_conv_cfg ? sram_a_data_last  : 1'b0;
205
206 // Conv wt port (through prepad inserter)
207 wire      conv_wt_sram_valid = route_conv_wt ? sram_a_data_valid : 1'b0;
208 wire      conv_wt_sram_last  = route_conv_wt ? sram_a_data_last  : 1'b0;
209
210 // Conv in port (L1/L3 data from SRAM_A; L2 data from SRAM_B handled below)
211 wire      conv_in_from_a_valid = route_conv_in ? sram_a_data_valid : 1'b0;
212 wire      conv_in_from_a_last  = route_conv_in ? sram_a_data_last  : 1'b0;
213
214 // FC cfg port
215 wire      fc_cfg_valid = route_fc_cfg ? sram_a_data_valid : 1'b0;
216 wire [31:0] fc_cfg_data = sram_a_read_data;
217 wire      fc_cfg_last  = route_fc_cfg ? sram_a_data_last  : 1'b0;
218
219 // SRAM_A data_ready: OR of all active consumers on the shared 32-bit bus.
220 assign sram_a_data_ready = preload_mode ? 1'b0 :
221     (route_conv_cfg ? conv_cfg_ready      : 1'b0) |
222     (route_conv_wt ? wt_prepad_up_ready   : 1'b0) |
223     (route_conv_in ? conv_in_adapter_up_ready : 1'b0) |
224     (route_fc_cfg ? fc_cfg_ready          : 1'b0);
225
226 // =====
227 // SRAM_B read stream routing
228 // =====
229 // L2: SRAM_B -> conv in port
230 // FC: SRAM_B -> fc_data_adapter data port
231 // L1/L3 must not see SRAM_B read-side signals while SRAM_B is the pool sink.
232 wire route_conv_in_b = !active_is_fc && (runner_layer_sel == 2'b01);
233 wire route_fc_data_b = active_is_fc;
234
235 wire      conv_in_from_b_valid = route_conv_in_b ? sram_b_data_valid : 1'b0;
236 wire      conv_in_from_b_last  = route_conv_in_b ? sram_b_data_last  : 1'b0;
237
238 wire      fc_data_valid = route_fc_data_b ? sram_b_data_valid : 1'b0;
239 wire      fc_data_last  = route_fc_data_b ? sram_b_data_last  : 1'b0;
240
241 assign sram_b_data_ready = (route_conv_in_b ? conv_in_adapter_up_ready : 1'b0) |
242     (route_fc_data_b ? fc_adapter_data_ready : 1'b0);
243
244 // =====
245 // Conv input MUX: L1 from external load_data (bypass SRAM_A),
246 //                  L2 from SRAM_B, L3 from SRAM_A
247 // =====
248 wire      conv_in_from_ext_valid = pixel_stream_active ? load_valid : 1'b0;
249 wire      conv_in_from_ext_last  = pixel_stream_active ? load_last  : 1'b0;
250
251 wire      conv_in_raw_valid = conv_in_from_ext_valid
252     | conv_in_from_a_valid
253     | conv_in_from_b_valid;

```

```

254 wire [31:0] conv_in_raw_data = conv_in_from_ext_valid ? load_data      :
255                               conv_in_from_a_valid  ? sram_a_read_data :
256                               sram_b_read_data;
257 wire          conv_in_raw_last = conv_in_from_ext_last
258                               | conv_in_from_a_last
259                               | conv_in_from_b_last;
260
261 // conv_data_adapter: L1 byte unpack (1 word -> 4 beats), L2/L3 pass-through
262 wire          conv_in_valid;
263 wire [31:0] conv_in_data;
264 wire [3:0] conv_in_byte_en;
265 wire          conv_in_last;
266
267 // =====
268 // SRAM_B write port: pool output from Conv (L1/L3 writeback)
269 // =====
270 // Only connected during L1/L3 STREAM phase; SRAM_B controller
271 // ignores pool_valid when not in write mode.
272 wire [31:0] sram_b_pool_data = conv_pool_data;
273 wire          sram_b_pool_valid = conv_pool_valid;
274 wire          sram_b_pool_last = conv_pool_last;
275
276 // Conv pool_ready: from whichever SRAM is the write sink
277 // L1 (00) / L3 (10): SRAM_B is sink -> pool_ready from SRAM_B
278 // L2 (01): SRAM_A is sink -> pool_ready from SRAM_A (when not in preload)
279 wire          pool_sink_is_b = (runner_layer_sel != 2'b01); // L1/L3 -> B; L2 -> A
280 wire          conv_pool_ready = preload_mode ? 1'b0 :
281                               (pool_sink_is_b ? sram_b_pool_ready
282                               : sram_a_pool_ready);
283
284 // =====
285 // Module instantiations
286 // =====
287
288 // --- TopFSM ---
289 top_fsm #(.ACC_W(32), .OUT_CHANNELS(OUT_CHANNELS)) u_top_fsm (
290     .clk(clk), .rst_n(rst_n),
291     .load_sel(load_sel), .load_valid(load_valid),
292     .load_data(load_data), .load_last(load_last),
293     .load_ready(load_ready),
294     .busy(busy), .predict_valid(predict_valid), .predict_class(predict_class),
295     .profile_l1_cycles(profile_l1_cycles),
296     .profile_l2_p0_cycles(profile_l2_p0_cycles),
297     .profile_l2_p1_cycles(profile_l2_p1_cycles),
298     .profile_l3_p0_cycles(profile_l3_p0_cycles),
299     .profile_l3_p1_cycles(profile_l3_p1_cycles),
300     .profile_fc_cycles(profile_fc_cycles),
301     .profile_argmax_cycles(profile_argmax_cycles),
302     .profile_total_cycles(profile_total_cycles),
303     .runner_start(runner_start),
304     .runner_layer_sel(runner_layer_sel),
305     .runner_pass_id(runner_pass_id),
306     .runner_is_fc(runner_is_fc),
307     .runner_done(runner_done),

```

```

308     .sram_a_start(top_sram_a_start),
309     .sram_a_layer_sel(top_sram_a_layer_sel),
310     .sram_a_data_sel(top_sram_a_data_sel),
311     .sram_a_done(sram_a_done),
312     .preload_wr_valid(preload_wr_valid),
313     .preload_wr_data(preload_wr_data),
314     .preload_wr_ready(preload_wr_ready),
315     .pixel_stream_active(pixel_stream_active),
316     .conv_adapter_up_ready(conv_in_adapter_up_ready),
317     .conv_frame_rearm(conv_frame_rearm_done),
318     .fc_acc_vec(fc_acc_vec)
319 );
320
321 // --- LayerRunnerFSM ---
322 layer_runner_fsm u_runner (
323     .clk(clk), .rst_n(rst_n),
324     .start(runner_start),
325     .layer_sel(runner_layer_sel),
326     .pass_id(runner_pass_id),
327     .is_fc(runner_is_fc),
328     .busy(), .done(runner_done),
329     .sram_a_start(run_sram_a_start),
330     .sram_a_layer_sel(run_sram_a_layer_sel),
331     .sram_a_data_sel(run_sram_a_data_sel),
332     .sram_a_pass_id(run_sram_a_pass_id),
333     .sram_a_done(sram_a_done),
334     .sram_b_start(run_sram_b_start),
335     .sram_b_layer_sel(run_sram_b_layer_sel),
336     .sram_b_data_sel(run_sram_b_data_sel),
337     .sram_b_pass_id(run_sram_b_pass_id),
338     .sram_b_done(sram_b_done),
339     .cfg_load_done(runner_is_fc ? fc_param_load_done : conv_cfg_load_done),
340     .wt_load_done(conv_wt_load_done),
341     .pool_frame_done(conv_pool_frame_done),
342     .conv_frame_rearm(conv_frame_rearm_done),
343     .fc_done(fc_done)
344 );
345
346 // --- SRAM_A (shared 32-bit pool + SRAM_FCW 80-bit peer for FC weights) ---
347 top_sram_A u_sram_a (
348     .clk(clk), .rst_n(rst_n),
349     .layer_sel(sram_a_layer_sel),
350     .data_sel(sram_a_data_sel),
351     .pass_id(sram_a_pass_id),
352     .start(sram_a_start),
353     .busy(), .done(sram_a_done),
354     .data_ready(sram_a_data_ready),
355     .read_data(sram_a_read_data),
356     .data_valid(sram_a_data_valid),
357     .data_last(sram_a_data_last),
358     .pool_ready(sram_a_pool_ready),
359     .pool_data(sram_a_pool_data),
360     .pool_valid(sram_a_pool_valid),
361     .pool_last(sram_a_pool_last),

```

```

362     .fc_wt_ready(fc_adapter_wt_ready),
363     .fc_wt_read_data(fc_wt_read_data),
364     .fc_wt_valid(fc_wt_stream_valid),
365     .fc_wt_last(fc_wt_stream_last)
366 );
367
368 // --- SRAM_B (from teammate) ---
369 top_sram_B u_sram_b (
370     .clk(clk), .rst_n(rst_n),
371     .layer_sel(run_sram_b_layer_sel),
372     .data_sel(run_sram_b_data_sel),
373     .pass_id(run_sram_b_pass_id),
374     .start(run_sram_b_start),
375     .busy(), .done(sram_b_done),
376     .data_ready(sram_b_data_ready),
377     .read_data(sram_b_read_data),
378     .data_valid(sram_b_data_valid),
379     .data_last(sram_b_data_last),
380     .pool_ready(sram_b_pool_ready),
381     .pool_data(sram_b_pool_data),
382     .pool_valid(sram_b_pool_valid),
383     .pool_last(sram_b_pool_last)
384 );
385
386 // --- WT Prepad Inserter (between SRAM_A read and Conv wt port) ---
387 wt_prepad_inserter u_wt_prepad (
388     .clk(clk), .rst_n(rst_n),
389     .is_wt_read(route_conv_wt),
390     .up_valid(conv_wt_sram_valid),
391     .up_data(sram_a_read_data),
392     .up_last(conv_wt_sram_last),
393     .up_ready(wt_prepad_up_ready),
394     .dn_valid(wt_prepad_dn_valid),
395     .dn_data(wt_prepad_dn_data),
396     .dn_last(wt_prepad_dn_last),
397     .dn_ready(wt_prepad_dn_ready)
398 );
399
400 // --- Conv Data Adapter (L1 byte unpack, L2/L3 pass-through) ---
401 conv_data_adapter u_conv_data_adapter (
402     .clk(clk), .rst_n(rst_n),
403     .layer_sel(runner_layer_sel),
404     .up_valid(conv_in_raw_valid),
405     .up_data(conv_in_raw_data),
406     .up_last(conv_in_raw_last),
407     .up_ready(conv_in_adapter_up_ready),
408     .dn_valid(conv_in_valid),
409     .dn_data(conv_in_data),
410     .dn_byte_en(conv_in_byte_en),
411     .dn_last(conv_in_last),
412     .dn_ready(conv_in_ready)
413 );
414
415 // --- Conv-Quant-Pool ---

```

```

416 conv_quant_pool u_conv (
417     .layer_sel(runner_layer_sel),
418     .clk(clk), .rst_n(rst_n),
419     // Image stream (from conv_data_adapter)
420     .in_valid(conv_in_valid),
421     .in_data(conv_in_data),
422     .in_byte_en(conv_in_byte_en),
423     .in_last(conv_in_last),
424     .in_ready(conv_in_ready),
425     // Weight stream (from prepad inserter)
426     .wt_valid(wt_prepad_dn_valid),
427     .wt_ready(wt_prepad_dn_ready),
428     .wt_last(wt_prepad_dn_last),
429     .wt_data(wt_prepad_dn_data),
430     // Cfg stream
431     .cfg_valid(conv_cfg_valid),
432     .cfg_ready(conv_cfg_ready),
433     .cfg_data(conv_cfg_data),
434     .cfg_last(conv_cfg_last),
435     // Pool output
436     .pool_ready(conv_pool_ready),
437     .pool_data(conv_pool_data),
438     .pool_valid(conv_pool_valid),
439     .pool_last(conv_pool_last),
440     // Control observation
441     .cfg_load_done(conv_cfg_load_done),
442     .wt_load_done(conv_wt_load_done),
443     .pool_frame_done(conv_pool_frame_done),
444     .conv_frame_rearm_out(conv_frame_rearm_done)
445 );
446
447 // --- FC Data Adapter ---
448 fc_data_adapter #(
449     .OUT_CHANNELS(OUT_CHANNELS),
450     .PIX_W(8), .KER_W(8)
451 ) u_fc_adapter (
452     .clk(clk), .rst_n(rst_n),
453     // 80-bit weight stream from SRAM_FCW (via top_sram_A FCW peer)
454     .wt_valid(fc_wt_stream_valid),
455     .wt_data(fc_wt_read_data),
456     .wt_last(fc_wt_stream_last),
457     .wt_ready(fc_adapter_wt_ready),
458     // Data stream from SRAM_B (packed pixels, 32b)
459     .data_valid(fc_data_valid),
460     .data_data(sram_b_read_data),
461     .data_last(fc_data_last),
462     .data_ready(fc_adapter_data_ready),
463     // To FC
464     .mul_en(fc_mul_en),
465     .pixel_vec(fc_pixel_vec),
466     .kernel_vec(fc_kernel_vec),
467     .all_done(fc_all_done),
468     .fc_done(fc_done)
469 );

```

```

470
471 // --- FC ---
472 FC #(
473     .PIX_W(8), .KER_W(8), .K(288), .ACC_W(32),
474     .OUT_CHANNELS(OUT_CHANNELS)
475 ) u_fc (
476     .clk(clk), .rst_n(rst_n),
477     .cfg_valid(fc_cfg_valid),
478     .cfg_ready(fc_cfg_ready),
479     .cfg_data(fc_cfg_data),
480     .cfg_last(fc_cfg_last),
481     .param_load_done(fc_param_load_done),
482     .mul_en(fc_mul_en),
483     .pixel_vec(fc_pixel_vec),
484     .kernel_vec(fc_kernel_vec),
485     .acc_vec(fc_acc_vec),
486     .mul_done_vec(fc_mul_done_vec),
487     .all_done(fc_all_done)
488 );
489
490 endmodule

```

## A.7.4 Network-Level FSM

*Controls the high-level model-load and inference schedule across convolution, FC, and argmax phases.*

```
code/input/RTL/fsm/top_fsm.v
```

```

1 // TopFSM: network-level sequencer.
2 //
3 // MODEL_LOAD (one-time): host sends cfg -> weight -> image via load_data bus
4 // INFER (repeatable): L1 -> L2p0 -> L2p1 -> L3p0 -> L3p1 -> FC -> ARGMAX -> READY
5 //
6 // Preload follows SRAM_CONNECTION_NOTES layer0 convention:
7 //   PL_CFG   -> sram_a_start(layer0, CFG_READ)
8 //   PL_WT    -> sram_a_start(layer0, WT_READ)
9 //   PL_PIXEL -> sram_a_start(layer0, DATA_READ)
10 // Each segment ends when host asserts load_last on the valid&&ready beat.
11
12 module top_fsm #(
13     parameter integer ACC_W      = 32, // FC accumulator width (matches SRAM/cfg 32-bit path)
14     parameter integer OUT_CHANNELS = 10 // FC output channels (0-9 gesture classes)
15 )(
16     input wire      clk,
17     input wire      rst_n,
18
19     // --- External data port (32-bit parallel, valid/ready) ---
20     input wire      load_sel,      // 0=model, 1=image
21     input wire      load_valid,

```

```

22  input wire [31:0] load_data,
23  input wire      load_last,
24  output wire     load_ready,
25
26  // --- External status ---
27  output wire     busy,
28  output reg      predict_valid,
29  output reg [3:0] predict_class, // 10-class gesture id (0-9)
30  output reg [31:0] profile_l1_cycles,
31  output reg [31:0] profile_l2_p0_cycles,
32  output reg [31:0] profile_l2_p1_cycles,
33  output reg [31:0] profile_l3_p0_cycles,
34  output reg [31:0] profile_l3_p1_cycles,
35  output reg [31:0] profile_fc_cycles,
36  output reg [31:0] profile_argmax_cycles,
37  output reg [31:0] profile_total_cycles,
38
39  // --- LayerRunnerFSM control ---
40  output reg      runner_start,
41  output reg [1:0] runner_layer_sel,
42  output reg      runner_pass_id,
43  output reg      runner_is_fc,
44  input wire     runner_done,
45
46  // --- SRAM_A preload control (active during MODEL_LOAD / LOAD_IMAGE) ---
47  // During inference, SRAM_A is driven by LayerRunnerFSM via system_top MUX.
48  output reg      sram_a_start,
49  output reg [2:0] sram_a_layer_sel,
50  output reg [1:0] sram_a_data_sel, // 0=CFG, 1=WT, 2=DATA
51  input wire     sram_a_done,
52
53  // --- Preload write data (directly forwarded to SRAM_A wrapper write port) ---
54  output wire     preload_wr_valid,
55  output wire [31:0] preload_wr_data,
56  input wire     preload_wr_ready, // backpressure from SRAM_A wrapper
57
58  // --- L1 pixel bypass: load_data streams directly to conv ---
59  output reg      pixel_stream_active, // high during L1 pixel streaming
60  input wire     conv_adapter_up_ready, // backpressure from conv_data_adapter
61  input wire     conv_frame_rearm,    // from Conv1_top: all pixels processed
62
63  // --- FC accumulator results (packed vector, LSB = channel 0) ---
64  input wire [OUT_CHANNELS*ACC_W-1:0] fc_acc_vec
65 );
66
67  // -----
68  // State encoding
69  // -----
70  localparam [4:0] ST_IDLE      = 5'd0,
71                  ST_PL_CFG    = 5'd1, // MODEL_LOAD sub-state: conv cfg
72                  ST_PL_CFG_W  = 5'd2,
73                  ST_PL_WT     = 5'd3, // MODEL_LOAD sub-state: conv weight
74                  ST_PL_WT_W   = 5'd4,
75                  ST_PL_PIXEL  = 5'd5,

```

```

76         ST_PL_PIXEL_W = 5'd6,
77         ST_READY      = 5'd7,
78         ST_L1         = 5'd8,
79         ST_L2_P0      = 5'd9,
80         ST_L2_P1      = 5'd10,
81         ST_L3_P0      = 5'd11,
82         ST_L3_P1      = 5'd12,
83         ST_FC         = 5'd13,
84         ST_ARGMAX     = 5'd14,
85         ST_PL_FC_CFG  = 5'd15, // MODEL_LOAD: 10 FC bias words into SRAM_A@0x111
86         ST_PL_FC_CFG_W = 5'd16,
87         ST_PL_FCW     = 5'd17, // MODEL_LOAD: 864 host words into SRAM_FCW packer
88         ST_PL_FCW_W   = 5'd18;
89
90 // SRAM_A layer0 preload encodings
91 localparam [2:0] LSEL_LAYER0 = 3'b000;
92 localparam [1:0] SEL_CFG    = 2'd0,
93                 SEL_WT     = 2'd1,
94                 SEL_DATA   = 2'd2,
95                 SEL_FCW    = 2'd3;
96
97 reg [4:0] state;
98 reg      preload_active; // data forwarding enabled
99 reg      load_last_seen; // current segment's load_last received
100 reg      model_loaded;   // model has been loaded at least once
101 reg      preload_sram_done; // latched sram_a_done during preload wait
102 reg      runner_started; // prevents repeated runner_start pulses
103 reg      runner_prepared; // gives downstream logic one cycle to latch layer config
104
105 // -----
106 // Preload data forwarding / L1 pixel bypass
107 // -----
108 // CFG/WT preload: forward load_data to SRAM_A write port
109 // L1 pixel streaming: load_data bypasses SRAM_A, goes directly to conv
110 assign preload_wr_valid = preload_active && load_valid;
111 assign preload_wr_data  = load_data;
112 assign load_ready      = (preload_active && preload_wr_ready)
113                        | (pixel_stream_active && conv_adapter_up_ready);
114
115 assign busy = (state != ST_IDLE) && (state != ST_READY);
116
117 // -----
118 // Iterative argmax state
119 // -----
120 reg [3:0] argmax_scan_idx;
121 reg [3:0] argmax_best_idx;
122 reg signed [ACC_W-1:0] argmax_best_val;
123 wire signed [ACC_W-1:0] argmax_scan_val =
124     $signed(fc_acc_vec[argmax_scan_idx*ACC_W +: ACC_W]);
125
126 // -----
127 // FSM
128 // -----
129 always @(posedge clk or negedge rst_n) begin

```

```

130     if (!rst_n) begin
131         state             <= ST_IDLE;
132         predict_valid    <= 1'b0;
133         predict_class    <= 4'd0;
134         runner_start     <= 1'b0;
135         runner_layer_sel <= 2'b00;
136         runner_pass_id   <= 1'b0;
137         runner_is_fc     <= 1'b0;
138         profile_l1_cycles <= 32'd0;
139         profile_l2_p0_cycles <= 32'd0;
140         profile_l2_p1_cycles <= 32'd0;
141         profile_l3_p0_cycles <= 32'd0;
142         profile_l3_p1_cycles <= 32'd0;
143         profile_fc_cycles <= 32'd0;
144         profile_argmax_cycles <= 32'd0;
145         profile_total_cycles <= 32'd0;
146         sram_a_start     <= 1'b0;
147         sram_a_layer_sel <= 3'd0;
148         sram_a_data_sel  <= 2'b00;
149         preload_active   <= 1'b0;
150         load_last_seen   <= 1'b0;
151         model_loaded     <= 1'b0;
152         preload_sram_done <= 1'b0;
153         runner_started   <= 1'b0;
154         runner_prepared  <= 1'b0;
155         pixel_stream_active <= 1'b0;
156         argmax_scan_idx  <= 4'd0;
157         argmax_best_idx  <= 4'd0;
158         argmax_best_val  <= {ACC_W{1'b0}};
159     end else begin
160         // Default: clear one-cycle pulses
161         sram_a_start <= 1'b0;
162         runner_start <= 1'b0;
163         predict_valid <= 1'b0;
164
165         // Track load_last during preload (fires on valid && load_last)
166         if (preload_active && load_valid && load_last)
167             load_last_seen <= 1'b1;
168
169         // Latch sram_a_done during preload wait (done is a 1-cycle pulse)
170         if (preload_active && sram_a_done)
171             preload_sram_done <= 1'b1;
172
173         if (state != ST_IDLE && state != ST_READY &&
174             state != ST_PL_CFG && state != ST_PL_CFG_W &&
175             state != ST_PL_WT && state != ST_PL_WT_W &&
176             state != ST_PL_PIXEL && state != ST_PL_PIXEL_W &&
177             state != ST_PL_FC_CFG && state != ST_PL_FC_CFG_W &&
178             state != ST_PL_FCW && state != ST_PL_FCW_W) begin
179             profile_total_cycles <= profile_total_cycles + 32'd1;
180         end
181
182         case (state)
183             ST_L1:    profile_l1_cycles <= profile_l1_cycles + 32'd1;

```

```

184     ST_L2_P0: profile_l2_p0_cycles <= profile_l2_p0_cycles + 32'd1;
185     ST_L2_P1: profile_l2_p1_cycles <= profile_l2_p1_cycles + 32'd1;
186     ST_L3_P0: profile_l3_p0_cycles <= profile_l3_p0_cycles + 32'd1;
187     ST_L3_P1: profile_l3_p1_cycles <= profile_l3_p1_cycles + 32'd1;
188     ST_FC:    profile_fc_cycles <= profile_fc_cycles + 32'd1;
189     ST_ARGMAX: profile_argmax_cycles <= profile_argmax_cycles + 32'd1;
190     default: begin
191         end
192     endcase
193
194     case (state)
195         // -----
196         ST_IDLE: begin
197             preload_active <= 1'b0;
198             if (load_valid && load_sel == 1'b0) begin
199                 // Host wants to load model -> start MODEL_LOAD
200                 state <= ST_PL_CFG;
201             end
202         end
203
204         // -----
205         // MODEL_LOAD: 3 sub-states (PL_CFG -> PL_WT -> PL_PIXEL)
206         // -----
207         ST_PL_CFG: begin
208             sram_a_start      <= 1'b1;
209             sram_a_layer_sel  <= LSEL_LAYER0;
210             sram_a_data_sel   <= SEL_CFG;
211             preload_active    <= 1'b1;
212             load_last_seen    <= 1'b0;
213             preload_sram_done <= 1'b0;
214             state             <= ST_PL_CFG_W;
215         end
216
217         ST_PL_CFG_W: begin
218             if (load_last_seen && (sram_a_done || preload_sram_done)) begin
219                 preload_active <= 1'b0;
220                 if (!load_valid) begin
221                     load_last_seen <= 1'b0;
222                     state          <= ST_PL_WT;
223                 end
224             end
225         end
226
227         ST_PL_WT: begin
228             sram_a_start      <= 1'b1;
229             sram_a_layer_sel  <= LSEL_LAYER0;
230             sram_a_data_sel   <= SEL_WT;
231             preload_active    <= 1'b1;
232             load_last_seen    <= 1'b0;
233             preload_sram_done <= 1'b0;
234             state             <= ST_PL_WT_W;
235         end
236
237         ST_PL_WT_W: begin

```

```

238     if (load_last_seen && (sram_a_done || preload_sram_done)) begin
239         preload_active <= 1'b0;
240         if (!load_valid) begin
241             load_last_seen <= 1'b0;
242             // Continue preload: FC bias words, then FCW weight packed stream.
243             state <= ST_PL_FC_CFG;
244         end
245     end
246 end
247
248 // FC bias: 10 x 32-bit words written to SRAM_A at FC_CFG_BASE (0x111).
249 // Uses LAYER_FC + SEL_CFG; controller enables both read/write modes,
250 // and preload_mode in system_top gates the write path while blocking
251 // conv_pool_valid so only the host write actually fires.
252 ST_PL_FC_CFG: begin
253     sram_a_start <= 1'b1;
254     sram_a_layer_sel <= 3'd4; // LAYER_FC
255     sram_a_data_sel <= SEL_CFG;
256     preload_active <= 1'b1;
257     load_last_seen <= 1'b0;
258     preload_sram_done <= 1'b0;
259     state <= ST_PL_FC_CFG_W;
260 end
261
262 ST_PL_FC_CFG_W: begin
263     if (load_last_seen && (sram_a_done || preload_sram_done)) begin
264         preload_active <= 1'b0;
265         if (!load_valid) begin
266             load_last_seen <= 1'b0;
267             state <= ST_PL_FC_W;
268         end
269     end
270 end
271
272 // FC weight: 864 x 32-bit words packed by fcw_preload_packer into
273 // 288 x 80-bit words inside SRAM_FCW.
274 ST_PL_FC_W: begin
275     sram_a_start <= 1'b1;
276     sram_a_layer_sel <= LSEL_LAYER0; // LAYER_PRELOAD
277     sram_a_data_sel <= SEL_FC_W;
278     preload_active <= 1'b1;
279     load_last_seen <= 1'b0;
280     preload_sram_done <= 1'b0;
281     state <= ST_PL_FC_W_W;
282 end
283
284 ST_PL_FC_W_W: begin
285     if (load_last_seen && (sram_a_done || preload_sram_done)) begin
286         preload_active <= 1'b0;
287         if (!load_valid) begin
288             load_last_seen <= 1'b0;
289             // MODEL_LOAD complete.
290             model_loaded <= 1'b1;
291             state <= ST_READY;

```

```

292     end
293     end
294 end
295
296 // ST_PL_PIXEL / ST_PL_PIXEL_W: removed - pixels now stream directly
297 // to conv_data_adapter via system_top during ST_L1.
298
299 // -----
300 ST_READY: begin
301     preload_active <= 1'b0;
302     if (load_valid && !load_last && load_sel == 1'b1) begin
303         profile_l1_cycles <= 32'd0;
304         profile_l2_p0_cycles <= 32'd0;
305         profile_l2_p1_cycles <= 32'd0;
306         profile_l3_p0_cycles <= 32'd0;
307         profile_l3_p1_cycles <= 32'd0;
308         profile_fc_cycles <= 32'd0;
309         profile_argmax_cycles <= 32'd0;
310         profile_total_cycles <= 32'd0;
311         // Host wants to send image -> enter L1 (pixels stream directly)
312         state <= ST_L1;
313     end else if (load_valid && !load_last && load_sel == 1'b0) begin
314         // Host wants to reload model
315         model_loaded <= 1'b0;
316         state <= ST_PL_CFG;
317     end
318 end
319
320 // -----
321 // Inference: L1 -> L2p0 -> L2p1 -> L3p0 -> L3p1 -> FC -> ARGMAX
322 // Each state: pulse runner_start once, then wait for runner_done.
323 // runner_started prevents repeated start pulses while waiting.
324 // -----
325 ST_L1: begin
326     if (!runner_prepared) begin
327         runner_layer_sel <= 2'b00;
328         runner_pass_id <= 1'b0;
329         runner_is_fc <= 1'b0;
330         runner_prepared <= 1'b1;
331     end else if (!runner_started) begin
332         runner_start <= 1'b1;
333         runner_started <= 1'b1;
334         pixel_stream_active <= 1'b1; // route load_data directly to conv
335     end
336     if (runner_done) begin
337         state <= ST_L2_P0;
338         runner_started <= 1'b0;
339         runner_prepared <= 1'b0;
340         pixel_stream_active <= 1'b0;
341     end
342 end
343
344 ST_L2_P0: begin
345     if (!runner_prepared) begin

```

```

346     runner_layer_sel <= 2'b01;
347     runner_pass_id  <= 1'b0;
348     runner_is_fc    <= 1'b0;
349     runner_prepared <= 1'b1;
350     end else if (!runner_started) begin
351         runner_start    <= 1'b1;
352         runner_started  <= 1'b1;
353     end
354     if (runner_done) begin state <= ST_L2_P1; runner_started <= 1'b0; runner_prepared <= 1'b0; end
355 end
356
357 ST_L2_P1: begin
358     if (!runner_prepared) begin
359         runner_layer_sel <= 2'b01;
360         runner_pass_id  <= 1'b1;
361         runner_is_fc    <= 1'b0;
362         runner_prepared <= 1'b1;
363     end else if (!runner_started) begin
364         runner_start    <= 1'b1;
365         runner_started  <= 1'b1;
366     end
367     if (runner_done) begin state <= ST_L3_P0; runner_started <= 1'b0; runner_prepared <= 1'b0; end
368 end
369
370 ST_L3_P0: begin
371     if (!runner_prepared) begin
372         runner_layer_sel <= 2'b10;
373         runner_pass_id  <= 1'b0;
374         runner_is_fc    <= 1'b0;
375         runner_prepared <= 1'b1;
376     end else if (!runner_started) begin
377         runner_start    <= 1'b1;
378         runner_started  <= 1'b1;
379     end
380     if (runner_done) begin state <= ST_L3_P1; runner_started <= 1'b0; runner_prepared <= 1'b0; end
381 end
382
383 ST_L3_P1: begin
384     if (!runner_prepared) begin
385         runner_layer_sel <= 2'b10;
386         runner_pass_id  <= 1'b1;
387         runner_is_fc    <= 1'b0;
388         runner_prepared <= 1'b1;
389     end else if (!runner_started) begin
390         runner_start    <= 1'b1;
391         runner_started  <= 1'b1;
392     end
393     if (runner_done) begin state <= ST_FC; runner_started <= 1'b0; runner_prepared <= 1'b0; end
394 end
395
396 ST_FC: begin
397     if (!runner_prepared) begin
398         runner_layer_sel <= 2'b00;
399         runner_pass_id  <= 1'b0;

```

```

400     runner_is_fc    <= 1'b1;
401     runner_prepared <= 1'b1;
402 end else if (!runner_started) begin
403     runner_start    <= 1'b1;
404     runner_started  <= 1'b1;
405 end
406 if (runner_done) begin
407     state          <= ST_ARGMAX;
408     runner_started <= 1'b0;
409     runner_prepared <= 1'b0;
410     argmax_scan_idx <= 4'd1;
411     argmax_best_idx <= 4'd0;
412     argmax_best_val <= $signed(fc_acc_vec[0 +: ACC_W]);
413 end
414 end
415
416 // -----
417 ST_ARGMAX: begin
418     // Compare one class per cycle. This removes the long 10-way
419     // comparator chain from the fabric critical path.
420     if (argmax_scan_val > argmax_best_val) begin
421         argmax_best_val <= argmax_scan_val;
422         argmax_best_idx <= argmax_scan_idx;
423     end
424
425     if (argmax_scan_idx == OUT_CHANNELS - 1) begin
426         predict_class <= (argmax_scan_val > argmax_best_val)
427             ? argmax_scan_idx
428             : argmax_best_idx;
429         predict_valid <= 1'b1;
430         state          <= ST_READY;
431     end else begin
432         argmax_scan_idx <= argmax_scan_idx + 4'd1;
433     end
434 end
435
436 default: state <= ST_IDLE;
437 endcase
438 end
439 end
440
441 endmodule

```

### A.7.5 Layer Runner FSM

*Sequences per-layer configuration, weights, input data, output writes, and completion for the CNN layers.*

```

1 // LayerRunnerFSM: single-layer / single-pass transaction controller.
2 //
3 // Conv mode (is_fc=0):
4 //  IDLE -> LOAD_CFG -> WAIT_CFG -> LOAD_WT -> WAIT_WT -> STREAM -> WAIT_DONE -> DONE_PULSE
5 //
6 // FC mode (is_fc=1):
7 //  IDLE -> LOAD_CFG -> WAIT_CFG -> STREAM -> WAIT_DONE -> DONE_PULSE
8 //
9 // SRAM_A interface uses data_sel (not op_type):
10 //  data_sel = 0 (CFG), 1 (WT), 2 (DATA)
11 //  Read/write direction derived by wrapper from layer_sel + data_sel.
12 //
13 // All state transitions driven by SRAM wrapper done signals.
14
15 module layer_runner_fsm (
16     input wire      clk,
17     input wire      rst_n,
18
19     // --- TopFSM control ---
20     input wire      start,
21     input wire [1:0] layer_sel,      // 00=L1, 01=L2, 10=L3
22     input wire      pass_id,
23     input wire      is_fc,
24     output wire     busy,
25     output reg      done,
26
27     // --- SRAM_A wrapper control ---
28     output reg      sram_a_start,
29     output reg [2:0] sram_a_layer_sel,
30     output reg [1:0] sram_a_data_sel, // 0=CFG, 1=WT, 2=DATA
31     output reg      sram_a_pass_id,
32     input wire      sram_a_done,
33
34     // --- SRAM_B wrapper control ---
35     output reg      sram_b_start,
36     output reg [2:0] sram_b_layer_sel,
37     output reg [1:0] sram_b_data_sel,
38     output reg      sram_b_pass_id,
39     input wire      sram_b_done,
40
41     // --- Conv path observation (auxiliary, NOT for transitions) ---
42     input wire      cfg_load_done,
43     input wire      wt_load_done,
44     input wire      pool_frame_done,
45     input wire      conv_frame_rearm,
46
47     // --- FC path feedback ---
48     input wire      fc_done
49 );
50
51 localparam [2:0] ST_IDLE      = 3'd0,
52                 ST_LOAD_CFG  = 3'd1,

```

```

53             ST_WAIT_CFG = 3'd2,
54             ST_LOAD_WT  = 3'd3,
55             ST_WAIT_WT  = 3'd4,
56             ST_STREAM   = 3'd5,
57             ST_WAIT_DONE = 3'd6,
58             ST_DONE     = 3'd7;
59
60     localparam [1:0] SEL_CFG = 2'd0,
61                   SEL_WT   = 2'd1,
62                   SEL_DATA = 2'd2;
63
64     localparam [2:0] LSEL_L1 = 3'd1,
65                   LSEL_L2 = 3'd2,
66                   LSEL_L3 = 3'd3,
67                   LSEL_FC = 3'd4;
68
69     reg [2:0] state;
70
71     // Debug-visible mirrors kept for TB/wave compatibility. They are no longer
72     // used to drive transactions.
73     reg [1:0] layer_sel_r;
74     reg      pass_id_r;
75     reg      is_fc_r;
76
77     // TopFSM keeps layer_sel / pass_id / is_fc stable while a layer is running.
78     // Use the live inputs for transaction control so the runner cannot start the
79     // next layer with stale latched parameters.
80
81     function [2:0] to_lsel;
82         input [1:0] ls;
83         input      fc;
84         begin
85             if (fc) to_lsel = LSEL_FC;
86             else case (ls)
87                 2'b00: to_lsel = LSEL_L1;
88                 2'b01: to_lsel = LSEL_L2;
89                 2'b10: to_lsel = LSEL_L3;
90                 default: to_lsel = LSEL_L1;
91             endcase
92         end
93     endfunction
94
95     // WAIT_DONE: accumulate done signals from both wrappers + pool
96     reg sram_a_done_seen, sram_b_done_seen;
97     reg pool_done_seen;
98     reg frame_rearm_seen;
99
100    assign busy = (state != ST_IDLE);
101
102    always @(posedge clk or negedge rst_n) begin
103        if (!rst_n) begin
104            state      <= ST_IDLE;
105            done       <= 1'b0;
106            sram_a_start <= 1'b0;

```

```

107     sram_a_layer_sel <= 3'd0;
108     sram_a_data_sel  <= 2'd0;
109     sram_a_pass_id   <= 1'b0;
110     sram_b_start     <= 1'b0;
111     sram_b_layer_sel <= 3'd0;
112     sram_b_data_sel  <= 2'd0;
113     sram_b_pass_id   <= 1'b0;
114     layer_sel_r      <= 2'b00;
115     pass_id_r        <= 1'b0;
116     is_fc_r          <= 1'b0;
117     sram_a_done_seen <= 1'b0;
118     sram_b_done_seen <= 1'b0;
119     pool_done_seen   <= 1'b0;
120     frame_rearm_seen <= 1'b0;
121 end else begin
122     // Clear one-cycle pulses
123     sram_a_start <= 1'b0;
124     sram_b_start <= 1'b0;
125     done        <= 1'b0;
126
127     case (state)
128     // -----
129     ST_IDLE: begin
130         if (start) begin
131             layer_sel_r <= layer_sel;
132             pass_id_r   <= pass_id;
133             is_fc_r     <= is_fc;
134             state       <= ST_LOAD_CFG;
135         end
136     end
137
138     // -----
139     ST_LOAD_CFG: begin
140         sram_a_start <= 1'b1;
141         sram_a_layer_sel <= to_lsel(layer_sel, is_fc);
142         sram_a_data_sel <= SEL_CFG;
143         sram_a_pass_id <= pass_id;
144         sram_a_done_seen <= 1'b0;
145         state          <= ST_WAIT_CFG;
146     end
147
148     ST_WAIT_CFG: begin
149         if (sram_a_done) sram_a_done_seen <= 1'b1;
150         if (is_fc) begin
151             if ((sram_a_done_seen || sram_a_done) && cfg_load_done) begin
152                 sram_a_done_seen <= 1'b0;
153                 state <= ST_STREAM; // FC skips LOAD_WT
154             end
155         end else begin
156             if (sram_a_done)
157                 state <= ST_LOAD_WT;
158         end
159     end
160

```

```

161 // -----
162 ST_LOAD_WT: begin
163     sram_a_start      <= 1'b1;
164     sram_a_layer_sel  <= to_lsel(layer_sel, 1'b0);
165     sram_a_data_sel   <= SEL_WT;
166     sram_a_pass_id    <= pass_id;
167     state             <= ST_WAIT_WT;
168 end
169
170 ST_WAIT_WT: begin
171     if (sram_a_done)
172         state <= ST_STREAM;
173 end
174
175 // -----
176 ST_STREAM: begin
177     sram_a_done_seen <= 1'b0;
178     sram_b_done_seen <= 1'b0;
179     pool_done_seen  <= 1'b0;
180     frame_rearm_seen <= 1'b0;
181
182     if (is_fc) begin
183         // FC: SRAM_A reads interleaved weight, SRAM_B reads packed data
184         sram_a_start      <= 1'b1;
185         sram_a_layer_sel  <= LSEL_FC;
186         sram_a_data_sel   <= SEL_WT;
187         sram_a_pass_id    <= 1'b0;
188
189         sram_b_start      <= 1'b1;
190         sram_b_layer_sel  <= LSEL_FC;
191         sram_b_data_sel   <= SEL_DATA;
192         sram_b_pass_id    <= 1'b0;
193     end else begin
194         // Conv: start DATA transaction on both SRAMs
195         // L1: pixels stream directly from external load_data (no SRAM_A read)
196         // L2/L3: read DATA from SRAM_A (or SRAM_B)
197         if (layer_sel != 2'b00) begin
198             sram_a_start      <= 1'b1;
199             sram_a_layer_sel  <= to_lsel(layer_sel, 1'b0);
200             sram_a_data_sel   <= SEL_DATA;
201             sram_a_pass_id    <= pass_id;
202         end
203
204         sram_b_start      <= 1'b1;
205         sram_b_layer_sel  <= to_lsel(layer_sel, 1'b0);
206         sram_b_data_sel   <= SEL_DATA;
207         sram_b_pass_id    <= pass_id;
208     end
209
210     state <= ST_WAIT_DONE;
211 end
212
213 // -----
214 ST_WAIT_DONE: begin

```

```

215     if (sram_a_done) sram_a_done_seen <= 1'b1;
216     if (sram_b_done) sram_b_done_seen <= 1'b1;
217     if (pool_frame_done) pool_done_seen <= 1'b1;
218     if (conv_frame_rearm) frame_rearm_seen <= 1'b1;
219
220     // Both wrapper dones must arrive
221     if (is_fc) begin
222         if (fc_done &&
223             (sram_a_done_seen || sram_a_done) &&
224             (sram_b_done_seen || sram_b_done))
225             state <= ST_DONE;
226     end else if (layer_sel == 2'b00) begin
227         // L1: pixels stream from external (no SRAM_A read), only wait
228         // for SRAM_B (pool writeback) and conv_frame_rearm.
229         if ((sram_b_done_seen || sram_b_done) &&
230             (frame_rearm_seen || conv_frame_rearm))
231             state <= ST_DONE;
232     end else begin
233         if ((sram_a_done_seen || sram_a_done) &&
234             (sram_b_done_seen || sram_b_done) &&
235             (frame_rearm_seen || conv_frame_rearm))
236             state <= ST_DONE;
237     end
238 end
239
240 // -----
241 ST_DONE: begin
242     done <= 1'b1;
243     state <= ST_IDLE;
244 end
245
246 default: state <= ST_IDLE;
247 endcase
248 end
249 end
250
251 endmodule

```

## A.7.6 Convolution Data Adapter

*Adapts SRAM and MMIO image words into the byte/word stream expected by convolution layers.*

```
code/input/RTL/fsm/conv_data_adapter.v
```

```

1 // conv_data_adapter: unpack 32-bit SRAM words into Conv input beats.
2 //
3 // L1 (layer_sel=00): 1 SRAM word -> 4 byte beats (1-word holding register)
4 //   SRAM outputs 1024 words, Conv receives 4096 beats x 8-bit, in_byte_en=0001
5 // L2/L3 also use a 1-word holding register so the adapter can safely absorb
6 // the SRAM wrapper's 1-cycle read pipeline under downstream backpressure.
7

```

```

8 module conv_data_adapter (
9     input wire      clk,
10    input wire      rst_n,
11    input wire [1:0] layer_sel,
12
13    // Upstream: from SRAM read stream
14    input wire      up_valid,
15    input wire [31:0] up_data,
16    input wire      up_last,
17    output wire     up_ready,
18
19    // Downstream: to Conv in port
20    output wire     dn_valid,
21    output wire [31:0] dn_data,
22    output wire [3:0] dn_byte_en,
23    output wire     dn_last,
24    input wire      dn_ready
25 );
26
27 wire is_l1 = (layer_sel == 2'b00);
28
29 // -----
30 // L1: byte unpack with 1-word holding register
31 // L2/L3: buffered word pass-through
32 // -----
33
34 reg [31:0] hold_reg;
35 reg      hold_valid;
36 reg [1:0] byte_sel; // 0,1,2,3
37 reg      hold_is_last; // the held word was the last SRAM word
38
39 wire [31:0] cur_word_l1 = hold_valid ? hold_reg : up_data;
40 wire [1:0]  cur_byte_sel = hold_valid ? byte_sel : 2'd0;
41 wire      cur_is_last_l1 = hold_valid ? hold_is_last : up_last;
42
43 // Current pixel byte from active word (held word or fresh bypass)
44 wire [7:0] pixel_byte = (cur_byte_sel == 2'd0) ? cur_word_l1[ 7: 0] :
45                        (cur_byte_sel == 2'd1) ? cur_word_l1[15: 8] :
46                        (cur_byte_sel == 2'd2) ? cur_word_l1[23:16] :
47                        cur_word_l1[31:24];
48
49 // L1 downstream signals
50 wire l1_dn_valid = hold_valid || up_valid;
51 wire l1_dn_last  = l1_dn_valid && cur_is_last_l1 && (cur_byte_sel == 2'd3);
52
53 // L1 up_ready: assert only when adapter WILL BE free next cycle.
54 // SRAM wrapper has 1-cycle read pipeline (read_en -> data_valid next cycle),
55 // so up_ready must predict availability 1 cycle ahead.
56 // Case A: hold_valid=1, last byte being consumed -> free next cycle
57 // Case B: hold_valid=0 AND no data arriving -> already free
58 // When data is arriving (up_valid=1, hold_valid=0), the adapter is latching
59 // this cycle - NOT free for another read next cycle.
60 wire l1_will_be_free = (hold_valid && dn_ready && byte_sel == 2'd3)
61                        || (!hold_valid && !up_valid);

```

```

62 wire l1_up_ready = l1_will_be_free;
63
64 // L2/L3: one-word buffer to absorb the upstream read pipeline.
65 // Byte-swap: pool_data stores {ch0,ch1,ch2,ch3} with ch0 at MSB,
66 // but Conv pixel_serializer expects ch0 at LSB (in_data[7:0]=ch0).
67 wire [31:0] nonl1_raw = hold_valid ? hold_reg : up_data;
68 wire [31:0] nonl1_word = {nonl1_raw[7:0], nonl1_raw[15:8],
69                          nonl1_raw[23:16], nonl1_raw[31:24]};
70 wire        nonl1_is_last = hold_valid ? hold_is_last : up_last;
71 wire nonl1_dn_valid = hold_valid || up_valid;
72 wire nonl1_dn_last  = nonl1_dn_valid && nonl1_is_last;
73 wire nonl1_will_be_free = (hold_valid && dn_ready)
74                          || (!hold_valid && (!up_valid || dn_ready));
75 wire nonl1_up_ready = nonl1_will_be_free;
76
77 // Muxed outputs
78 assign dn_valid  = is_l1 ? l1_dn_valid : nonl1_dn_valid;
79 assign dn_data   = is_l1 ? {24'd0, pixel_byte} : nonl1_word;
80 assign dn_byte_en = is_l1 ? 4'b0001 : 4'b1111;
81 assign dn_last   = is_l1 ? l1_dn_last : nonl1_dn_last;
82 assign up_ready  = is_l1 ? l1_up_ready : nonl1_up_ready;
83
84 // -----
85 // L1 holding register and byte_sel control
86 // -----
87 always @(posedge clk or negedge rst_n) begin
88     if (!rst_n) begin
89         hold_reg    <= 32'd0;
90         hold_valid  <= 1'b0;
91         byte_sel    <= 2'd0;
92         hold_is_last <= 1'b0;
93     end else if (is_l1) begin
94         if (!hold_valid) begin
95             if (up_valid) begin
96                 hold_reg    <= up_data;
97                 hold_is_last <= up_last;
98                 hold_valid  <= 1'b1;
99                 // Fresh word can emit byte0 immediately; keep byte1 queued next.
100                byte_sel    <= dn_ready ? 2'd1 : 2'd0;
101            end
102        end else if (dn_ready) begin
103            if (byte_sel == 2'd3) begin
104                // Finished current word
105                hold_valid <= 1'b0;
106                byte_sel  <= 2'd0;
107            end else begin
108                byte_sel <= byte_sel + 2'd1;
109            end
110        end
111    end else begin
112        if (!hold_valid) begin
113            if (up_valid) begin
114                if (!dn_ready) begin
115                    hold_reg    <= up_data;

```

```

116         hold_valid  <= 1'b1;
117         hold_is_last <= up_last;
118     end
119     byte_sel    <= 2'd0;
120 end
121 end else if (dn_ready) begin
122     hold_valid <= 1'b0;
123     byte_sel   <= 2'd0;
124 end
125 end
126 end
127
128 endmodule

```

### A.7.7 Weight Prepad Inserter

*Pads and aligns convolution weight streams before they enter the convolution datapath.*

```
code/input/RTL/fsm/wt_prepad_inserter.v
```

```

1 // wt_prepad_inserter: insert 3 zero words before real weight data.
2 //
3 // Sits between SRAM_A read stream and Conv wt port.
4 // When enabled (is_wt_read=1), inserts 3 zero beats then forwards real data.
5 // When disabled (is_wt_read=0), passes through transparently.
6 //
7 // Conv weight_buffer expects DOT_K+3 beats: 3 prepad zeros + DOT_K real weights.
8 // SRAM only stores DOT_K real weights. This module bridges the gap.
9
10 module wt_prepad_inserter (
11     input wire    clk,
12     input wire    rst_n,
13
14     // Control: assert when current SRAM_A transaction is WT_READ
15     input wire    is_wt_read,
16
17     // Upstream: SRAM_A read stream
18     input wire    up_valid,
19     input wire [31:0] up_data,
20     input wire    up_last,
21     output wire   up_ready,
22
23     // Downstream: to Conv wt port
24     output wire   dn_valid,
25     output wire [31:0] dn_data,
26     output wire   dn_last,
27     input wire    dn_ready
28 );
29
30 localparam [1:0] ST_IDLE   = 2'd0,
31                 ST_PREPAD = 2'd1,

```

```

32             ST_PASS    = 2'd2;
33
34 reg [1:0] state;
35 reg [1:0] pad_cnt; // counts 0,1,2 (3 zero beats)
36
37 // Rising-edge detection on is_wt_read to avoid re-entering prepad
38 reg is_wt_read_d;
39 wire is_wt_read_rise = is_wt_read && !is_wt_read_d;
40 always @(posedge clk or negedge rst_n) begin
41     if (!rst_n) is_wt_read_d <= 1'b0;
42     else        is_wt_read_d <= is_wt_read;
43 end
44
45 // Pass-through when not in WT_READ mode
46 wire passthrough = !is_wt_read || (state == ST_PASS);
47
48 // SRAM wrapper read is 1-cycle pipelined. On the last zero-prepad beat,
49 // request the first real weight so ST_PASS can start without a bubble.
50 assign up_ready = passthrough ? dn_ready
51                        : ((state == ST_PREPAD) && (pad_cnt == 2'd2) && dn_ready);
52 assign dn_valid = passthrough ? up_valid : (state == ST_PREPAD);
53 assign dn_data  = passthrough ? up_data  : 32'd0;
54 assign dn_last  = passthrough ? up_last  : 1'b0;
55
56 always @(posedge clk or negedge rst_n) begin
57     if (!rst_n) begin
58         state <= ST_IDLE;
59         pad_cnt <= 2'd0;
60     end else begin
61         case (state)
62             ST_IDLE: begin
63                 if (is_wt_read_rise) begin
64                     // is_wt_read just asserted: insert 3 prepad zeros before real data
65                     state <= ST_PREPAD;
66                     pad_cnt <= 2'd0;
67                 end
68             end
69
70             ST_PREPAD: begin
71                 if (dn_ready) begin
72                     if (pad_cnt == 2'd2) begin
73                         state <= ST_PASS; // prepad done, now forward real data
74                     end else begin
75                         pad_cnt <= pad_cnt + 2'd1;
76                     end
77                 end
78             end
79
80             ST_PASS: begin
81                 // Forward until upstream last beat is consumed
82                 if (up_valid && up_ready && up_last) begin
83                     state <= ST_IDLE;
84                 end
85             end

```

```

86
87     default: state <= ST_IDLE;
88     endcase
89     end
90 end
91
92 endmodule

```

## A.8 Convolution Core RTL

### A.8.1 Convolution Top

*Connects input alignment, line buffering, convolution buffering, weight buffering, SA launch, and output streaming.*

code/input/RTL/conv\_core/conv\_top.v

```

1 module conv_top #(
2     parameter integer W      = 64,    // max line buffer width (L1)
3     parameter integer DATA_W = 8,
4     parameter integer ROWS   = 16,
5     parameter integer COLS   = 4,
6     parameter integer DOT_K  = 72,    // max dot product length (L3: 9*8)
7     parameter integer ACC_W  = 23,
8     parameter integer IMG_H  = 64,
9     parameter integer IMG_W  = 64,
10    parameter integer C_IN   = 8     // max input channels (L3)
11 ) (
12     input wire [1:0]        layer_sel,
13     input wire clk,
14     input wire rst_n,
15
16     // Pixel stream input (ready/valid handshake)
17     input wire              in_valid,
18     input wire [31:0]      in_data,
19     input wire [3:0]       in_byte_en,
20     input wire              in_last,
21     output wire            in_ready,
22
23     // Weight load channel
24     input wire              wt_valid,
25     output wire             wt_ready,
26     input wire              wt_last,
27     input wire signed [COLS*DATA_W-1:0] wt_data,
28
29     // Weight load done (1-cycle pulse on weights_loaded rising edge)
30     output wire             wt_load_done,
31
32     // Systolic-array output
33     output wire             sa_done,

```

```

34 // Per-column real-time stream output
35 output wire signed [COLS*ACC_W-1:0] c_out_col_stream_flat, // 4*23 bits, one output row
36 output wire [COLS-1:0] c_out_col_valid,
37 output wire [COLS-1:0] c_out_col_last,
38
39 // Frame rearm: Conv backend fully drained, ready for next frame/pass
40 output wire frame_rearm_out,
41 // Frame done: all input pixels received, SA may still be draining
42 output wire frame_done_out
43 );
44
45 // Debug-visible state encoding kept in top hierarchy for TB/wave compatibility.
46 localparam [1:0] ST_IDLE = 2'd0;
47
48 // Runtime IMG_PIXELS decode from layer_sel
49 // Counts input BEATS (not bytes). Each beat is one pix_accept.
50 // L1: 4096 beats (1 byte/beat, 64*64*1 bytes)
51 // L2: 961 beats (4 bytes/beat, 31*31 wide pixels)
52 // L3: 392 beats (4 bytes/beat, 14*14*8/4 half-pixels)
53 reg [1:0] active_layer_sel_q;
54 reg [31:0] eff_img_pixels;
55 always @(posedge clk or negedge rst_n) begin
56     if (!rst_n)
57         active_layer_sel_q <= 2'b00;
58     else
59         active_layer_sel_q <= layer_sel;
60 end
61
62 always @* begin
63     case (active_layer_sel_q)
64         2'b01: eff_img_pixels = 32'd961; // 31*31
65         2'b10: eff_img_pixels = 32'd392; // 14*14*2
66         default: eff_img_pixels = 32'd4096; // 64*64
67     endcase
68 end
69
70 reg [31:0] pix_cnt;
71 reg frame_done;
72
73 wire [1:0] state;
74 wire [8:0] feed_cnt;
75 wire [8:0] wt_load_cnt;
76 wire [8:0] rd_col_r;
77 wire rd_en_r;
78 wire rd_bank_r;
79 wire wb_ring_r;
80 wire start_pulse_r;
81 wire weights_loaded;
82 reg weights_loaded_d;
83 assign wt_load_done = weights_loaded && !weights_loaded_d;
84 always @(posedge clk or negedge rst_n) begin
85     if (!rst_n) weights_loaded_d <= 1'b0;
86     else weights_loaded_d <= weights_loaded;
87 end

```

```

88  wire      wt_load_active;
89  wire      consume_ready_bank;
90  wire      consume_bank_sel;
91  wire [4:0] valid_rows_A, valid_rows_B;
92  wire      partial_pending;
93  // Frame re-arm: declared here, driven below after backend idle detection.
94  wire frame_rearm;
95
96  // Internal observation signals kept hierarchical for TB and wave scripts.
97  wire      pix3_valid;
98  wire signed [C_IN*DATA_W-1:0] row0_pix;
99  wire signed [C_IN*DATA_W-1:0] row1_pix;
100 wire signed [C_IN*DATA_W-1:0] row2_pix;
101 wire signed [DATA_W-1:0] wb_out0, wb_out1, wb_out2, wb_out3;
102 wire [COLS-1:0] wb_out_valid;
103 wire      bank_can_write;
104
105 // Input-side flow control.
106 wire pix_accept = in_valid && in_ready;
107 assign in_ready = !frame_done && bank_can_write;
108
109 // Bank write-side control.
110 wire rd_active = (state != ST_IDLE);
111
112 // Launch arbitration and weight-side handshake.
113 wire have_ready_block;
114 wire wt_fire = wt_valid && wt_ready;
115 wire launch_from_A;
116 wire have_launchable_block = have_ready_block && weights_loaded;
117 // Prevent SA launch after frame_done (avoid stale bank data between layers/passes)
118 wire can_launch = have_launchable_block && !wt_fire && !frame_done;
119 wire [1:0] sa_mode_cfg = active_layer_sel_q;
120 reg [1:0] sa_mode_cfg_q;
121 reg [4:0] sa_valid_rows_cfg_q;
122 reg      sa_start_pulse_q;
123 reg signed [ROWS*DATA_W-1:0] sa_a_in_flat_q;
124 reg signed [COLS*DATA_W-1:0] sa_b_in_flat_q;
125
126 // SA input gating.
127 wire signed [ROWS*DATA_W-1:0] a_in_flat;
128 wire signed [ROWS*DATA_W-1:0] a_in_from_bank;
129 wire signed [ROWS*DATA_W-1:0] raw_firstcol_flat;
130 wire signed [ROWS*COLS*ACC_W-1:0] c_out_raw_flat;
131 wire signed [DATA_W-1:0] wb_col0 = wb_out_valid[0] ? wb_out0 : {DATA_W{1'b0}};
132 wire signed [DATA_W-1:0] wb_col1 = wb_out_valid[1] ? wb_out1 : {DATA_W{1'b0}};
133 wire signed [DATA_W-1:0] wb_col2 = wb_out_valid[2] ? wb_out2 : {DATA_W{1'b0}};
134 wire signed [DATA_W-1:0] wb_col3 = wb_out_valid[3] ? wb_out3 : {DATA_W{1'b0}};
135 wire sa_a_gate = rd_en_r || ((state == ST_IDLE) && can_launch);
136 wire rd_bank_sel = ((state == ST_IDLE) && can_launch) ? (launch_from_A ? 1'b0 : 1'b1) : rd_bank_r;
137 wire signed [COLS*DATA_W-1:0] b_in_to_sa =
138     rd_en_r ? {wb_col3, wb_col2, wb_col1, wb_col0} : {(COLS*DATA_W){1'b0}};
139 // in_last is consumed by frame_done logic below (no longer unused).
140
141 // Gate A inputs to zero while idle to avoid accumulating stale values.

```

```

142 assign a_in_flat = sa_a_gate ? a_in_from_bank : {ROWS*DATA_W{1'b0}};
143
144 // One-cycle input pipeline into the systolic array. This cuts the long
145 // layer-select-dependent Conv_Buffer -> PE MAC path.
146 always @(posedge clk or negedge rst_n) begin
147     if (!rst_n) begin
148         sa_mode_cfg_q         <= 2'b00;
149         sa_valid_rows_cfg_q <= 5'd0;
150         sa_start_pulse_q     <= 1'b0;
151         sa_a_in_flat_q       <= {ROWS*DATA_W{1'b0}};
152         sa_b_in_flat_q       <= {COLS*DATA_W{1'b0}};
153     end else begin
154         sa_mode_cfg_q         <= sa_mode_cfg;
155         sa_valid_rows_cfg_q <= rd_bank_sel ? valid_rows_B : valid_rows_A;
156         sa_start_pulse_q     <= start_pulse_r;
157         sa_a_in_flat_q       <= a_in_flat;
158         sa_b_in_flat_q       <= b_in_to_sa;
159     end
160 end
161
162 input_row_aligner #(
163     .W      (W),
164     .DW     (DATA_W),
165     .C_IN   (C_IN)
166 ) u_input_row_aligner (
167     .layer_sel (active_layer_sel_q),
168     .clk       (clk),
169     .rst_n     (rst_n),
170     .frame_rearm(frame_rearm),
171     .in_valid  (pix_accept),
172     .in_data   (in_data),
173     .in_byte_en (in_byte_en),
174     .pix3_valid (pix3_valid),
175     .row0_out  (row0_pix),
176     .row1_out  (row1_pix),
177     .row2_out  (row2_pix)
178 );
179
180 weight_buffer #(
181     .DATA_W(COLS*DATA_W),
182     .COLS   (COLS),
183     .L0     (DOT_K),
184     .L1     (DOT_K + 1),
185     .L2     (DOT_K + 2),
186     .L3     (DOT_K + 3)
187 ) u_weight_buffer (
188     .layer_sel(active_layer_sel_q),
189     .clk       (clk),
190     .rst_n     (rst_n),
191     .ring      (wb_ring_r),
192     .wr_en     (wt_fire),
193     .in_data   (wt_data),
194     .out_ready(rd_en_r),
195     .out0      (wb_out0),

```

```

196     .out1      (wb_out1),
197     .out2      (wb_out2),
198     .out3      (wb_out3),
199     .out_valid(wb_out_valid)
200 );
201
202 Conv_Buffer #(
203     .DATA_W(DATA_W),
204     .ROWS   (ROWS),
205     .DOT_K  (DOT_K),
206     .C_IN   (C_IN)
207 ) u_conv_buffer (
208     .layer_sel      (active_layer_sel_q),
209     .clk            (clk),
210     .rst_n         (rst_n),
211     .frame_rearm   (frame_rearm),
212     .pix3_valid    (pix3_valid),
213     .row0_pix      (row0_pix),
214     .row1_pix      (row1_pix),
215     .row2_pix      (row2_pix),
216     .rd_active     (rd_active),
217     .rd_en         (rd_en_r),
218     .rd_bank       (rd_bank_sel),
219     .rd_col        (rd_col_r),
220     .consume_ready_bank(consume_ready_bank),
221     .consume_bank_sel (consume_bank_sel),
222     .bank_can_write  (bank_can_write),
223     .have_ready_block (have_ready_block),
224     .launch_from_A   (launch_from_A),
225     .valid_rows_A    (valid_rows_A),
226     .valid_rows_B    (valid_rows_B),
227     .partial_pending (partial_pending),
228     .raw_firstcol_flat (raw_firstcol_flat)
229 );
230
231 sa_skew_feeder #(
232     .ROWS(ROWS)
233 ) u_sa_skew_feeder (
234     .clk            (clk),
235     .rst_n         (rst_n),
236     .frame_rearm   (frame_rearm),
237     .rd_en         (rd_en_r),
238     .in_firstcol_flat(raw_firstcol_flat),
239     .out_firstcol_flat(a_in_from_bank)
240 );
241
242 conv_engine_ctrl #(
243     .DOT_K(DOT_K),
244     .COLS (COLS),
245     .ROWS (ROWS)
246 ) u_conv_engine_ctrl (
247     .layer_sel      (active_layer_sel_q),
248     .clk            (clk),
249     .rst_n         (rst_n),

```

```

250     .wt_valid      (wt_valid),
251     .wt_last      (wt_last),
252     .have_ready_block (have_ready_block),
253     .launch_from_A (launch_from_A),
254     .launch_valid_rows (launch_from_A ? valid_rows_A : valid_rows_B),
255     .frame_rearm    (frame_rearm),
256     .sa_done       (sa_done),
257     .wt_ready      (wt_ready),
258     .weights_loaded (weights_loaded),
259     .wt_load_active (wt_load_active),
260     .wt_load_cnt   (wt_load_cnt),
261     .state_dbg     (state),
262     .feed_cnt_dbg  (feed_cnt),
263     .consume_ready_bank(consume_ready_bank),
264     .consume_bank_sel (consume_bank_sel),
265     .rd_en         (rd_en_r),
266     .rd_bank       (rd_bank_r),
267     .rd_col        (rd_col_r),
268     .wb_ring       (wb_ring_r),
269     .start_pulse   (start_pulse_r)
270 );
271
272 // Backend drain detection: all SA work complete, window FIFO truly empty,
273 // and no ready/partial block pending.
274 wire backend_idle = frame_done
275     && (state == ST_IDLE)
276     && !have_ready_block
277     && !partial_pending;
278
279 // Re-arm: one cycle after backend_idle, clear frame state for next transaction.
280 reg backend_idle_d;
281 assign frame_rearm = backend_idle && !backend_idle_d;
282 assign frame_rearm_out = frame_rearm;
283 assign frame_done_out = frame_done;
284
285 always @(posedge clk or negedge rst_n) begin
286     if (!rst_n)
287         backend_idle_d <= 1'b0;
288     else
289         backend_idle_d <= backend_idle;
290 end
291
292 // Frame ends when the beat count reaches eff_img_pixels or in_last is asserted.
293 // frame_rearm clears frame state after the backend drains, allowing back-to-back transactions.
294 always @(posedge clk or negedge rst_n) begin
295     if (!rst_n) begin
296         pix_cnt <= 32'd0;
297         frame_done <= 1'b0;
298     end else begin
299         if (frame_rearm) begin
300             pix_cnt <= 32'd0;
301             frame_done <= 1'b0;
302         end else if (pix_accept && !frame_done) begin
303             if (pix_cnt == (eff_img_pixels - 32'd1) || in_last) begin

```

```

304     frame_done <= 1'b1;
305     end
306     pix_cnt <= pix_cnt + 32'd1;
307     end
308     end
309     end
310
311     systolic_array_top #(
312         .DATA_W(DATA_W),
313         .ROWS (ROWS),
314         .COLS (COLS),
315         .DOT_K (DOT_K),
316         .ACC_W (ACC_W)
317     ) u_systolic_array_top (
318         .clk      (clk),
319         .rst_n    (rst_n),
320         .start_pulse(sa_start_pulse_q),
321         .mode_cfg  (sa_mode_cfg_q),
322         .valid_rows_cfg(sa_valid_rows_cfg_q),
323         .a_in_flat (sa_a_in_flat_q),
324         .b_in_flat (sa_b_in_flat_q),
325         .c_out_flat (c_out_raw_flat),
326         .done      (sa_done),
327         .col_stream_data_flat(c_out_col_stream_flat),
328         .col_stream_valid (c_out_col_valid),
329         .col_stream_last (c_out_col_last)
330     );
331
332     endmodule

```

## A.8.2 Convolution Engine Controller

*Controls weight loading, ready-bank consumption, systolic-array launch, feed, and drain sequencing.*

```
code/input/RTL/conv_core/conv_engine_ctrl.v
```

```

1  module conv_engine_ctrl #(
2      parameter integer DOT_K = 72,
3      parameter integer COLS  = 4,
4      parameter integer ROWS  = 16
5  )(
6      input wire [1:0] layer_sel,
7      input wire      clk,
8      input wire      rst_n,
9      input wire      wt_valid,
10     input wire      wt_last,
11     input wire      have_ready_block,
12     input wire      launch_from_A,
13     input wire [4:0] launch_valid_rows,
14     input wire      frame_rearm,
15     input wire      sa_done,

```

```

16  output wire      wt_ready,
17  output wire      weights_loaded,
18  output wire      wt_load_active,
19  output wire [8:0] wt_load_cnt,
20  output wire [1:0] state_dbg,
21  output wire [8:0] feed_cnt_dbg,
22  output wire      consume_ready_bank,
23  output wire      consume_bank_sel,
24  output wire      rd_en,
25  output wire      rd_bank,
26  output wire [8:0] rd_col,
27  output wire      wb_ring,
28  output wire      start_pulse
29 );
30
31  localparam [1:0] ST_IDLE      = 2'd0;
32  localparam [1:0] ST_PRERING  = 2'd1;
33  localparam [1:0] ST_FEED     = 2'd2;
34  localparam [1:0] ST_WAITDONE = 2'd3;
35
36  // Runtime decode: eff_dot_k, FEED_CYCLES, WT_LOAD_CYCLES
37  reg [8:0] eff_dot_k;
38  reg [8:0] eff_feed_last_base;
39  reg [8:0] eff_wt_load_cycles;
40  always @* begin
41      case (layer_sel)
42          2'b01: begin eff_dot_k = 9'd36; eff_feed_last_base = 9'd50; eff_wt_load_cycles = 9'd39; end //
36+16-1-1=50, 36+4-1=39
43          2'b10: begin eff_dot_k = 9'd72; eff_feed_last_base = 9'd86; eff_wt_load_cycles = 9'd75; end //
72+16-1-1=86, 72+4-1=75
44          default: begin eff_dot_k = 9'd9; eff_feed_last_base = 9'd23; eff_wt_load_cycles = 9'd12; end //
9+16-1-1=23, 9+4-1=12
45      endcase
46  end
47
48  reg [1:0] state_r;
49  reg [8:0] feed_cnt_r;
50  reg [8:0] wt_load_cnt_r;
51  reg      rd_en_r;
52  reg      rd_bank_r;
53  reg [8:0] rd_col_r;
54  reg      wb_ring_r;
55  reg      start_pulse_r;
56  reg      weights_loaded_r;
57  reg      wt_load_active_r;
58  reg [4:0] active_valid_rows_r;
59  reg [8:0] active_feed_last_r;
60
61  wire wt_fire = wt_valid && wt_ready;
62  wire have_launchable_block = have_ready_block && weights_loaded_r;
63  wire can_launch = have_launchable_block && !wt_fire;
64
65  assign wt_ready      = (state_r == ST_IDLE);
66  assign weights_loaded = weights_loaded_r;

```

```

67 assign wt_load_active    = wt_load_active_r;
68 assign wt_load_cnt      = wt_load_cnt_r;
69 assign state_dbg        = state_r;
70 assign feed_cnt_dbg     = feed_cnt_r;
71 assign rd_en            = rd_en_r;
72 assign rd_bank          = rd_bank_r;
73 assign rd_col           = rd_col_r;
74 assign wb_ring          = wb_ring_r;
75 assign start_pulse      = start_pulse_r;
76 assign consume_ready_bank =
77     ((state_r == ST_IDLE) && can_launch) ||
78     ((state_r == ST_FEED) && (feed_cnt_r == active_feed_last_r) && have_launchable_block) ||
79     ((state_r == ST_WAITDONE) && sa_done && have_launchable_block);
80 assign consume_bank_sel  = launch_from_A ? 1'b0 : 1'b1;
81
82 // Weight load protocol:
83 // 1) After wt_valid/wt_ready, weight_buffer accepts one new weight beat.
84 // 2) wt_last marks the last beat in the current weight group.
85 // 3) SA launch is allowed from a ready bank only after weights_loaded=1.
86 always @(posedge clk or negedge rst_n) begin
87     if (!rst_n) begin
88         weights_loaded_r <= 1'b0;
89         wt_load_active_r <= 1'b0;
90         wt_load_cnt_r    <= 9'd0;
91     end else if (frame_rearm) begin
92         weights_loaded_r <= 1'b0;
93         wt_load_active_r <= 1'b0;
94         wt_load_cnt_r    <= 9'd0;
95     end else if (wt_fire) begin
96         if (!wt_load_active_r) begin
97             weights_loaded_r <= wt_last;
98             wt_load_active_r <= !wt_last;
99             wt_load_cnt_r    <= 9'd1;
100        end else begin
101            weights_loaded_r <= wt_last;
102            wt_load_active_r <= !wt_last;
103            if (wt_load_cnt_r < eff_wt_load_cycles)
104                wt_load_cnt_r <= wt_load_cnt_r + 9'd1;
105            else
106                wt_load_cnt_r <= wt_load_cnt_r;
107        end
108    end
109 end
110
111 // SA launch control:
112 // ST_IDLE, the ST_FEED tail beat, or ST_WAITDONE can enter ST_PRERING
113 // when a ready block is available.
114 // The next ST_PRERING cycle asserts start_pulse and rd_en together.
115 // ST_FEED sends columns 0..DOT_K-1, then sends out-of-range columns to flush.
116 always @(posedge clk or negedge rst_n) begin
117     if (!rst_n) begin
118         state_r        <= ST_IDLE;
119         feed_cnt_r     <= 9'd0;
120         rd_col_r       <= 9'd0;

```

```

121     rd_en_r      <= 1'b0;
122     rd_bank_r    <= 1'b0;
123     wb_ring_r    <= 1'b0;
124     start_pulse_r <= 1'b0;
125     active_valid_rows_r <= 5'd0;
126     active_feed_last_r <= 9'd0;
127     end else if (frame_rearm) begin
128         state_r      <= ST_IDLE;
129         feed_cnt_r    <= 9'd0;
130         rd_col_r      <= 9'd0;
131         rd_en_r      <= 1'b0;
132         rd_bank_r    <= 1'b0;
133         wb_ring_r    <= 1'b0;
134         start_pulse_r <= 1'b0;
135         active_valid_rows_r <= 5'd0;
136         active_feed_last_r <= 9'd0;
137     end else begin
138         rd_en_r      <= 1'b0;
139         wb_ring_r    <= 1'b0;
140         start_pulse_r <= 1'b0;
141
142     case (state_r)
143         ST_IDLE: begin
144             feed_cnt_r <= 9'd0;
145             rd_col_r   <= 9'd0;
146             if (can_launch) begin
147                 rd_bank_r <= launch_from_A ? 1'b0 : 1'b1;
148                 active_valid_rows_r <= launch_valid_rows;
149                 if (launch_valid_rows > 5'd0)
150                     active_feed_last_r <= eff_dot_k + {{4{1'b0}}}, launch_valid_rows} - 9'd2;
151                 else
152                     active_feed_last_r <= eff_feed_last_base;
153                 wb_ring_r <= 1'b1;
154                 state_r   <= ST_PRERING;
155             end
156         end
157
158         ST_PRERING: begin
159             rd_col_r   <= 9'd0;
160             rd_en_r    <= 1'b1;
161             start_pulse_r <= 1'b1;
162             feed_cnt_r <= 9'd1;
163             state_r    <= ST_FEED;
164         end
165
166         ST_FEED: begin
167             rd_en_r <= 1'b1;
168             if (feed_cnt_r < eff_dot_k) begin
169                 rd_col_r <= feed_cnt_r;
170             end else begin
171                 rd_col_r <= 9'd511; // out-of-range -> get_byte returns 0 (flush)
172             end
173
174             if (feed_cnt_r == active_feed_last_r) begin

```

```

175     feed_cnt_r <= 9'd0;
176     rd_col_r   <= 9'd0;
177     if (have_launchable_block) begin
178         rd_bank_r <= launch_from_A ? 1'b0 : 1'b1;
179         active_valid_rows_r <= launch_valid_rows;
180         if (launch_valid_rows > 5'd0)
181             active_feed_last_r <= eff_dot_k + {{4{1'b0}}, launch_valid_rows} - 9'd2;
182         else
183             active_feed_last_r <= eff_feed_last_base;
184         wb_ring_r <= 1'b1;
185         state_r   <= ST_PRERING;
186     end else begin
187         state_r <= ST_WAITDONE;
188     end
189 end else begin
190     feed_cnt_r <= feed_cnt_r + 9'd1;
191 end
192 end
193
194 ST_WAITDONE: begin
195     if (sa_done) begin
196         if (have_launchable_block) begin
197             rd_bank_r <= launch_from_A ? 1'b0 : 1'b1;
198             active_valid_rows_r <= launch_valid_rows;
199             if (launch_valid_rows > 5'd0)
200                 active_feed_last_r <= eff_dot_k + {{4{1'b0}}, launch_valid_rows} - 9'd2;
201             else
202                 active_feed_last_r <= eff_feed_last_base;
203             wb_ring_r <= 1'b1;
204             state_r   <= ST_PRERING;
205         end else begin
206             state_r <= ST_IDLE;
207         end
208     end
209 end
210
211 default: begin
212     state_r <= ST_IDLE;
213 end
214 endcase
215 end
216 end
217
218 endmodule

```

### A.8.3 Convolution Buffer

*Stores incoming convolution windows, tracks ready banks, and feeds rows into the systolic array.*

```

1 // Conv_Buffer - dual-bank stripe buffer for conv_top.
2 //
3 // Storage: flat 8-bit regs (2 x 432 entries) - cheap write-enable per register.
4 // Read path: hardwired concat into 24 x 144-bit "group" wires (zero gate cost),
5 // then a single 24:1 MUX + 3:1 barrel shift extract 16 consecutive bytes.
6 // Address decode: counter-based (ch_r, dx_r, dy_r) advances with rd_col - no
7 // runtime division/modulo.
8 // zero_fill_tail: eliminated; read path masks out unfilled columns via fill_cols.
9 //
10 // Synthesis: ~84K um^2 (vs ~114K baseline with 16x 432:1 MUX + div/mod).
11
12 module Conv_Buffer #(
13     parameter integer DATA_W = 8,
14     parameter integer ROWS    = 16,
15     parameter integer DOT_K   = 72,
16     parameter integer C_IN    = 8
17 ) (
18     input wire [1:0]          layer_sel,
19     input wire                clk,
20     input wire                rst_n,
21     input wire                frame_rearm,
22     input wire                pix3_valid,
23     input wire signed [C_IN*DATA_W-1:0] row0_pix,
24     input wire signed [C_IN*DATA_W-1:0] row1_pix,
25     input wire signed [C_IN*DATA_W-1:0] row2_pix,
26     input wire                rd_active,
27     input wire                rd_en,
28     input wire                rd_bank,
29     input wire [8:0]          rd_col,
30     input wire                consume_ready_bank,
31     input wire                consume_bank_sel,
32     output wire               bank_can_write,
33     output wire               have_ready_block,
34     output wire               launch_from_A,
35     output wire [4:0]         valid_rows_A,
36     output wire [4:0]         valid_rows_B,
37     output wire               partial_pending,
38     output wire signed [ROWS*DATA_W-1:0] raw_firstcol_flat
39 );
40
41 // -----
42 // Constants
43 // -----
44 localparam integer PIX_W      = C_IN * DATA_W;
45 localparam integer STRIPE_ROWS = 3;
46 localparam integer STRIPE_COLS = 18;
47 localparam integer STRIPE_STEP = 16;
48 localparam integer STRIPE_BYTES = STRIPE_ROWS * STRIPE_COLS * C_IN;
49 localparam integer WORD_W      = STRIPE_COLS * DATA_W; // 144
50 localparam integer MAX_WORDS   = STRIPE_ROWS * C_IN; // 24
51
52 // -----

```

```

53 // Flat 8-bit storage (same as Conv_Buffer - cheap writes)
54 // -----
55 reg signed [7:0] bank_A [0:STRIPE_BYTES-1];
56 reg signed [7:0] bank_B [0:STRIPE_BYTES-1];
57
58 // -----
59 // Control registers
60 // -----
61 reg [6:0] x_pos_r;
62 reg [6:0] fill_base_x_r;
63 reg [4:0] fill_cols_r;
64 reg      wr_bank_r;
65 reg      ready_A_r;
66 reg      ready_B_r;
67 reg [4:0] valid_rows_A_r;
68 reg [4:0] valid_rows_B_r;
69 reg      seed_pending_r;
70 reg      seed_src_bank_r;
71 // Per-bank fill column count (replaces zero_fill_tail)
72 reg [4:0] fill_cols_A_r;
73 reg [4:0] fill_cols_B_r;
74
75 // -----
76 // Layer-dependent parameters
77 // -----
78 reg [6:0] eff_w;
79 reg [3:0] eff_c;
80 reg [6:0] eff_out_w;
81 reg [8:0] eff_dot_k;
82
83 always @* begin
84     case (layer_sel)
85         2'b01: begin
86             eff_w      = 7'd31;
87             eff_c      = 4'd4;
88             eff_out_w = 7'd29;
89             eff_dot_k = 9'd36;
90         end
91         2'b10: begin
92             eff_w      = 7'd14;
93             eff_c      = 4'd8;
94             eff_out_w = 7'd12;
95             eff_dot_k = 9'd72;
96         end
97         default: begin
98             eff_w      = 7'd64;
99             eff_c      = 4'd1;
100            eff_out_w = 7'd62;
101            eff_dot_k = 9'd9;
102        end
103     endcase
104 end
105
106 // -----

```

```

107 // TB-visible diagnostic wires (same as Conv_Buffer)
108 // -----
109 wire stripe_window_open =
110     (x_pos_r >= fill_base_x_r) && (x_pos_r < (fill_base_x_r + 7'd18));
111 wire write_bank_blocked = ready_A_r && ready_B_r;
112 wire write_read_conflict = rd_active && (wr_bank_r == rd_bank);
113 wire input_stall = write_bank_blocked || write_read_conflict;
114 wire accepting_input_col = pix3_valid && !input_stall && stripe_window_open;
115 wire signed [DOT_K*DATA_W-1:0] write_sample_flat =
116     {{(DOT_K*DATA_W - 3*PIX_W){1'b0}}, row2_pix, row1_pix, row0_pix};
117
118 // -----
119 // Functions
120 // -----
121 function integer stripe_idx;
122     input integer row_idx;
123     input integer col_idx;
124     input integer ch_idx;
125     begin
126         stripe_idx = ((row_idx * STRIPE_COLS) + col_idx) * C_IN + ch_idx;
127     end
128 endfunction
129
130 function [4:0] calc_valid_rows;
131     input [6:0] base_x;
132     input [6:0] out_w;
133     reg [6:0] remaining;
134     begin
135         if (base_x >= out_w)
136             calc_valid_rows = 5'd0;
137         else begin
138             remaining = out_w - base_x;
139             if (remaining >= ROWS[6:0])
140                 calc_valid_rows = ROWS[4:0];
141             else
142                 calc_valid_rows = {1'b0, remaining[3:0]};
143         end
144     end
145 endfunction
146
147 // -----
148 // Write tasks (from Conv_Buffer - simple flat-reg writes)
149 // -----
150 integer col_off_i;
151 reg [4:0] new_valid_rows;
152 integer real_cols_i;
153
154 task write_input_column;
155     input dst_bank;
156     input integer col_idx;
157     integer ch_i;
158     begin
159         for (ch_i = 0; ch_i < C_IN; ch_i = ch_i + 1) begin
160             if (dst_bank) begin

```

```

161     bank_B[stripe_idx(0, col_idx, ch_i)] <= row0_pix[ch_i*DATA_W +: DATA_W];
162     bank_B[stripe_idx(1, col_idx, ch_i)] <= row1_pix[ch_i*DATA_W +: DATA_W];
163     bank_B[stripe_idx(2, col_idx, ch_i)] <= row2_pix[ch_i*DATA_W +: DATA_W];
164     end else begin
165         bank_A[stripe_idx(0, col_idx, ch_i)] <= row0_pix[ch_i*DATA_W +: DATA_W];
166         bank_A[stripe_idx(1, col_idx, ch_i)] <= row1_pix[ch_i*DATA_W +: DATA_W];
167         bank_A[stripe_idx(2, col_idx, ch_i)] <= row2_pix[ch_i*DATA_W +: DATA_W];
168     end
169 end
170 end
171 endtask
172
173 task copy_overlap_cols;
174     input dst_bank;
175     input src_bank;
176     integer row_i;
177     integer ch_i;
178     begin
179         for (row_i = 0; row_i < STRIPE_ROWS; row_i = row_i + 1) begin
180             for (ch_i = 0; ch_i < C_IN; ch_i = ch_i + 1) begin
181                 if (dst_bank) begin
182                     if (src_bank) begin
183                         bank_B[stripe_idx(row_i, 0, ch_i)] <= bank_B[stripe_idx(row_i, 16, ch_i)];
184                         bank_B[stripe_idx(row_i, 1, ch_i)] <= bank_B[stripe_idx(row_i, 17, ch_i)];
185                     end else begin
186                         bank_B[stripe_idx(row_i, 0, ch_i)] <= bank_A[stripe_idx(row_i, 16, ch_i)];
187                         bank_B[stripe_idx(row_i, 1, ch_i)] <= bank_A[stripe_idx(row_i, 17, ch_i)];
188                     end
189                 end else begin
190                     if (src_bank) begin
191                         bank_A[stripe_idx(row_i, 0, ch_i)] <= bank_B[stripe_idx(row_i, 16, ch_i)];
192                         bank_A[stripe_idx(row_i, 1, ch_i)] <= bank_B[stripe_idx(row_i, 17, ch_i)];
193                     end else begin
194                         bank_A[stripe_idx(row_i, 0, ch_i)] <= bank_A[stripe_idx(row_i, 16, ch_i)];
195                         bank_A[stripe_idx(row_i, 1, ch_i)] <= bank_A[stripe_idx(row_i, 17, ch_i)];
196                     end
197                 end
198             end
199         end
200     end
201 endtask
202
203 // zero_fill_tail: ELIMINATED - read-side fill_cols mask handles this
204
205 task mark_ready_bank;
206     input dst_bank;
207     input [4:0] valid_rows;
208     input [4:0] fill_cols;
209     begin
210         if (dst_bank) begin
211             ready_B_r      <= (valid_rows != 5'd0);
212             valid_rows_B_r <= valid_rows;
213             fill_cols_B_r  <= fill_cols;
214         end else begin

```

```

215     ready_A_r      <= (valid_rows != 5'd0);
216     valid_rows_A_r <= valid_rows;
217     fill_cols_A_r <= fill_cols;
218     end
219     end
220 endtask
221
222 // -----
223 // Output assignments
224 // -----
225 assign bank_can_write  = !input_stall;
226 assign have_ready_block = ready_A_r || ready_B_r;
227 assign launch_from_A   = ready_A_r;
228 assign valid_rows_A    = valid_rows_A_r;
229 assign valid_rows_B    = valid_rows_B_r;
230 assign partial_pending = seed_pending_r || (fill_cols_r != 5'd0);
231
232 // =====
233 // READ PATH - counter-based decode + hardwired concatenation
234 // =====
235
236 // Hardwired 144-bit group wires: zero gate cost, pure wiring.
237 // group_X[dy*C_IN+ch] concatenates 18 flat regs into a 144-bit word.
238 wire [WORD_W-1:0] group_A [0:MAX_WORDS-1];
239 wire [WORD_W-1:0] group_B [0:MAX_WORDS-1];
240
241 genvar g_dy, g_ch, g_col;
242 generate
243     for (g_dy = 0; g_dy < STRIPE_ROWS; g_dy = g_dy + 1) begin : gen_dy
244         for (g_ch = 0; g_ch < C_IN; g_ch = g_ch + 1) begin : gen_ch
245             for (g_col = 0; g_col < STRIPE_COLS; g_col = g_col + 1) begin : gen_col
246                 localparam integer WADDR = g_dy * C_IN + g_ch;
247                 localparam integer SIDX  = (g_dy * STRIPE_COLS + g_col) * C_IN + g_ch;
248                 assign group_A[WADDR][g_col*DATA_W +: DATA_W] = bank_A[SIDX];
249                 assign group_B[WADDR][g_col*DATA_W +: DATA_W] = bank_B[SIDX];
250             end
251         end
252     end
253 endgenerate
254
255 // -----
256 // Counter-based im2col decode (replaces runtime div/mod on rd_col)
257 // -----
258 reg [3:0] ch_r;
259 reg [1:0] dx_r;
260 reg [1:0] dy_r;
261
262 always @(posedge clk or negedge rst_n) begin
263     if (!rst_n) begin
264         ch_r <= 4'd0;
265         dx_r <= 2'd0;
266         dy_r <= 2'd0;
267     end else if (frame_rearm) begin
268         ch_r <= 4'd0;

```

```

269     dx_r <= 2'd0;
270     dy_r <= 2'd0;
271 end else if (rd_en && rd_col < eff_dot_k) begin
272     if (rd_col == 9'd0) begin
273         if (eff_c == 4'd1) begin
274             ch_r <= 4'd0;
275             dx_r <= 2'd1;
276             dy_r <= 2'd0;
277         end else begin
278             ch_r <= 4'd1;
279             dx_r <= 2'd0;
280             dy_r <= 2'd0;
281         end
282     end else begin
283         if (ch_r == eff_c - 4'd1) begin
284             ch_r <= 4'd0;
285             if (dx_r == 2'd2) begin
286                 dx_r <= 2'd0;
287                 dy_r <= dy_r + 2'd1;
288             end else begin
289                 dx_r <= dx_r + 2'd1;
290             end
291         end else begin
292             ch_r <= ch_r + 4'd1;
293         end
294     end
295 end else if (!rd_en) begin
296     ch_r <= 4'd0;
297     dx_r <= 2'd0;
298     dy_r <= 2'd0;
299 end
300 end
301
302 // Current-cycle combinational decode
303 reg [3:0] cur_ch;
304 reg [1:0] cur_dx;
305 reg [1:0] cur_dy;
306 always @* begin
307     if (!rd_en || rd_col >= eff_dot_k) begin
308         cur_ch = 4'd0;
309         cur_dx = 2'd0;
310         cur_dy = 2'd0;
311     end else if (rd_col == 9'd0) begin
312         cur_ch = 4'd0;
313         cur_dx = 2'd0;
314         cur_dy = 2'd0;
315     end else begin
316         cur_ch = ch_r;
317         cur_dx = dx_r;
318         cur_dy = dy_r;
319     end
320 end
321
322 // -----

```

```

323 // Read MUX: 24:1 group select -> shift extract -> fill_cols mask
324 // -----
325 reg [WORD_W-1:0] rd_word;
326 reg [4:0] selected_valid_rows;
327 reg [4:0] rd_fill_cols;
328 reg signed [ROWS*DATA_W-1:0] raw_col_r;
329 integer ri;
330
331 // Word address - always use C_IN stride to match generate block indexing
332 reg [4:0] rd_word_addr;
333 always @* begin
334     rd_word_addr = {cur_dy, 3'b000} + {1'b0, cur_ch}; // dy * 8 + ch
335 end
336
337 always @* begin
338     selected_valid_rows = rd_bank ? valid_rows_B_r : valid_rows_A_r;
339     rd_fill_cols        = rd_bank ? fill_cols_B_r  : fill_cols_A_r;
340     rd_word             = rd_bank ? group_B[rd_word_addr] : group_A[rd_word_addr];
341
342     raw_col_r = {(ROWS*DATA_W){1'b0}};
343     if (rd_en && rd_col < eff_dot_k) begin
344         for (ri = 0; ri < ROWS; ri = ri + 1) begin
345             if (ri < selected_valid_rows && (ri + cur_dx) < rd_fill_cols)
346                 raw_col_r[ri*DATA_W +: DATA_W] =
347                     rd_word[(ri + cur_dx) * DATA_W +: DATA_W];
348         end
349     end
350 end
351
352 assign raw_firstcol_flat = raw_col_r;
353
354 // -----
355 // Main sequential block - stripe fill FSM
356 // -----
357 always @(posedge clk or negedge rst_n) begin
358     if (!rst_n) begin
359         x_pos_r          <= 7'd0;
360         fill_base_x_r    <= 7'd0;
361         fill_cols_r      <= 5'd0;
362         wr_bank_r        <= 1'b0;
363         ready_A_r        <= 1'b0;
364         ready_B_r        <= 1'b0;
365         valid_rows_A_r   <= 5'd0;
366         valid_rows_B_r   <= 5'd0;
367         seed_pending_r   <= 1'b0;
368         seed_src_bank_r  <= 1'b0;
369         fill_cols_A_r    <= 5'd18;
370         fill_cols_B_r    <= 5'd18;
371     end else if (frame_rearm) begin
372         x_pos_r          <= 7'd0;
373         fill_base_x_r    <= 7'd0;
374         fill_cols_r      <= 5'd0;
375         wr_bank_r        <= 1'b0;
376         ready_A_r        <= 1'b0;

```

```

377     ready_B_r      <= 1'b0;
378     valid_rows_A_r <= 5'd0;
379     valid_rows_B_r <= 5'd0;
380     seed_pending_r <= 1'b0;
381     seed_src_bank_r <= 1'b0;
382     fill_cols_A_r  <= 5'd18;
383     fill_cols_B_r  <= 5'd18;
384 end else begin
385     if (consume_ready_bank) begin
386         if (consume_bank_sel == 1'b0)
387             ready_A_r <= 1'b0;
388         else
389             ready_B_r <= 1'b0;
390     end
391
392     if (seed_pending_r && !pix3_valid && !input_stall) begin
393         copy_overlap_cols(wr_bank_r, seed_src_bank_r);
394         fill_cols_r      <= 5'd2;
395         seed_pending_r <= 1'b0;
396     end else if (pix3_valid && !input_stall) begin
397         if (seed_pending_r) begin
398             copy_overlap_cols(wr_bank_r, seed_src_bank_r);
399             fill_cols_r      <= 5'd2;
400             seed_pending_r <= 1'b0;
401         end
402
403         col_off_i = x_pos_r - fill_base_x_r;
404         if (stripe_window_open) begin
405             write_input_column(wr_bank_r, col_off_i);
406             fill_cols_r <= col_off_i[4:0] + 5'd1;
407         end
408
409         if (x_pos_r == eff_w - 7'd1) begin
410             if (fill_base_x_r < eff_out_w) begin
411                 real_cols_i = x_pos_r - fill_base_x_r + 1;
412                 // zero_fill_tail REMOVED - read-side fill_cols mask handles it
413                 new_valid_rows = calc_valid_rows(fill_base_x_r, eff_out_w);
414                 mark_ready_bank(wr_bank_r, new_valid_rows, real_cols_i[4:0]);
415                 wr_bank_r <= ~wr_bank_r;
416             end
417
418             x_pos_r      <= 7'd0;
419             fill_base_x_r <= 7'd0;
420             fill_cols_r  <= 5'd0;
421             seed_pending_r <= 1'b0;
422         end else if (x_pos_r == fill_base_x_r + 7'd17) begin
423             new_valid_rows = calc_valid_rows(fill_base_x_r, eff_out_w);
424             mark_ready_bank(wr_bank_r, new_valid_rows, 5'd18);
425
426             seed_src_bank_r <= wr_bank_r;
427             wr_bank_r      <= ~wr_bank_r;
428             fill_base_x_r  <= fill_base_x_r + STRIPE_STEP[6:0];
429             fill_cols_r    <= 5'd0;
430             seed_pending_r <= ((fill_base_x_r + STRIPE_STEP[6:0]) < eff_out_w);

```

```

431         x_pos_r      <= x_pos_r + 7'd1;
432     end else begin
433         x_pos_r <= x_pos_r + 7'd1;
434     end
435 end
436 end
437 end
438
439 endmodule

```

#### A.8.4 Line Buffer

*Implements the layer-dependent row delay needed to assemble 3-row convolution windows.*

code/input/RTL/conv\_core/Line\_Buffer.v

```

1  `timescale 1ns/1ns
2
3  // Line_Buffer: runtime-compatible version.
4  // Fixed 124-byte shift register (992 bits).
5  // layer_sel selects effective depth and stride:
6  //  L1 (2'b00): stride=1 byte, depth=64 -> tap at byte 63
7  //  L2 (2'b01): stride=4 bytes, depth=31 -> tap at byte 120..123
8  //  L3 (2'b10): stride=8 bytes, depth=14 -> tap at byte 104..111
9  module Line_Buffer #(
10     parameter integer W = 64,
11     parameter integer DW = 64 // max pixel bit-width
12 ) (
13     input wire      [1:0] layer_sel,
14     input wire      clk,
15     input wire      rst_n,
16     input wire      frame_rearm, // clears the shift register and fill counter
17     input wire      in_valid,
18     input wire signed [DW-1:0] in_data,
19     output wire      out_valid,
20     output wire signed [DW-1:0] out_data // selected shift-register tap for layer_sel
21 );
22
23     localparam integer N = 124; // total bytes
24
25     reg [7:0] shreg [0:N-1];
26
27     // Runtime parameters
28     reg [6:0] eff_w; // effective depth (pixels)
29     reg [3:0] eff_stride; // bytes per pixel
30     always @* begin
31         case (layer_sel)
32             2'b01: begin eff_w = 7'd31; eff_stride = 4'd4; end // L2
33             2'b10: begin eff_w = 7'd14; eff_stride = 4'd8; end // L3
34             default: begin eff_w = 7'd64; eff_stride = 4'd1; end // L1
35         endcase

```

```

36     end
37
38     // Fill counter
39     reg [6:0] fill_cnt;
40     assign out_valid = in_valid && (fill_cnt == eff_w);
41
42     // Output tap: read eff_stride bytes from the oldest position.
43     // Use byte-level MUX to avoid negative replication for any DW.
44     reg signed [DW-1:0] out_r;
45     integer ob;
46     always @* begin
47         out_r = {DW{1'b0}};
48         case (layer_sel)
49             2'b01: begin // L2: 4 bytes at [120..123]
50                 for (ob = 0; ob < DW/8 && ob < 4; ob = ob + 1)
51                     out_r[ob*8 +: 8] = shreg[120 + ob];
52             end
53             2'b10: begin // L3: 8 bytes at [104..111]
54                 for (ob = 0; ob < DW/8 && ob < 8; ob = ob + 1)
55                     out_r[ob*8 +: 8] = shreg[104 + ob];
56             end
57             default: begin // L1: 1 byte at [63]
58                 out_r[7:0] = shreg[63];
59             end
60         endcase
61     end
62     assign out_data = out_r;
63
64     // Shift register: on each in_valid, shift by eff_stride bytes.
65     integer i;
66     always @(posedge clk or negedge rst_n) begin
67         if (!rst_n) begin
68             fill_cnt <= 7'd0;
69             for (i = 0; i < N; i = i + 1)
70                 shreg[i] <= 8'd0;
71         end else if (frame_rearm) begin
72             fill_cnt <= 7'd0;
73             for (i = 0; i < N; i = i + 1)
74                 shreg[i] <= 8'd0;
75         end else if (in_valid) begin
76             case (eff_stride)
77                 4'd4: begin
78                     for (i = N-1; i >= 4; i = i - 1) shreg[i] <= shreg[i-4];
79                     shreg[3] <= in_data[31:24];
80                     shreg[2] <= in_data[23:16];
81                     shreg[1] <= in_data[15:8];
82                     shreg[0] <= in_data[7:0];
83                 end
84                 4'd8: begin
85                     for (i = N-1; i >= 8; i = i - 1) shreg[i] <= shreg[i-8];
86                     for (ob = 0; ob < 8; ob = ob + 1) begin
87                         if (ob < DW/8)
88                             shreg[ob] <= in_data[ob*8 +: 8];
89                         else

```

```

90         shreg[ob] <= 8'd0;
91     end
92 end
93 default: begin // stride=1
94     for (i = N-1; i >= 1; i = i - 1) shreg[i] <= shreg[i-1];
95     shreg[0] <= in_data[7:0];
96 end
97 endcase
98
99     if (fill_cnt < eff_w)
100         fill_cnt <= fill_cnt + 7'd1;
101     end
102 end
103
104 endmodule

```

### A.8.5 Input Row Aligner

*Unpacks and aligns incoming pixels into the row streams consumed by the convolution buffer.*

```
code/input/RTL/conv_core/input_row_aligner.v
```

```

1  `timescale 1ns/1ns
2
3  module input_row_aligner #(
4      parameter integer W    = 64,
5      parameter integer DW   = 8,
6      parameter integer C_IN = 8
7  )(
8      input  wire          [1:0]    layer_sel,
9      input  wire          clk,
10     input  wire          rst_n,
11     input  wire          frame_rearm,
12     input  wire          in_valid,
13     input  wire          [31:0]   in_data,
14     input  wire          [3:0]   in_byte_en,
15     output wire          pix3_valid,
16     output wire signed  [C_IN*DW-1:0] row0_out,
17     output wire signed  [C_IN*DW-1:0] row1_out,
18     output wire signed  [C_IN*DW-1:0] row2_out
19 );
20
21     localparam integer PIX_W = C_IN * DW;
22
23     reg [31:0] saved_lo;
24     reg half_r;
25     reg ser_valid_r;
26     reg signed [PIX_W-1:0] ser_data_r;
27
28     wire [31:0] in_masked = { in_byte_en[3] ? in_data[31:24] : 8'd0,
29                             in_byte_en[2] ? in_data[23:16] : 8'd0,

```

```

30             in_byte_en[1] ? in_data[15:8] : 8'd0,
31             in_byte_en[0] ? in_data[7:0]  : 8'd0 };
32
33 integer pb;
34 always @* begin
35     ser_valid_r = 1'b0;
36     ser_data_r  = {PIX_W{1'b0}};
37     case (layer_sel)
38         2'b00: begin
39             ser_valid_r  = in_valid;
40             ser_data_r[7:0] = in_masked[7:0];
41         end
42         2'b01: begin
43             ser_valid_r = in_valid;
44             for (pb = 0; pb < 4 && pb < PIX_W/8; pb = pb + 1)
45                 ser_data_r[pb*8 +: 8] = in_masked[pb*8 +: 8];
46         end
47         default: begin
48             ser_valid_r = in_valid && half_r;
49             for (pb = 0; pb < 4 && pb < PIX_W/8; pb = pb + 1)
50                 ser_data_r[pb*8 +: 8] = saved_lo[pb*8 +: 8];
51             for (pb = 4; pb < 8 && pb < PIX_W/8; pb = pb + 1)
52                 ser_data_r[pb*8 +: 8] = in_masked[(pb-4)*8 +: 8];
53         end
54     endcase
55 end
56
57 always @(posedge clk or negedge rst_n) begin
58     if (!rst_n) begin
59         saved_lo <= 32'd0;
60         half_r   <= 1'b0;
61     end else if (frame_rearm) begin
62         saved_lo <= 32'd0;
63         half_r   <= 1'b0;
64     end else if (in_valid) begin
65         case (layer_sel)
66             2'b10: begin
67                 if (!half_r)
68                     saved_lo <= in_masked;
69                 half_r <= ~half_r;
70             end
71             default: begin
72                 half_r <= 1'b0;
73             end
74         endcase
75     end
76 end
77
78 wire signed [PIX_W-1:0] row0;
79 wire signed [PIX_W-1:0] row1;
80 wire signed [PIX_W-1:0] row2;
81 wire row0_valid;
82 wire row1_valid;
83 wire row2_valid;

```

```

84
85 assign row0      = ser_data_r;
86 assign row0_valid = ser_valid_r;
87
88 Line_Buffer #(
89     .W(W),
90     .DW(PIX_W)
91 ) Line_Buffer_row1 (
92     .layer_sel (layer_sel),
93     .clk       (clk),
94     .rst_n     (rst_n),
95     .frame_rearm(frame_rearm),
96     .in_valid  (row0_valid),
97     .in_data   (row0),
98     .out_valid (row1_valid),
99     .out_data  (row1)
100 );
101
102 Line_Buffer #(
103     .W(W),
104     .DW(PIX_W)
105 ) Line_Buffer_row2 (
106     .layer_sel (layer_sel),
107     .clk       (clk),
108     .rst_n     (rst_n),
109     .frame_rearm(frame_rearm),
110     .in_valid  (row1_valid),
111     .in_data   (row1),
112     .out_valid (row2_valid),
113     .out_data  (row2)
114 );
115
116 assign pix3_valid = row2_valid & row1_valid & row0_valid;
117 assign row0_out   = row0;
118 assign row1_out   = row1;
119 assign row2_out   = row2;
120
121 endmodule

```

## A.8.6 Systolic-Array Skew Feeder

*Applies row-wise skewing so activation rows arrive at the systolic array with the required timing.*

```
code/input/RTL/conv_core/sa_skew_feeder.v
```

```

1 module sa_skew_feeder #(
2     parameter integer ROWS = 16
3 )(
4     input wire          clk,
5     input wire          rst_n,
6     input wire          frame_rearm,

```

```

7     input wire          rd_en,
8     input wire signed [ROWS*8-1:0] in_firstcol_flat,
9     output wire signed [ROWS*8-1:0] out_firstcol_flat
10 );
11
12 // Skew-only stage: row r is delayed by r cycles to create the systolic
13 // wavefront. The large full-window banks now live in Conv_Buffer instead.
14 reg signed [7:0] delay [0:119];
15 reg          rd_en_d;
16
17 function integer base;
18     input integer r;
19     begin
20         base = (r * (r - 1)) / 2;
21     end
22 endfunction
23
24 wire rd_start = rd_en & ~rd_en_d;
25
26 wire signed [7:0] raw0 = in_firstcol_flat[7:0];
27 wire signed [7:0] raw1 = in_firstcol_flat[15:8];
28 wire signed [7:0] raw2 = in_firstcol_flat[23:16];
29 wire signed [7:0] raw3 = in_firstcol_flat[31:24];
30 wire signed [7:0] raw4 = in_firstcol_flat[39:32];
31 wire signed [7:0] raw5 = in_firstcol_flat[47:40];
32 wire signed [7:0] raw6 = in_firstcol_flat[55:48];
33 wire signed [7:0] raw7 = in_firstcol_flat[63:56];
34 wire signed [7:0] raw8 = in_firstcol_flat[71:64];
35 wire signed [7:0] raw9 = in_firstcol_flat[79:72];
36 wire signed [7:0] raw10 = in_firstcol_flat[87:80];
37 wire signed [7:0] raw11 = in_firstcol_flat[95:88];
38 wire signed [7:0] raw12 = in_firstcol_flat[103:96];
39 wire signed [7:0] raw13 = in_firstcol_flat[111:104];
40 wire signed [7:0] raw14 = in_firstcol_flat[119:112];
41 wire signed [7:0] raw15 = in_firstcol_flat[127:120];
42
43 integer x;
44 integer s;
45 integer b;
46 always @(posedge clk or negedge rst_n) begin
47     if (!rst_n) begin
48         rd_en_d <= 1'b0;
49         for (x = 0; x < 120; x = x + 1)
50             delay[x] <= 8'h00;
51     end else if (frame_rearm) begin
52         rd_en_d <= 1'b0;
53         for (x = 0; x < 120; x = x + 1)
54             delay[x] <= 8'h00;
55     end else begin
56         rd_en_d <= rd_en;
57         if (rd_en) begin
58             for (x = 1; x <= 15; x = x + 1) begin
59                 b = base(x);
60                 if (rd_start) begin

```

```

61         for (s = x - 1; s >= 1; s = s - 1)
62             delay[b + s] <= 8'h00;
63     end else begin
64         for (s = x - 1; s >= 1; s = s - 1)
65             delay[b + s] <= delay[b + s - 1];
66     end
67
68     case (x)
69         1: delay[b] <= raw1;
70         2: delay[b] <= raw2;
71         3: delay[b] <= raw3;
72         4: delay[b] <= raw4;
73         5: delay[b] <= raw5;
74         6: delay[b] <= raw6;
75         7: delay[b] <= raw7;
76         8: delay[b] <= raw8;
77         9: delay[b] <= raw9;
78         10: delay[b] <= raw10;
79         11: delay[b] <= raw11;
80         12: delay[b] <= raw12;
81         13: delay[b] <= raw13;
82         14: delay[b] <= raw14;
83         15: delay[b] <= raw15;
84         default: delay[b] <= 8'h00;
85     endcase
86     end
87 end
88 end
89 end
90
91 wire signed [7:0] out0 = rd_en ? raw0 : 8'h00;
92 wire signed [7:0] out1 = delay[base(1) + 0];
93 wire signed [7:0] out2 = delay[base(2) + 1];
94 wire signed [7:0] out3 = delay[base(3) + 2];
95 wire signed [7:0] out4 = delay[base(4) + 3];
96 wire signed [7:0] out5 = delay[base(5) + 4];
97 wire signed [7:0] out6 = delay[base(6) + 5];
98 wire signed [7:0] out7 = delay[base(7) + 6];
99 wire signed [7:0] out8 = delay[base(8) + 7];
100 wire signed [7:0] out9 = delay[base(9) + 8];
101 wire signed [7:0] out10 = delay[base(10) + 9];
102 wire signed [7:0] out11 = delay[base(11) + 10];
103 wire signed [7:0] out12 = delay[base(12) + 11];
104 wire signed [7:0] out13 = delay[base(13) + 12];
105 wire signed [7:0] out14 = delay[base(14) + 13];
106 wire signed [7:0] out15 = delay[base(15) + 14];
107
108 assign out_firstcol_flat = {
109     out15, out14, out13, out12, out11, out10, out9, out8,
110     out7, out6, out5, out4, out3, out2, out1, out0
111 };
112
113 endmodule

```

## A.8.7 Systolic Array Top

Implements the 16-by-4 MAC array and per-column output stream used by convolution and FC-compatible modes.

```
code/input/RTL/conv_core/systolic_array_top.v

1 // =====
2 // 16x4 systolic array top.
3 // - Keeps the flat input, flat output, and column-stream interface.
4 // - Uses mode_cfg and runtime cur_k for layer-dependent MAC length.
5 // - Conv1-compatible mode is selected with mode_cfg=2'b00.
6 // =====
7
8 module systolic_array_top #(
9     parameter integer DATA_W = 8,
10    parameter integer ROWS    = 16,
11    parameter integer COLS    = 4,
12    parameter integer DOT_K   = 9, // legacy compatibility only
13    parameter integer ACC_W   = 23
14 ) (
15     input wire clk,
16     input wire rst_n,
17     input wire start_pulse,
18     input wire [1:0] mode_cfg,
19     input wire [4:0] valid_rows_cfg,
20
21     // Flat inputs:
22     // a_in_flat[(r+1)*DATA_W-1 : r*DATA_W] is the left-edge input for row r.
23     input wire signed [ROWS*DATA_W-1:0] a_in_flat,
24     // b_in_flat[(c+1)*DATA_W-1 : c*DATA_W] is the top-edge input for column c.
25     input wire signed [COLS*DATA_W-1:0] b_in_flat,
26
27     // Flat outputs:
28     // idx = r*COLS + c
29     // c_out_flat[(idx+1)*ACC_W-1 : idx*ACC_W] maps to c_out[r][c].
30     output wire signed [ROWS*COLS*ACC_W-1:0] c_out_flat,
31     output wire done,
32
33     // Per-column real-time stream output.
34     output reg signed [COLS*ACC_W-1:0] col_stream_data_flat,
35     output reg [COLS-1:0] col_stream_valid,
36     output reg [COLS-1:0] col_stream_last
37 );
38
39 localparam [1:0] MODE_CONV1 = 2'b00;
40 localparam [1:0] MODE_CONV2 = 2'b01;
41 localparam [1:0] MODE_CONV3 = 2'b10;
42 localparam [1:0] MODE_FC    = 2'b11;
43 localparam integer MAX_DOT_K = 288;
44 localparam integer CNT_W     = $clog2(MAX_DOT_K + 1);
45 localparam integer START_STAGES = ROWS + COLS - 1;
46
```

```

47 integer k;
48 integer sc;
49 integer data_lo;
50
51 reg [START_STAGES-1:0] start_pipe;
52 reg [CNT_W-1:0] cur_k;
53 wire fc_mode;
54
55 reg [COLS-1:0] stream_active;
56 reg [COLS-1:0] stream_arm;
57 reg [7:0] stream_row [0:COLS-1];
58 reg [4:0] stream_rows_limit [0:COLS-1];
59 reg [4:0] arm_rows_limit [0:COLS-1];
60 reg [4:0] col_rows_cfg [0:COLS-1];
61 reg [4:0] valid_rows_pipe [0:COLS-1];
62 reg valid_rows_pipe_v [0:COLS-1];
63
64 wire signed [DATA_W-1:0] a_in [0:ROWS-1];
65 wire signed [DATA_W-1:0] b_in [0:COLS-1];
66 wire start_ij [0:ROWS-1][0:COLS-1];
67 wire finish_ij [0:ROWS-1][0:COLS-1];
68
69 wire signed [DATA_W-1:0] pe_left [0:ROWS-1][0:COLS-1];
70 wire signed [DATA_W-1:0] pe_up [0:ROWS-1][0:COLS-1];
71 wire signed [DATA_W-1:0] pe_right [0:ROWS-1][0:COLS-1];
72 wire signed [DATA_W-1:0] pe_bottom [0:ROWS-1][0:COLS-1];
73 wire signed [ACC_W-1:0] pe_result [0:ROWS-1][0:COLS-1];
74 reg conv_done_r;
75
76 assign fc_mode = (mode_cfg == MODE_FC);
77
78 always @(*) begin
79     case (mode_cfg)
80         MODE_CONV1: cur_k = 9;
81         MODE_CONV2: cur_k = 36;
82         MODE_CONV3: cur_k = 72;
83         default: cur_k = 288;
84     endcase
85 end
86
87 always @(posedge clk or negedge rst_n) begin
88     if (!rst_n) begin
89         start_pipe <= {START_STAGES{1'b0}};
90     end else begin
91         start_pipe[0] <= start_pulse;
92         for (k=1; k<START_STAGES; k=k+1) begin
93             start_pipe[k] <= start_pipe[k-1];
94         end
95     end
96 end
97
98 // FC mode uses row 0 PEs; convolution mode completes on the active block rows.
99 always @* begin
100     case (valid_rows_cfg)

```

```

101     5'd0:    conv_done_r = 1'b0;
102     5'd1:    conv_done_r = finish_ij[0][COLS-1];
103     5'd2:    conv_done_r = finish_ij[1][COLS-1];
104     5'd3:    conv_done_r = finish_ij[2][COLS-1];
105     5'd4:    conv_done_r = finish_ij[3][COLS-1];
106     5'd5:    conv_done_r = finish_ij[4][COLS-1];
107     5'd6:    conv_done_r = finish_ij[5][COLS-1];
108     5'd7:    conv_done_r = finish_ij[6][COLS-1];
109     5'd8:    conv_done_r = finish_ij[7][COLS-1];
110     5'd9:    conv_done_r = finish_ij[8][COLS-1];
111     5'd10:   conv_done_r = finish_ij[9][COLS-1];
112     5'd11:   conv_done_r = finish_ij[10][COLS-1];
113     5'd12:   conv_done_r = finish_ij[11][COLS-1];
114     5'd13:   conv_done_r = finish_ij[12][COLS-1];
115     5'd14:   conv_done_r = finish_ij[13][COLS-1];
116     5'd15:   conv_done_r = finish_ij[14][COLS-1];
117     default: conv_done_r = finish_ij[15][COLS-1];
118     endcase
119 end
120 assign done = fc_mode ? finish_ij[0][2] : conv_done_r;
121
122 always @(posedge clk or negedge rst_n) begin
123     if (!rst_n) begin
124         col_stream_data_flat <= {(COLS*ACC_W){1'b0}};
125         col_stream_valid     <= {COLS{1'b0}};
126         col_stream_last      <= {COLS{1'b0}};
127         stream_active        <= {COLS{1'b0}};
128         stream_arm           <= {COLS{1'b0}};
129         for (sc=0; sc<COLS; sc=sc+1) begin
130             stream_row[sc] <= 8'd0;
131             stream_rows_limit[sc] <= 5'd0;
132             arm_rows_limit[sc] <= 5'd0;
133             col_rows_cfg[sc] <= 5'd0;
134             valid_rows_pipe[sc] <= 5'd0;
135             valid_rows_pipe_v[sc] <= 1'b0;
136         end
137     end else begin
138         col_stream_valid <= {COLS{1'b0}};
139         col_stream_last <= {COLS{1'b0}};
140         for (sc=0; sc<COLS; sc=sc+1) begin
141             data_lo = sc*ACC_W;
142             col_stream_data_flat[data_lo +: ACC_W] <= {ACC_W{1'b0}};
143         end
144
145         valid_rows_pipe[0] <= valid_rows_cfg;
146         valid_rows_pipe_v[0] <= start_pulse;
147         for (sc=1; sc<COLS; sc=sc+1) begin
148             valid_rows_pipe[sc] <= valid_rows_pipe[sc-1];
149             valid_rows_pipe_v[sc] <= valid_rows_pipe_v[sc-1];
150         end
151
152         if (start_pulse)
153             col_rows_cfg[0] <= valid_rows_cfg;
154         for (sc=1; sc<COLS; sc=sc+1) begin

```

```

155     if (valid_rows_pipe_v[sc-1])
156         col_rows_cfg[sc] <= valid_rows_pipe[sc-1];
157     end
158
159     for (sc=0; sc<COLS; sc=sc+1) begin
160         if (!stream_arm[sc] && finish_ij[0][sc]) begin
161             stream_arm[sc] <= 1'b1;
162             arm_rows_limit[sc] <= col_rows_cfg[sc];
163         end
164
165         data_lo = sc*ACC_W;
166         if (!stream_active[sc] && stream_arm[sc]) begin
167             stream_arm[sc] <= 1'b0;
168             stream_rows_limit[sc] <= arm_rows_limit[sc];
169             if (arm_rows_limit[sc] != 5'd0) begin
170                 col_stream_valid[sc] <= 1'b1;
171                 col_stream_last[sc] <= (arm_rows_limit[sc] == 5'd1);
172                 col_stream_data_flat[data_lo +: ACC_W] <= pe_result[0][sc];
173             end
174
175             if (arm_rows_limit[sc] <= 5'd1) begin
176                 stream_active[sc] <= 1'b0;
177                 stream_row[sc] <= 8'd0;
178             end else begin
179                 stream_active[sc] <= 1'b1;
180                 stream_row[sc] <= 8'd1;
181             end
182         end else if (stream_active[sc]) begin
183             col_stream_valid[sc] <= (stream_row[sc] < stream_rows_limit[sc]);
184             if (stream_row[sc] < stream_rows_limit[sc])
185                 col_stream_data_flat[data_lo +: ACC_W] <= pe_result[stream_row[sc]][sc];
186
187             if (stream_row[sc] == (stream_rows_limit[sc]-1'b1)) begin
188                 col_stream_last[sc] <= 1'b1;
189                 stream_active[sc] <= 1'b0;
190                 stream_row[sc] <= 8'd0;
191             end else begin
192                 stream_row[sc] <= stream_row[sc] + 8'd1;
193             end
194         end
195     end
196 end
197 end
198
199 genvar ar, bc;
200 generate
201     for (ar=0; ar<ROWS; ar=ar+1) begin: GEN_A_IN
202         localparam integer A_LO = ar*DATA_W;
203         localparam integer A_HI = ar*DATA_W + (DATA_W-1);
204         assign a_in[ar] = a_in_flat[A_HI:A_LO];
205     end
206     for (bc=0; bc<COLS; bc=bc+1) begin: GEN_B_IN
207         localparam integer B_LO = bc*DATA_W;
208         localparam integer B_HI = bc*DATA_W + (DATA_W-1);

```

```

209     assign b_in[bc] = b_in_flat[B_HI:B_LO];
210     end
211 endgenerate
212
213 genvar r, c;
214 generate
215     for (r=0; r<ROWS; r=r+1) begin: GEN_ROW
216         for (c=0; c<COLS; c=c+1) begin: GEN_COL
217             localparam integer O_LO = (r*COLS + c)*ACC_W;
218             localparam integer O_HI = (r*COLS + c)*ACC_W + (ACC_W-1);
219
220             if (c == 0) begin: LEFT_EDGE
221                 assign pe_left[r][c] = a_in[r];
222             end else begin: LEFT_INNER
223                 assign pe_left[r][c] = pe_right[r][c-1];
224             end
225
226             if (r == 0) begin: UP_EDGE
227                 assign pe_up[r][c] = b_in[c];
228             end else begin: UP_INNER
229                 assign pe_up[r][c] = pe_bottom[r-1][c];
230             end
231
232             assign start_ij[r][c] = fc_mode ?
233                 ((r == 0) && (c < 3)) ? ((c == 0) ? start_pulse : start_pipe[c-1]) : 1'
234 b0) :
235                 ((r == 0) && (c == 0)) ? start_pulse : start_pipe[r+c-1];
236
237         single_PE #(
238             .DATA_W(DATA_W),
239             .CNT_W (CNT_W),
240             .ACC_W (ACC_W)
241         ) PE_ij (
242             .clk (clk),
243             .rst_n (rst_n),
244             .start (start_ij[r][c]),
245             .cur_k (cur_k),
246             .finish(finish_ij[r][c]),
247             .left (pe_left[r][c]),
248             .up (pe_up[r][c]),
249             .right (pe_right[r][c]),
250             .bottom(pe_bottom[r][c]),
251             .result(pe_result[r][c])
252         );
253
254         assign c_out_flat[O_HI:O_LO] = pe_result[r][c];
255     end
256 endgenerate
257
258 endmodule
259
260
261 module single_PE #(

```

```

262 parameter integer DATA_W = 8,
263 parameter integer CNT_W = 9,
264 parameter integer ACC_W = 23
265 )(
266 input wire clk,
267 input wire rst_n,
268 input wire start,
269 input wire [CNT_W-1:0] cur_k,
270 output reg finish,
271 input wire signed [DATA_W-1:0] left,
272 input wire signed [DATA_W-1:0] up,
273
274 output reg signed [DATA_W-1:0] right,
275 output reg signed [DATA_W-1:0] bottom,
276 output reg signed [ACC_W-1:0] result
277 );
278
279 reg signed [ACC_W-1:0] mem;
280 wire signed [2*DATA_W-1:0] product;
281 reg [CNT_W-1:0] cnt;
282
283 assign product = left * up;
284
285 always @(posedge clk or negedge rst_n) begin
286     if (!rst_n) begin
287         right <= {DATA_W{1'b0}};
288         bottom <= {DATA_W{1'b0}};
289         mem <= {ACC_W{1'b0}};
290         result <= {ACC_W{1'b0}};
291         finish <= 1'b0;
292         cnt <= {CNT_W{1'b0}};
293     end else begin
294         right <= left;
295         bottom <= up;
296         finish <= 1'b0;
297
298         if (start) begin
299             if (cur_k <= 1) begin
300                 result <= {{(ACC_W-2*DATA_W){product[2*DATA_W-1]}}, product};
301                 mem <= {ACC_W{1'b0}};
302                 cnt <= {CNT_W{1'b0}};
303                 finish <= 1'b1;
304             end else begin
305                 mem <= {{(ACC_W-2*DATA_W){product[2*DATA_W-1]}}, product};
306                 cnt <= {{(CNT_W-1){1'b0}}, 1'b1};
307             end
308         end else if (cnt != {CNT_W{1'b0}}) begin
309             if (cnt == cur_k - 1'b1) begin
310                 result <= mem + {{(ACC_W-2*DATA_W){product[2*DATA_W-1]}}, product};
311                 mem <= {ACC_W{1'b0}};
312                 cnt <= {CNT_W{1'b0}};
313                 finish <= 1'b1;
314             end else begin
315                 mem <= mem + {{(ACC_W-2*DATA_W){product[2*DATA_W-1]}}, product};

```

```

316         cnt <= cnt + 1'b1;
317     end
318 end
319 end
320 end
321
322 endmodule

```

### A.8.8 Weight Buffer

*Buffers and rotates convolution weights for the systolic-array columns.*

```
code/input/RTL/conv_core/weight_buffer.v
```

```

1
2
3 module ring_shift_reg #(
4     parameter Length = 75,    // max channel depth (DOT_K_MAX + 3)
5     parameter DATA_W = 8,
6     parameter VALID_LEN = 72, // max valid window (DOT_K_MAX)
7     parameter PERIOD = 75    // max period (= Length)
8 ) (
9     input wire [1:0] layer_sel,
10    input wire clk,
11    input wire rst_n,
12    input wire wr_en,
13    input wire signed [DATA_W-1:0] in_data,
14    input wire ring,
15    input wire out_ready,
16    output wire signed [DATA_W-1:0] out_data,
17    output wire out_valid
18 );
19
20 reg signed [DATA_W-1:0] sr [Length-1:0];
21 reg signed [DATA_W-1:0] saved_sr [Length-1:0];
22
23 // Runtime decode: effective length, valid_len, period per channel
24 // These are set by the parent weight_buffer based on which channel instance this is.
25 // However, since ring_shift_reg is generic, we decode from layer_sel + compile-time
26 // Length offset. The key relationship:
27 //   eff_len = eff_dot_k + (Length - VALID_LEN) -- channel offset preserved
28 //   eff_valid = eff_dot_k
29 //   eff_period = eff_dot_k + 3 -- always equal to max channel length
30 reg [6:0] eff_dot_k;
31 reg [6:0] eff_len;
32 reg [6:0] eff_valid;
33 reg [6:0] eff_period;
34
35 localparam integer CH_OFFSET = Length - VALID_LEN; // 0, 1, 2, or 3
36 localparam [6:0] LENGTH_U7 = Length;
37

```

```

38 always @* begin
39     case (layer_sel)
40         2'b01:   eff_dot_k = 7'd36;
41         2'b10:   eff_dot_k = 7'd72;
42         default: eff_dot_k = 7'd9;
43     endcase
44     eff_valid = eff_dot_k;
45     eff_len   = eff_dot_k + CH_OFFSET[6:0];
46     eff_period = eff_dot_k + 7'd3;
47 end
48
49 reg [$clog2(PERIOD)-1:0] phase;
50 reg [$clog2(PERIOD+1)-1:0] burst_cnt;
51 reg burst_active;
52 reg ring_d;
53
54 localparam integer ALIGN_MAX = Length - VALID_LEN;
55 wire [6:0] eff_align = eff_len - eff_valid;
56 wire [6:0] feedback_idx = eff_len - 7'd1;
57 wire feedback_idx_valid = (eff_len != 7'd0) && (eff_len <= LENGTH_U7);
58 wire in_valid_window = (phase >= eff_align) && (phase < eff_align + eff_valid);
59 wire ring_rise = ring && !ring_d;
60 wire step_en = burst_active && (!in_valid_window || out_ready);
61
62 // Runtime MUX: read from the effective end of the ring, not the physical end.
63 // When DOT_K < Length (e.g. L1: eff_len=9, Length=72), the data lives in
64 // sr[0..eff_len-1] and never reaches sr[Length-1] within the burst period.
65 reg signed [DATA_W-1:0] out_tap;
66 integer oi;
67 always @* begin
68     out_tap = 0;
69     for (oi = 0; oi < Length; oi = oi + 1)
70         if (oi[6:0] == eff_len - 7'd1)
71             out_tap = sr[oi];
72 end
73 assign out_data = (burst_active && in_valid_window) ? out_tap : 0;
74 assign out_valid = burst_active && in_valid_window;
75
76 integer i;
77 always@(posedge clk or negedge rst_n)begin
78     if(!rst_n)begin
79         for(i=0;i<Length;i=i+1) begin
80             sr[i] <= 0;
81             saved_sr[i] <= 0;
82         end
83         phase <= 0;
84         burst_cnt <= 0;
85         burst_active <= 1'b0;
86         ring_d <= 1'b0;
87     end
88     else begin
89         ring_d <= ring;
90         if(ring_rise) begin
91             for (i=0; i<Length; i=i+1) saved_sr[i] <= sr[i];

```

```

92     burst_active <= 1'b1;
93     burst_cnt <= 0;
94     phase <= 0;
95     end
96     else if(burst_active && step_en) begin
97         if(burst_cnt == eff_period) begin
98             for (i=0; i<Length; i=i+1) sr[i] <= saved_sr[i];
99             phase <= 0;
100            burst_cnt <= 0;
101            burst_active <= 1'b0;
102        end else begin
103            // Runtime-select the feedback tap so smaller Length instances
104            // never statically reference out-of-range sr[] entries.
105            if (feedback_idx_valid) sr[0] <= sr[feedback_idx];
106            else sr[0] <= 0;
107            for (i=1; i<Length; i=i+1) sr[i] <= sr[i-1];
108            phase <= phase + 1'b1;
109            burst_cnt <= burst_cnt + 1'b1;
110        end
111    end
112    else if(wr_en) begin
113        sr[0] <= in_data;
114        for (i=1; i<Length; i=i+1) sr[i] <= sr[i-1];
115        phase <= 0;
116        burst_cnt <= 0;
117    end
118 end
119 end
120 endmodule
121
122
123 module weight_buffer #(
124     parameter integer DATA_W = 32,
125     parameter integer COLS = 4,
126     parameter integer L0 = 72, L1 = 73, L2 = 74, L3 = 75 // max depths
127 ) (
128     input wire [1:0] layer_sel,
129     input wire clk,
130     input wire rst_n,
131     input wire ring,
132     input wire wr_en,
133     input wire signed [DATA_W-1:0] in_data,
134     input wire out_ready,
135
136     output wire signed [DATA_W/COLS-1:0] out0, out1, out2, out3,
137     output wire [COLS-1:0] out_valid
138 );
139
140 localparam integer CH_W = DATA_W/COLS;
141 wire signed [CH_W-1:0] in0, in1, in2, in3;
142 assign in0 = in_data[0*CH_W +: CH_W];
143 assign in1 = in_data[1*CH_W +: CH_W];
144 assign in2 = in_data[2*CH_W +: CH_W];
145 assign in3 = in_data[3*CH_W +: CH_W];

```

```

146
147 ring_shift_reg #(.Length(L0), .DATA_W(CH_W), .VALID_LEN(L0), .PERIOD(L3))
148   u0_ring_shift_reg (
149     .layer_sel(layer_sel), .clk(clk), .rst_n(rst_n),
150     .wr_en(wr_en), .in_data(in0), .ring(ring), .out_ready(out_ready),
151     .out_data(out0), .out_valid(out_valid[0])
152   );
153
154 ring_shift_reg #(.Length(L1), .DATA_W(CH_W), .VALID_LEN(L0), .PERIOD(L3))
155   u1_ring_shift_reg (
156     .layer_sel(layer_sel), .clk(clk), .rst_n(rst_n),
157     .wr_en(wr_en), .in_data(in1), .ring(ring), .out_ready(out_ready),
158     .out_data(out1), .out_valid(out_valid[1])
159   );
160
161 ring_shift_reg #(.Length(L2), .DATA_W(CH_W), .VALID_LEN(L0), .PERIOD(L3))
162   u2_ring_shift_reg (
163     .layer_sel(layer_sel), .clk(clk), .rst_n(rst_n),
164     .wr_en(wr_en), .in_data(in2), .ring(ring), .out_ready(out_ready),
165     .out_data(out2), .out_valid(out_valid[2])
166   );
167
168 ring_shift_reg #(.Length(L3), .DATA_W(CH_W), .VALID_LEN(L0), .PERIOD(L3))
169   u3_ring_shift_reg (
170     .layer_sel(layer_sel), .clk(clk), .rst_n(rst_n),
171     .wr_en(wr_en), .in_data(in3), .ring(ring), .out_ready(out_ready),
172     .out_data(out3), .out_valid(out_valid[3])
173   );
174
175 endmodule

```

## A.9 Quantization, Pooling, And FC RTL

### A.9.1 Convolution-Quantization-Pooling Wrapper

*Connects convolution output streams to quantization, pooling, and SRAM writeback flow control.*

```
code/input/RTL/conv_core/conv_quant_pool.v
```

```

1 module conv_quant_pool #(
2   parameter integer W           = 64,
3   parameter integer DATA_W    = 8,
4   parameter integer ROWS       = 16,
5   parameter integer COLS       = 4,
6   parameter integer DOT_K      = 72, // max (L3: 9*8)
7   parameter integer ACC_W      = 23,
8   parameter integer IMG_H      = 64,
9   parameter integer IMG_W      = 64,
10  parameter integer C_IN       = 8, // max (L3)
11  parameter integer OUT_H      = IMG_H - 2,

```

```

12 parameter integer OUT_W      = IMG_W - 2,
13 parameter integer POOL_DEPTH = 31 // max (L1)
14 )(
15 input wire [1:0]            layer_sel,
16 input wire                  clk,
17 input wire                  rst_n,
18 input wire                  in_valid,
19 input wire [31:0]          in_data,
20 input wire [3:0]           in_byte_en,
21 input wire                  in_last,
22 output wire                 in_ready,
23 input wire                  wt_valid,
24 output wire                 wt_ready,
25 input wire                  wt_last,
26 input wire signed [COLS*DATA_W-1:0] wt_data,
27 input wire                  cfg_valid,
28 output wire                 cfg_ready,
29 input wire [31:0]          cfg_data,
30 input wire                  cfg_last,
31 input wire                  pool_ready,
32 output wire signed [31:0]   pool_data,
33 output wire                 pool_valid,
34 output wire                 pool_last,
35 // Control observation signals for FSM integration
36 output wire                 cfg_load_done,
37 output wire                 wt_load_done,
38 output wire                 pool_frame_done,
39 output wire                 conv_frame_rearm_out
40 );
41
42 wire                          sa_done;
43 wire signed [COLS*ACC_W-1:0] c_out_col_stream_flat;
44 wire [COLS-1:0]              c_out_col_valid;
45 wire [COLS-1:0]              c_out_col_last;
46 wire                          conv_frame_rearm;
47 wire                          conv_frame_done;
48
49 wire signed [ACC_W-1:0]      qp_rso0;
50 wire signed [ACC_W-1:0]      qp_rso1;
51 wire signed [ACC_W-1:0]      qp_rso2;
52 wire signed [ACC_W-1:0]      qp_rso3;
53 wire                          in_valid1;
54 wire                          in_valid2;
55 wire                          in_valid3;
56 wire                          in_valid4;
57 wire                          start1;
58 wire                          start2;
59 wire                          start3;
60 wire                          start4;
61
62 wire signed [31:0]           bias_in_w;
63 wire signed [31:0]           M_in_w;
64 wire [31:0]                  sh_in_w;
65 wire                          load_bias1;

```

```

66  wire          load_bias2;
67  wire          load_bias3;
68  wire          load_bias4;
69  wire          load_M1;
70  wire          load_M2;
71  wire          load_M3;
72  wire          load_M4;
73  wire          load_sh;
74
75  wire signed [7:0] cut1;
76  wire signed [7:0] cut2;
77  wire signed [7:0] cut3;
78  wire signed [7:0] cut4;
79  wire          cut_valid1;
80  wire          cut_valid2;
81  wire          cut_valid3;
82  wire          cut_valid4;
83
84  wire signed [7:0] cut1_out;
85  wire signed [7:0] cut2_out;
86  wire signed [7:0] cut3_out;
87  wire signed [7:0] cut4_out;
88  wire          cut_valid1_out;
89  wire          cut_valid2_out;
90  wire          cut_valid3_out;
91  wire          cut_valid4_out;
92
93  conv_top #(
94    .W(W),
95    .DATA_W(DATA_W),
96    .ROWS(ROWS),
97    .COLS(COLS),
98    .DOT_K(DOT_K),
99    .ACC_W(ACC_W),
100  .IMG_H(IMG_H),
101  .IMG_W(IMG_W),
102  .C_IN(C_IN)
103 ) u_conv1 (
104  .layer_sel(layer_sel),
105  .clk(clk),
106  .rst_n(rst_n),
107  .in_valid(in_valid),
108  .in_data(in_data),
109  .in_byte_en(in_byte_en),
110  .in_last(in_last),
111  .in_ready(in_ready),
112  .wt_valid(wt_valid),
113  .wt_ready(wt_ready),
114  .wt_last(wt_last),
115  .wt_data(wt_data),
116  .wt_load_done(wt_load_done),
117  .sa_done(sa_done),
118  .c_out_col_stream_flat(c_out_col_stream_flat),
119  .c_out_col_valid(c_out_col_valid),

```

```

120     .c_out_col_last(c_out_col_last),
121     .frame_rearm_out(conv_frame_rearm),
122     .frame_done_out(conv_frame_done)
123 );
124
125 conv_quant_adapter #(
126     .ACC_W(ACC_W),
127     .COLS(COLS),
128     .ROWS(ROWS)
129 ) u_conv_quant_adapter (
130     .clk(clk),
131     .rst_n(rst_n),
132     .c_out_col_stream_flat(c_out_col_stream_flat),
133     .c_out_col_valid(c_out_col_valid),
134     .c_out_col_last(c_out_col_last),
135     .sa_done(sa_done),
136     .qp_rso0(qp_rso0),
137     .qp_rso1(qp_rso1),
138     .qp_rso2(qp_rso2),
139     .qp_rso3(qp_rso3),
140     .in_valid1(in_valid1),
141     .in_valid2(in_valid2),
142     .in_valid3(in_valid3),
143     .in_valid4(in_valid4),
144     .start1(start1),
145     .start2(start2),
146     .start3(start3),
147     .start4(start4),
148     .en1(),
149     .en2(),
150     .en3(),
151     .en4(),
152     .bias_in(),
153     .M_in(),
154     .sh_in(),
155     .load_bias1(),
156     .load_bias2(),
157     .load_bias3(),
158     .load_bias4(),
159     .load_M1(),
160     .load_M2(),
161     .load_M3(),
162     .load_M4(),
163     .load_sh(),
164     .param_bias0(32'sd0),
165     .param_bias1(32'sd0),
166     .param_bias2(32'sd0),
167     .param_bias3(32'sd0),
168     .param_M0(32'sd0),
169     .param_M1(32'sd0),
170     .param_M2(32'sd0),
171     .param_M3(32'sd0),
172     .param_sh0(8'd0),
173     .param_sh1(8'd0),

```

```

174     .param_sh2(8'd0),
175     .param_sh3(8'd0),
176     .param_load_start(1'b0),
177     .param_load_done()
178 );
179
180 quant_param_loader u_quant_param_loader (
181     .clk(clk),
182     .rst_n(rst_n),
183     .cfg_valid(cfg_valid),
184     .cfg_ready(cfg_ready),
185     .cfg_data(cfg_data),
186     .cfg_last(cfg_last),
187     .param_load_done(cfg_load_done),
188     .bias_in(bias_in_w),
189     .M_in(M_in_w),
190     .sh_in(sh_in_w),
191     .load_bias1(load_bias1),
192     .load_bias2(load_bias2),
193     .load_bias3(load_bias3),
194     .load_bias4(load_bias4),
195     .load_M1(load_M1),
196     .load_M2(load_M2),
197     .load_M3(load_M3),
198     .load_M4(load_M4),
199     .load_sh(load_sh)
200 );
201
202 Quantization_Top u_quant (
203     .clk(clk),
204     .rst_n(rst_n),
205     .start1(start1),
206     .start2(start2),
207     .start3(start3),
208     .start4(start4),
209     .in_valid1(in_valid1),
210     .in_valid2(in_valid2),
211     .in_valid3(in_valid3),
212     .in_valid4(in_valid4),
213     .load_bias1(load_bias1),
214     .load_bias2(load_bias2),
215     .load_bias3(load_bias3),
216     .load_bias4(load_bias4),
217     .load_M1(load_M1),
218     .load_M2(load_M2),
219     .load_M3(load_M3),
220     .load_M4(load_M4),
221     .load_sh(load_sh),
222     .bias_in(bias_in_w),
223     .M_in(M_in_w),
224     .sh_in(sh_in_w),
225     .rso0(qp_rso0),
226     .rso1(qp_rso1),
227     .rso2(qp_rso2),

```

```

228     .rso3(qp_rso3),
229     .cut1(cut1),
230     .cut2(cut2),
231     .cut3(cut3),
232     .cut4(cut4),
233     .cut_valid1(cut_valid1),
234     .cut_valid2(cut_valid2),
235     .cut_valid3(cut_valid3),
236     .cut_valid4(cut_valid4)
237 );
238
239 quant_pool_adapter u_quant_pool_adapter (
240     .clk(clk),
241     .rst_n(rst_n),
242     .cut1_in(cut1),
243     .cut2_in(cut2),
244     .cut3_in(cut3),
245     .cut4_in(cut4),
246     .cut_valid1(cut_valid1),
247     .cut_valid2(cut_valid2),
248     .cut_valid3(cut_valid3),
249     .cut_valid4(cut_valid4),
250     .cut1_out(cut1_out),
251     .cut2_out(cut2_out),
252     .cut3_out(cut3_out),
253     .cut4_out(cut4_out),
254     .cut_valid1_out(cut_valid1_out),
255     .cut_valid2_out(cut_valid2_out),
256     .cut_valid3_out(cut_valid3_out),
257     .cut_valid4_out(cut_valid4_out)
258 );
259
260 pool_stream_top #(
261     .POOL_DEPTH(POOL_DEPTH)
262 ) u_pool (
263     .layer_sel(layer_sel),
264     .clk(clk),
265     .rst_n(rst_n),
266     .cut1(cut1_out),
267     .cut2(cut2_out),
268     .cut3(cut3_out),
269     .cut4(cut4_out),
270     .cut_valid1(cut_valid1_out),
271     .cut_valid2(cut_valid2_out),
272     .cut_valid3(cut_valid3_out),
273     .cut_valid4(cut_valid4_out),
274     .frame_rearm(conv_frame_rearm),
275     .pool_ready(pool_ready),
276     .pool_data(pool_data),
277     .pool_valid(pool_valid),
278     .pool_last(pool_last),
279     .pool_frame_done(pool_frame_done)
280 );
281

```

```

282   assign conv_frame_rearm_out = conv_frame_rearm;
283
284 endmodule

```

## A.9.2 Convolution-Quantization Adapter

*Adapts convolution stream timing and layout into the quantization block interface.*

code/input/RTL/quant\_pool/integration/conv\_quant\_adapter.v

```

1  module conv_quant_adapter #(
2     parameter integer ACC_W = 23,
3     parameter integer COLS = 4,
4     parameter integer ROWS = 16
5  )(
6     input wire clk,
7     input wire rst_n,
8
9     // From Conv1_top
10    input wire signed [COLS*ACC_W-1:0] c_out_col_stream_flat,
11    input wire [COLS-1:0] c_out_col_valid,
12    input wire [COLS-1:0] c_out_col_last,
13    input wire sa_done,
14
15    // To Quantization_Top
16    output wire signed [ACC_W-1:0] qp_rso0,
17    output wire signed [ACC_W-1:0] qp_rso1,
18    output wire signed [ACC_W-1:0] qp_rso2,
19    output wire signed [ACC_W-1:0] qp_rso3,
20    output wire in_valid1,
21    output wire in_valid2,
22    output wire in_valid3,
23    output wire in_valid4,
24    output reg start1,
25    output reg start2,
26    output reg start3,
27    output reg start4,
28    // Deprecated: en1..4 are not consumed by any active-path module.
29    // Retained only for legacy TB compatibility; do not rely on in new logic.
30    output wire en1,
31    output wire en2,
32    output wire en3,
33    output wire en4,
34
35    // Quant parameter loading
36    output wire signed [31:0] bias_in,
37    output wire signed [31:0] M_in,
38    output wire [31:0] sh_in,
39    output wire load_bias1,
40    output wire load_bias2,
41    output wire load_bias3,

```

```

42  output wire                load_bias4,
43  output wire                load_M1,
44  output wire                load_M2,
45  output wire                load_M3,
46  output wire                load_M4,
47  output wire                load_sh,
48  input  wire signed [31:0]  param_bias0,
49  input  wire signed [31:0]  param_bias1,
50  input  wire signed [31:0]  param_bias2,
51  input  wire signed [31:0]  param_bias3,
52  input  wire signed [31:0]  param_M0,
53  input  wire signed [31:0]  param_M1,
54  input  wire signed [31:0]  param_M2,
55  input  wire signed [31:0]  param_M3,
56  input  wire            [7:0]  param_sh0,
57  input  wire            [7:0]  param_sh1,
58  input  wire            [7:0]  param_sh2,
59  input  wire            [7:0]  param_sh3,
60  input  wire                param_load_start,
61  output wire                param_load_done
62 );
63
64 // Reserved for future boundary cleanup when c_out_col_last / sa_done become
65 // part of the explicit Conv->Quant transaction contract.
66 wire unused_last = ^c_out_col_last;
67 wire unused_done = sa_done;
68 wire _unused_ok  = unused_last ^ unused_done;
69
70 wire signed [ACC_W-1:0] col0 = c_out_col_stream_flat[0*ACC_W +: ACC_W];
71 wire signed [ACC_W-1:0] col1 = c_out_col_stream_flat[1*ACC_W +: ACC_W];
72 wire signed [ACC_W-1:0] col2 = c_out_col_stream_flat[2*ACC_W +: ACC_W];
73 wire signed [ACC_W-1:0] col3 = c_out_col_stream_flat[3*ACC_W +: ACC_W];
74
75 reg signed [ACC_W-1:0] qp_rso0_r;
76 reg signed [ACC_W-1:0] qp_rso1_r;
77 reg signed [ACC_W-1:0] qp_rso2_r;
78 reg signed [ACC_W-1:0] qp_rso3_r;
79 reg [COLS-1:0] col_valid_d1;
80 reg [COLS-1:0] col_valid_d2;
81
82 reg          cfg_compat_active;
83 reg [3:0]    cfg_compat_idx;
84 wire        cfg_valid_i;
85 wire        cfg_ready_i;
86 wire [31:0] cfg_data_i;
87 wire        cfg_last_i;
88
89 reg [3:0]    en1_sr, en2_sr, en3_sr, en4_sr;
90 reg          en1_raw, en2_raw, en3_raw, en4_raw;
91
92 assign qp_rso0 = qp_rso0_r;
93 assign qp_rso1 = qp_rso1_r;
94 assign qp_rso2 = qp_rso2_r;
95 assign qp_rso3 = qp_rso3_r;

```

```

96  assign in_valid1 = col_valid_d1[0];
97  assign in_valid2 = col_valid_d1[1];
98  assign in_valid3 = col_valid_d1[2];
99  assign in_valid4 = col_valid_d1[3];
100 assign en1      = en1_raw;
101 assign en2      = en2_raw;
102 assign en3      = en3_raw;
103 assign en4      = en4_raw;
104
105 assign cfg_valid_i = cfg_compat_active;
106 assign cfg_last_i  = (cfg_compat_idx == 4'd8);
107 assign cfg_data_i  =
108     (cfg_compat_idx == 4'd0) ? param_bias0 :
109     (cfg_compat_idx == 4'd1) ? param_bias1 :
110     (cfg_compat_idx == 4'd2) ? param_bias2 :
111     (cfg_compat_idx == 4'd3) ? param_bias3 :
112     (cfg_compat_idx == 4'd4) ? param_M0   :
113     (cfg_compat_idx == 4'd5) ? param_M1   :
114     (cfg_compat_idx == 4'd6) ? param_M2   :
115     (cfg_compat_idx == 4'd7) ? param_M3   :
116     {param_sh0, param_sh1, param_sh2, param_sh3};
117
118 always @(posedge clk or negedge rst_n) begin
119     if (!rst_n) begin
120         qp_rso0_r <= {ACC_W{1'b0}};
121         qp_rso1_r <= {ACC_W{1'b0}};
122         qp_rso2_r <= {ACC_W{1'b0}};
123         qp_rso3_r <= {ACC_W{1'b0}};
124         col_valid_d1 <= {COLS{1'b0}};
125         col_valid_d2 <= {COLS{1'b0}};
126         start1      <= 1'b0;
127         start2      <= 1'b0;
128         start3      <= 1'b0;
129         start4      <= 1'b0;
130         en1_sr      <= 4'b0;
131         en2_sr      <= 4'b0;
132         en3_sr      <= 4'b0;
133         en4_sr      <= 4'b0;
134         en1_raw     <= 1'b0;
135         en2_raw     <= 1'b0;
136         en3_raw     <= 1'b0;
137         en4_raw     <= 1'b0;
138         cfg_compat_active <= 1'b0;
139         cfg_compat_idx  <= 4'd0;
140     end else begin
141         qp_rso0_r <= col0;
142         qp_rso1_r <= col1;
143         qp_rso2_r <= col2;
144         qp_rso3_r <= col3;
145
146         col_valid_d1 <= c_out_col_valid;
147         col_valid_d2 <= col_valid_d1;
148         start1 <= col_valid_d1[0] | col_valid_d2[0];
149         start2 <= col_valid_d1[1] | col_valid_d2[1];

```

```

150     start3 <= col_valid_d1[2] | col_valid_d2[2];
151     start4 <= col_valid_d1[3] | col_valid_d2[3];
152
153     en1_sr <= {en1_sr[2:0], col_valid_d1[0]};
154     en2_sr <= {en2_sr[2:0], col_valid_d1[1]};
155     en3_sr <= {en3_sr[2:0], col_valid_d1[2]};
156     en4_sr <= {en4_sr[2:0], col_valid_d1[3]};
157     en1_raw <= en1_sr[3];
158     en2_raw <= en2_sr[3];
159     en3_raw <= en3_sr[3];
160     en4_raw <= en4_sr[3];
161
162     if (!cfg_compat_active) begin
163         if (param_load_start) begin
164             cfg_compat_active <= 1'b1;
165             cfg_compat_idx <= 4'd0;
166         end
167     end else if (cfg_valid_i && cfg_ready_i) begin
168         if (cfg_compat_idx == 4'd8) begin
169             cfg_compat_active <= 1'b0;
170             cfg_compat_idx <= 4'd0;
171         end else begin
172             cfg_compat_idx <= cfg_compat_idx + 4'd1;
173         end
174     end
175 end
176 end
177
178 quant_param_loader u_quant_param_loader (
179     .clk(clk),
180     .rst_n(rst_n),
181     .cfg_valid(cfg_valid_i),
182     .cfg_ready(cfg_ready_i),
183     .cfg_data(cfg_data_i),
184     .cfg_last(cfg_last_i),
185     .param_load_done(param_load_done),
186     .bias_in(bias_in),
187     .M_in(M_in),
188     .sh_in(sh_in),
189     .load_bias1(load_bias1),
190     .load_bias2(load_bias2),
191     .load_bias3(load_bias3),
192     .load_bias4(load_bias4),
193     .load_M1(load_M1),
194     .load_M2(load_M2),
195     .load_M3(load_M3),
196     .load_M4(load_M4),
197     .load_sh(load_sh)
198 );
199
200 endmodule

```

### A.9.3 Quantization Parameter Loader

*Loads per-layer multiplier, shift, bias, and zero-point parameters into the quantization path.*

```
code/input/RTL/quant_pool/integration/quant_param_loader.v
```

```
1 module quant_param_loader(  
2   input wire      clk,  
3   input wire      rst_n,  
4  
5   // Handshake config stream: fixed 9-word packet  
6   // 0..3: bias0..3  
7   // 4..7: M0..3  
8   // 8   : {sh0, sh1, sh2, sh3}  
9   input wire      cfg_valid,  
10  output wire     cfg_ready,  
11  input wire [31:0] cfg_data,  
12  input wire      cfg_last,  
13  
14  output wire     param_load_done,  
15  output reg signed [31:0] bias_in,  
16  output reg signed [31:0] M_in,  
17  output reg      [31:0] sh_in,  
18  output reg      load_bias1,  
19  output reg      load_bias2,  
20  output reg      load_bias3,  
21  output reg      load_bias4,  
22  output reg      load_M1,  
23  output reg      load_M2,  
24  output reg      load_M3,  
25  output reg      load_M4,  
26  output reg      load_sh  
27 );  
28  
29 localparam integer CFG_WORDS = 9;  
30  
31 localparam [3:0] PL_IDLE = 4'd0,  
32                 PL_BIAS0 = 4'd1,  
33                 PL_BIAS1 = 4'd2,  
34                 PL_BIAS2 = 4'd3,  
35                 PL_BIAS3 = 4'd4,  
36                 PL_M0    = 4'd5,  
37                 PL_M1    = 4'd6,  
38                 PL_M2    = 4'd7,  
39                 PL_M3    = 4'd8,  
40                 PL_SH    = 4'd9;  
41  
42 reg signed [31:0] bias_reg0, bias_reg1, bias_reg2, bias_reg3;  
43 reg signed [31:0] M_reg0,   M_reg1,   M_reg2,   M_reg3;  
44 reg          [31:0] sh_reg;  
45  
46 reg [3:0] rx_idx;  
47 reg [3:0] pl_state;  
48 reg      preload_active;
```

```

49  reg        load_done_r;
50
51  wire cfg_fire = cfg_valid && cfg_ready;
52
53  assign cfg_ready = !preload_active && (rx_idx < CFG_WORDS);
54  assign param_load_done = load_done_r;
55
56  always @(posedge clk or negedge rst_n) begin
57      if (!rst_n) begin
58          bias_reg0 <= 32'sd0;
59          bias_reg1 <= 32'sd0;
60          bias_reg2 <= 32'sd0;
61          bias_reg3 <= 32'sd0;
62          M_reg0    <= 32'sd0;
63          M_reg1    <= 32'sd0;
64          M_reg2    <= 32'sd0;
65          M_reg3    <= 32'sd0;
66          sh_reg    <= 32'd0;
67          bias_in   <= 32'sd0;
68          M_in      <= 32'sd0;
69          sh_in     <= 32'd0;
70          load_bias1 <= 1'b0;
71          load_bias2 <= 1'b0;
72          load_bias3 <= 1'b0;
73          load_bias4 <= 1'b0;
74          load_M1    <= 1'b0;
75          load_M2    <= 1'b0;
76          load_M3    <= 1'b0;
77          load_M4    <= 1'b0;
78          load_sh    <= 1'b0;
79          rx_idx     <= 4'd0;
80          pl_state   <= PL_IDLE;
81          preload_active <= 1'b0;
82          load_done_r <= 1'b0;
83      end else begin
84          load_bias1 <= 1'b0;
85          load_bias2 <= 1'b0;
86          load_bias3 <= 1'b0;
87          load_bias4 <= 1'b0;
88          load_M1    <= 1'b0;
89          load_M2    <= 1'b0;
90          load_M3    <= 1'b0;
91          load_M4    <= 1'b0;
92          load_sh    <= 1'b0;
93
94          if (cfg_fire) begin
95              load_done_r <= 1'b0;
96              case (rx_idx)
97                  4'd0: bias_reg0 <= $signed(cfg_data);
98                  4'd1: bias_reg1 <= $signed(cfg_data);
99                  4'd2: bias_reg2 <= $signed(cfg_data);
100                 4'd3: bias_reg3 <= $signed(cfg_data);
101                 4'd4: M_reg0    <= $signed(cfg_data);
102                 4'd5: M_reg1    <= $signed(cfg_data);

```

```

103     4'd6: M_reg2    <= $signed(cfg_data);
104     4'd7: M_reg3    <= $signed(cfg_data);
105     4'd8: sh_reg    <= cfg_data;
106     default: ;
107 endcase
108
109     if (rx_idx == CFG_WORDS - 1) begin
110         rx_idx        <= 4'd0;
111         preload_active <= 1'b1;
112         pl_state      <= PL_BIAS0;
113     end else begin
114         rx_idx <= rx_idx + 4'd1;
115     end
116 end
117
118 if (preload_active) begin
119     case (pl_state)
120     PL_BIAS0: begin
121         bias_in    <= bias_reg0;
122         load_bias1 <= 1'b1;
123         pl_state   <= PL_BIAS1;
124     end
125     PL_BIAS1: begin
126         bias_in    <= bias_reg1;
127         load_bias2 <= 1'b1;
128         pl_state   <= PL_BIAS2;
129     end
130     PL_BIAS2: begin
131         bias_in    <= bias_reg2;
132         load_bias3 <= 1'b1;
133         pl_state   <= PL_BIAS3;
134     end
135     PL_BIAS3: begin
136         bias_in    <= bias_reg3;
137         load_bias4 <= 1'b1;
138         pl_state   <= PL_M0;
139     end
140     PL_M0: begin
141         M_in       <= M_reg0;
142         load_M1    <= 1'b1;
143         pl_state   <= PL_M1;
144     end
145     PL_M1: begin
146         M_in       <= M_reg1;
147         load_M2    <= 1'b1;
148         pl_state   <= PL_M2;
149     end
150     PL_M2: begin
151         M_in       <= M_reg2;
152         load_M3    <= 1'b1;
153         pl_state   <= PL_M3;
154     end
155     PL_M3: begin
156         M_in       <= M_reg3;

```

```

157         load_M4    <= 1'b1;
158         pl_state   <= PL_SH;
159     end
160     PL_SH: begin
161         sh_in       <= sh_reg;
162         load_sh     <= 1'b1;
163         pl_state    <= PL_IDLE;
164         preload_active <= 1'b0;
165         load_done_r <= 1'b1;
166     end
167     default: begin
168         pl_state    <= PL_IDLE;
169         preload_active <= 1'b0;
170     end
171 endcase
172 end
173
174 // Keep cfg_last in the interface contract even though the current
175 // fixed 9-word packet format is ordered by index.
176 if (cfg_fire && cfg_last && rx_idx != CFG_WORDS - 1) begin
177     rx_idx        <= 4'd0;
178     load_done_r   <= 1'b0;
179 end
180 end
181 end
182
183 endmodule

```

#### A.9.4 Quantization-Pooling Adapter

*Bridges quantized outputs into the pooling stream and output writeback path.*

```
code/input/RTL/quant_pool/integration/quant_pool_adapter.v
```

```

1 module quant_pool_adapter #(
2     parameter integer OUT_W      = 62,
3     parameter integer POOL_DEPTH = 31
4 ) (
5     input wire          clk,
6     input wire          rst_n,
7     input wire signed [7:0] cut1_in,
8     input wire signed [7:0] cut2_in,
9     input wire signed [7:0] cut3_in,
10    input wire signed [7:0] cut4_in,
11    input wire          cut_valid1,
12    input wire          cut_valid2,
13    input wire          cut_valid3,
14    input wire          cut_valid4,
15    output wire signed [7:0] cut1_out,
16    output wire signed [7:0] cut2_out,
17    output wire signed [7:0] cut3_out,

```

```

18  output wire signed [7:0] cut4_out,
19  output wire          cut_valid1_out,
20  output wire          cut_valid2_out,
21  output wire          cut_valid3_out,
22  output wire          cut_valid4_out
23 );
24
25  // Reserved for future continuity smoothing or skid buffering if a later
26  // streaming pool revision adds input backpressure. The active stream-top
27  // path currently accepts continuous lane traffic directly.
28  wire unused_clk = clk;
29  wire unused_rst = rst_n;
30  wire _unused_ok = unused_clk ^ unused_rst;
31
32  assign cut1_out      = cut1_in;
33  assign cut2_out      = cut2_in;
34  assign cut3_out      = cut3_in;
35  assign cut4_out      = cut4_in;
36  assign cut_valid1_out = cut_valid1;
37  assign cut_valid2_out = cut_valid2;
38  assign cut_valid3_out = cut_valid3;
39  assign cut_valid4_out = cut_valid4;
40
41  endmodule

```

### A.9.5 Quantization Processing Element

*Implements one per-channel fixed-point requantization lane.*

code/input/RTL/quant\_pool/quant/Quantization\_PE.v

```

1  module Quantization_PE(
2      input clk, rst_n,
3      input start,
4      input in_valid,
5      input load_bias, load_M, load_sh,
6      input signed [31:0] bias_in,
7      input signed [31:0] M_in,
8      input [7:0] sh_in,
9      input signed [22:0] rso,
10     output reg signed [7:0] cut,
11     output reg cut_valid
12 );
13
14  reg load_delay; // created by load_sh by one cycle delay to activate round_term calculation
15  reg [3:0] valid_sr;
16  reg signed [31:0] bias;
17  reg signed [31:0] next_bias;
18  reg signed [31:0] M;
19  reg signed [31:0] next_M;
20  reg [7:0] sh;

```

```

21 reg [7:0] next_sh;
22 reg signed [63:0] round_term;
23 reg signed [63:0] next_round_term;
24 reg signed [31:0] acc;
25 reg signed [31:0] next_acc;
26 reg signed [63:0] mul;
27 reg signed [63:0] next_mul;
28 reg signed [63:0] correct;
29 reg signed [63:0] next_correct;
30 reg signed [63:0] shiftt;
31 reg signed [63:0] next_shiftt;
32 reg signed [63:0] y;
33 reg signed [7:0] next_cut;
34 wire signed [31:0] MM;
35
36 // create load_delay by one cycle delay from load_sh
37 always @(posedge clk or negedge rst_n) begin
38     if (!rst_n) load_delay <= 1'd0;
39     else load_delay <= load_sh;
40 end
41
42 // Output-valid tracks true input samples, not the tail-extended start signal.
43 always @(posedge clk or negedge rst_n) begin
44     if (!rst_n) begin
45         valid_sr <= 4'b0;
46         cut_valid <= 1'b0;
47     end
48     else if (load_sh) begin
49         valid_sr <= 4'b0;
50         cut_valid <= 1'b0;
51     end
52     else begin
53         valid_sr <= {valid_sr[2:0], in_valid};
54         cut_valid <= valid_sr[3];
55     end
56 end
57
58 // Bias(32-bits) loaded from outside to bias branch
59 always @* begin
60     next_bias = bias;
61     if (load_bias) begin
62         next_bias = bias_in;
63     end
64 end
65
66 always @(posedge clk or negedge rst_n) begin
67     if (!rst_n) begin
68         bias <= 32'sd0;
69     end
70     else begin
71         bias <= next_bias;
72     end
73 end
74

```

```

75 // M(32-bits) loaded from outside to M branch
76 always @* begin
77     next_M = M;
78     if (load_M) begin
79         next_M = M_in;
80     end
81 end
82
83 always @(posedge clk or negedge rst_n) begin
84     if (!rst_n) begin
85         M <= 32'sd0;
86     end
87     else begin
88         M <= next_M;
89     end
90 end
91
92 assign MM = start ? M : 32'sd0;
93
94 // sh loaded from outside to sh(shifter) branch
95 always @* begin
96     next_sh = sh;
97     if (load_sh) begin
98         next_sh = sh_in;
99     end
100 end
101
102 always @(posedge clk or negedge rst_n) begin
103     if (!rst_n) begin
104         sh <= 8'd0;
105     end
106     else begin
107         sh <= next_sh;
108     end
109 end
110
111 // one cycle delay to generate round_term
112 always @* begin
113     if (!load_delay) begin
114         next_round_term = round_term;
115     end
116     else begin
117         next_round_term = 64'sd1 << (sh - 8'd1);
118     end
119 end
120
121 always @(posedge clk or negedge rst_n) begin
122     if (!rst_n) begin
123         round_term <= 64'sd0;
124     end
125     else if (load_sh) begin
126         round_term <= 64'sd0;
127     end
128     else begin

```

```

129     round_term <= next_round_term;
130     end
131 end
132
133 // main branch
134 // step 1. adding bias
135 always @* begin
136     next_acc = rso + bias;
137 end
138
139 always @(posedge clk or negedge rst_n) begin
140     if (!rst_n) begin
141         acc <= 32'sd0;
142     end
143     else if (load_sh) begin
144         acc <= 32'sd0;
145     end
146     else begin
147         acc <= next_acc;
148     end
149 end
150
151 // step 2. acc * M
152 always @* begin
153     next_mul = acc * MM;
154 end
155
156 always @(posedge clk or negedge rst_n) begin
157     if (!rst_n) begin
158         mul <= 64'sd0;
159     end
160     else if (load_sh) begin
161         mul <= 64'sd0;
162     end
163     else begin
164         mul <= next_mul;
165     end
166 end
167
168 // step 3. Correctness: adding round_term & negative correctness
169 always @* begin
170     next_correct = mul + round_term - mul[63];
171 end
172
173 always @(posedge clk or negedge rst_n) begin
174     if (!rst_n) begin
175         correct <= 64'sd0;
176     end
177     else if (load_sh) begin
178         correct <= 64'sd0;
179     end
180     else begin
181         correct <= next_correct;
182     end

```

```

183 end
184
185 // step 4. shifting sh-bits
186 always @* begin
187     next_shiftt = correct >>> sh;
188 end
189
190 always @(posedge clk or negedge rst_n) begin
191     if (!rst_n) begin
192         shiftt <= 64'sd0;
193     end
194     else if (load_sh) begin
195         shiftt <= 64'sd0;
196     end
197     else begin
198         shiftt <= next_shiftt;
199     end
200 end
201
202 // step 5. sdding zeropoint & cutting LSB 8 bits
203 always @* begin
204     y = shiftt - 64'sd128; // 128 is the zp
205
206     // saturating clamp
207     if (y > 64'sd127) next_cut = 8'sd127;
208     else if (y < -64'sd128) next_cut = -8'sd128;
209     else next_cut = y[7:0];
210 end
211
212 always @(posedge clk or negedge rst_n) begin
213     if (!rst_n) begin
214         cut <= 8'sd0;
215     end
216     else if (load_sh) begin
217         cut <= 8'sd0;
218     end
219     else begin
220         cut <= next_cut;
221     end
222 end
223
224 endmodule

```

## A.9.6 Quantization Top

*Applies per-channel bias, multiplier, shift, zero-point, and saturation to convolution accumulators.*

```

1 module Quantization_Top(
2     input clk, rst_n,
3
4     // triggered by its own layer's finish signal by 2 cycles delay.
5     input start1, start2, start3, start4,
6     input in_valid1, in_valid2, in_valid3, in_valid4,
7
8     // bias 32 bits, only one bias_in is loaded at a cycle; same M.
9     input load_bias1, load_bias2, load_bias3, load_bias4,
10    input load_M1, load_M2, load_M3, load_M4,
11
12    // sh1, sh2, sh3, sh4 are 8 bits, can loaded all of them at a same cycle as 32 bits.
13    input load_sh,
14
15    input signed [31:0] bias_in,
16    input signed [31:0] M_in,
17
18    // sh_in 32 bits [31:24] is sh_in1, [23:16] is sh_in2, ...
19    input [31:0] sh_in,
20
21    input signed [22:0] rso0,
22    input signed [22:0] rso1,
23    input signed [22:0] rso2,
24    input signed [22:0] rso3,
25    output signed [7:0] cut1, cut2, cut3, cut4,
26    output cut_valid1, cut_valid2, cut_valid3, cut_valid4
27 );
28
29 Quantization_PE u_QPE1(
30     .clk(clk), .rst_n(rst_n), .start(start1), .in_valid(in_valid1),
31     .load_bias(load_bias1), .load_M(load_M1), .load_sh(load_sh),
32     .bias_in(bias_in), .M_in(M_in), .sh_in(sh_in[31:24]),
33     .rso(rso0), .cut(cut1), .cut_valid(cut_valid1)
34 );
35
36 Quantization_PE u_QPE2(
37     .clk(clk), .rst_n(rst_n), .start(start2), .in_valid(in_valid2),
38     .load_bias(load_bias2), .load_M(load_M2), .load_sh(load_sh),
39     .bias_in(bias_in), .M_in(M_in), .sh_in(sh_in[23:16]),
40     .rso(rso1), .cut(cut2), .cut_valid(cut_valid2)
41 );
42
43 Quantization_PE u_QPE3(
44     .clk(clk), .rst_n(rst_n), .start(start3), .in_valid(in_valid3),
45     .load_bias(load_bias3), .load_M(load_M3), .load_sh(load_sh),
46     .bias_in(bias_in), .M_in(M_in), .sh_in(sh_in[15:8]),
47     .rso(rso2), .cut(cut3), .cut_valid(cut_valid3)
48 );
49
50 Quantization_PE u_QPE4(
51     .clk(clk), .rst_n(rst_n), .start(start4), .in_valid(in_valid4),
52     .load_bias(load_bias4), .load_M(load_M4), .load_sh(load_sh),

```

```

53     .bias_in(bias_in), .M_in(M_in), .sh_in(sh_in[7:0]),
54     .rso(rso3), .cut(cut4), .cut_valid(cut_valid4)
55 );
56
57 endmodule

```

## A.9.7 Pooling Core

*Implements the streaming 2-by-2 max-pooling datapath.*

code/input/RTL/quant\_pool/pool/pool\_core.v

```

1  module pool_core_lane #(
2      parameter integer POOL_DEPTH = 31 // max across all layers
3  )(
4      input wire      [1:0] layer_sel,
5      input wire      clk,
6      input wire      rst_n,
7      input wire      lane_reset,
8      input wire signed [7:0] cut,
9      input wire      cut_valid,
10     input wire      emit_step,
11     input wire      emit_done,
12     output wire     bank_ready,
13     output wire signed [7:0] emit_head
14 );
15
16 // --- runtime decode from layer_sel ---
17 reg [4:0] eff_pool_depth;
18 reg [5:0] eff_in_w;
19 reg [5:0] eff_in_h;
20 always @* begin
21     case (layer_sel)
22         2'b01: begin eff_pool_depth = 5'd14; eff_in_w = 6'd29; eff_in_h = 6'd29; end
23         2'b10: begin eff_pool_depth = 5'd6;  eff_in_w = 6'd12; eff_in_h = 6'd12; end
24         default: begin eff_pool_depth = 5'd31; eff_in_w = 6'd62; eff_in_h = 6'd62; end
25     endcase
26 end
27
28 reg      cnt1;
29 reg [4:0] cnt2;
30 reg      cnt3;
31 reg      ready_r;
32 reg signed [7:0] accum_bank [0:POOL_DEPTH-1];
33 reg signed [7:0] emit_bank  [0:POOL_DEPTH-1];
34 reg signed [7:0] next_accum_val;
35 integer i;
36
37 // --- orphan pixel suppression counters ---
38 reg [5:0] col_in_row;
39 reg [5:0] row_cnt;

```

```

40 wire pixel_ok = (col_in_row < {1'b0, eff_pool_depth} << 1)
41     && (row_cnt < {1'b0, eff_pool_depth} << 1);
42
43 wire tile_complete = cut_valid && pixel_ok
44     && cnt1 && (cnt2 == eff_pool_depth - 5'd1) && cnt3;
45
46 wire [4:0] cnt2_next =
47     (cnt1 && (cnt2 != eff_pool_depth - 5'd1)) ? (cnt2 + 5'd1) :
48     (cnt1 && (cnt2 == eff_pool_depth - 5'd1)) ? 5'd0 :
49     cnt2;
50 wire cnt3_next = (cnt1 && (cnt2 == eff_pool_depth - 5'd1)) ? ~cnt3 : cnt3;
51
52 always @* begin
53     if (cnt3)
54         next_accum_val = (cut > accum_bank[cnt2]) ? cut : accum_bank[cnt2];
55     else if (cnt1)
56         next_accum_val = (cut > accum_bank[cnt2]) ? cut : accum_bank[cnt2];
57     else
58         next_accum_val = cut;
59 end
60
61 always @(posedge clk or negedge rst_n) begin
62     if (!rst_n) begin
63         cnt1 <= 1'b0;
64         cnt2 <= 5'd0;
65         cnt3 <= 1'b0;
66         ready_r <= 1'b0;
67         col_in_row <= 6'd0;
68         row_cnt <= 6'd0;
69         for (i = 0; i < POOL_DEPTH; i = i + 1) begin
70             accum_bank[i] <= 8'sd0;
71             emit_bank[i] <= 8'sd0;
72         end
73     end else if (lane_reset) begin
74         cnt1 <= 1'b0;
75         cnt2 <= 5'd0;
76         cnt3 <= 1'b0;
77         ready_r <= 1'b0;
78         col_in_row <= 6'd0;
79         row_cnt <= 6'd0;
80         for (i = 0; i < POOL_DEPTH; i = i + 1) begin
81             accum_bank[i] <= 8'sd0;
82             emit_bank[i] <= 8'sd0;
83         end
84     end else begin
85         if (emit_step) begin
86             for (i = 0; i < POOL_DEPTH-1; i = i + 1)
87                 emit_bank[i] <= emit_bank[i+1];
88             emit_bank[POOL_DEPTH-1] <= 8'sd0;
89         end
90
91         if (emit_done)
92             ready_r <= 1'b0;
93

```

```

94     if (cut_valid) begin
95         // advance orphan counters; auto-wrap at frame boundary
96         if (col_in_row == eff_in_w - 6'd1) begin
97             col_in_row <= 6'd0;
98             if (row_cnt == eff_in_h - 6'd1)
99                 row_cnt <= 6'd0;          // frame complete -> ready for next
100            else
101                row_cnt <= row_cnt + 6'd1;
102        end else begin
103            col_in_row <= col_in_row + 6'd1;
104        end
105
106        if (pixel_ok) begin
107            if (tile_complete) begin
108                for (i = 0; i < POOL_DEPTH; i = i + 1) begin
109                    if (i == cnt2)
110                        emit_bank[i] <= next_accum_val;
111                    else
112                        emit_bank[i] <= accum_bank[i];
113                        accum_bank[i] <= 8'sd0;
114                    end
115                    ready_r <= 1'b1;
116                    cnt1    <= 1'b0;
117                    cnt2    <= 5'd0;
118                    cnt3    <= 1'b0;
119                end else begin
120                    accum_bank[cnt2] <= next_accum_val;
121                    cnt1             <= ~cnt1;
122                    cnt2             <= cnt2_next;
123                    cnt3             <= cnt3_next;
124                end
125            end
126            // !pixel_ok -> orphan pixel, only counters advanced, accum/cnt untouched
127        end
128    end
129 end
130
131 assign bank_ready = ready_r;
132 assign emit_head  = emit_bank[0];
133
134 endmodule
135
136 module pool_core #(
137     parameter integer POOL_DEPTH = 31
138 ) (
139     input wire      [1:0] layer_sel,
140     input wire      clk,
141     input wire      rst_n,
142     input wire signed [7:0] cut1,
143     input wire signed [7:0] cut2,
144     input wire signed [7:0] cut3,
145     input wire signed [7:0] cut4,
146     input wire      en1,
147     input wire      en2,

```

```

148 input wire          en3,
149 input wire          en4,
150 input wire          frame_rearm,
151 input wire          pool_ready,
152 output wire signed [31:0] pool_data,
153 output wire          pool_valid,
154 output wire          pool_last,
155 output wire          pool_frame_done
156 );
157
158 // runtime decode for emit_last
159 reg [4:0] eff_pool_depth;
160 always @* begin
161     case (layer_sel)
162         2'b01: eff_pool_depth = 5'd14;
163         2'b10: eff_pool_depth = 5'd6;
164         default: eff_pool_depth = 5'd31;
165     endcase
166 end
167
168 // Delay frame_rearm by 8 cycles so quant pipeline residual is flushed.
169 // Reset on layer_changed so stale frame_rearm from previous layer doesn't
170 // prematurely clear drop_until_idle in the new layer context.
171 reg emit_active;
172 reg [5:0] emit_count;
173 reg lane_reset_r;
174 reg drop_until_idle;
175 reg [1:0] layer_sel_d;
176 wire layer_changed;
177 wire l2_to_l3_change = (layer_sel_d == 2'b01) && (layer_sel == 2'b10);
178 reg [7:0] frame_rearm_sr;
179 reg frame_rearm_hold_active;
180 reg frame_rearm_pending;
181 wire frame_rearm_delayed = frame_rearm_sr[7];
182 always @(posedge clk or negedge rst_n) begin
183     if (!rst_n) frame_rearm_sr <= 8'd0;
184     else frame_rearm_sr <= {frame_rearm_sr[6:0], frame_rearm};
185 end
186
187 assign layer_changed = (layer_sel != layer_sel_d);
188 wire lane1_ready, lane2_ready, lane3_ready, lane4_ready;
189 wire signed [7:0] lane1_emit, lane2_emit, lane3_emit, lane4_emit;
190 wire all_ready = lane1_ready & lane2_ready & lane3_ready & lane4_ready;
191 wire emit_last = (emit_count == {1'b0, eff_pool_depth} - 6'd1);
192 wire emit_fire = emit_active && pool_ready;
193 wire emit_step = emit_fire;
194 wire emit_done = emit_fire && emit_last;
195 wire any_cut_valid = en1 | en2 | en3 | en4;
196 wire lane_cut_valid1 = drop_until_idle ? 1'b0 : en1;
197 wire lane_cut_valid2 = drop_until_idle ? 1'b0 : en2;
198 wire lane_cut_valid3 = drop_until_idle ? 1'b0 : en3;
199 wire lane_cut_valid4 = drop_until_idle ? 1'b0 : en4;
200
201 pool_core_lane #(.POOL_DEPTH(POOL_DEPTH)) u_lane1 (

```

```

202     .layer_sel(layer_sel),
203     .clk(clk), .rst_n(rst_n),
204     .lane_reset(lane_reset_r),
205     .cut(cut1), .cut_valid(lane_cut_valid1),
206     .emit_step(emit_step), .emit_done(emit_done),
207     .bank_ready(lane1_ready), .emit_head(lane1_emit)
208 );
209
210 pool_core_lane #(.POOL_DEPTH(POOL_DEPTH)) u_lane2 (
211     .layer_sel(layer_sel),
212     .clk(clk), .rst_n(rst_n),
213     .lane_reset(lane_reset_r),
214     .cut(cut2), .cut_valid(lane_cut_valid2),
215     .emit_step(emit_step), .emit_done(emit_done),
216     .bank_ready(lane2_ready), .emit_head(lane2_emit)
217 );
218
219 pool_core_lane #(.POOL_DEPTH(POOL_DEPTH)) u_lane3 (
220     .layer_sel(layer_sel),
221     .clk(clk), .rst_n(rst_n),
222     .lane_reset(lane_reset_r),
223     .cut(cut3), .cut_valid(lane_cut_valid3),
224     .emit_step(emit_step), .emit_done(emit_done),
225     .bank_ready(lane3_ready), .emit_head(lane3_emit)
226 );
227
228 pool_core_lane #(.POOL_DEPTH(POOL_DEPTH)) u_lane4 (
229     .layer_sel(layer_sel),
230     .clk(clk), .rst_n(rst_n),
231     .lane_reset(lane_reset_r),
232     .cut(cut4), .cut_valid(lane_cut_valid4),
233     .emit_step(emit_step), .emit_done(emit_done),
234     .bank_ready(lane4_ready), .emit_head(lane4_emit)
235 );
236
237 assign pool_data = {lane1_emit, lane2_emit, lane3_emit, lane4_emit};
238 assign pool_valid = emit_active;
239 assign pool_last = emit_active && emit_last;
240
241 // --- frame-level burst counter for pool_frame_done ---
242 reg [4:0] eff_burst_total;
243 always @* begin
244     case (layer_sel)
245         2'b01:    eff_burst_total = 5'd14; // L2
246         2'b10:    eff_burst_total = 5'd6;  // L3
247         default: eff_burst_total = 5'd31; // L1
248     endcase
249 end
250
251 reg [4:0] burst_cnt;
252 reg      frame_done_r;
253
254 assign pool_frame_done = frame_done_r;
255

```

```

256 always @(posedge clk or negedge rst_n) begin
257     if (!rst_n) begin
258         emit_active  <= 1'b0;
259         emit_count   <= 6'd0;
260         burst_cnt    <= 5'd0;
261         frame_done_r <= 1'b0;
262         lane_reset_r <= 1'b0;
263         drop_until_idle <= 1'b0;
264         layer_sel_d   <= 2'b00;
265         frame_rearm_hold_active <= 1'b0;
266         frame_rearm_pending <= 1'b0;
267     end else begin
268         frame_done_r <= 1'b0;
269         lane_reset_r <= 1'b0;
270         layer_sel_d  <= layer_sel;
271
272         if (frame_rearm)
273             frame_rearm_pending <= 1'b1;
274
275         if (layer_changed) begin
276             emit_active  <= 1'b0;
277             emit_count   <= 6'd0;
278             burst_cnt    <= 5'd0;
279             lane_reset_r <= 1'b1;
280             drop_until_idle <= l2_to_l3_change;
281             frame_rearm_hold_active <= l2_to_l3_change;
282             frame_rearm_pending <= 1'b0;
283         end else begin
284             if (frame_rearm_hold_active && frame_rearm_delayed) begin
285                 drop_until_idle <= 1'b0;
286                 lane_reset_r    <= 1'b1;
287                 frame_rearm_hold_active <= 1'b0;
288                 frame_rearm_pending <= 1'b0;
289             end
290
291             if (!emit_active) begin
292                 emit_count <= 6'd0;
293                 if (all_ready)
294                     emit_active <= 1'b1;
295             end else if (emit_fire) begin
296                 if (emit_last) begin
297                     emit_active <= 1'b0;
298                     emit_count  <= 6'd0;
299                     if (burst_cnt == eff_burst_total - 5'd1) begin
300                         burst_cnt    <= 5'd0;
301                         frame_done_r  <= 1'b1;
302                         lane_reset_r  <= 1'b1;
303                         if (frame_rearm_pending) begin
304                             frame_rearm_pending <= 1'b0;
305                         end else begin
306                             drop_until_idle      <= 1'b1;
307                             frame_rearm_hold_active <= 1'b1;
308                         end
309                     end else begin

```

```

310         burst_cnt <= burst_cnt + 5'd1;
311     end
312 end else begin
313     emit_count <= emit_count + 6'd1;
314 end
315 end
316 end
317 end
318 end
319
320 endmodule

```

### A.9.8 Pooling Stream Top

*Wraps the pooling core with stream-level valid, ready, and last handling.*

```
code/input/RTL/quant_pool/pool/pool_stream_top.v
```

```

1 module pool_stream_top #(
2     parameter integer POOL_DEPTH = 31
3 ) (
4     input wire        [1:0] layer_sel,
5     input wire        clk,
6     input wire        rst_n,
7     input wire signed [7:0] cut1,
8     input wire signed [7:0] cut2,
9     input wire signed [7:0] cut3,
10    input wire signed [7:0] cut4,
11    input wire        cut_valid1,
12    input wire        cut_valid2,
13    input wire        cut_valid3,
14    input wire        cut_valid4,
15    input wire        frame_rearm,
16    input wire        pool_ready,
17    output wire signed [31:0] pool_data,
18    output wire        pool_valid,
19    output wire        pool_last,
20    output wire        pool_frame_done
21 );
22
23 pool_core #(
24     .POOL_DEPTH(POOL_DEPTH)
25 ) u_pool_core (
26     .layer_sel(layer_sel),
27     .clk(clk),
28     .rst_n(rst_n),
29     .cut1(cut1),
30     .cut2(cut2),
31     .cut3(cut3),
32     .cut4(cut4),
33     .en1(cut_valid1),

```

```

34     .en2(cut_valid2),
35     .en3(cut_valid3),
36     .en4(cut_valid4),
37     .frame_rearm(frame_rearm),
38     .pool_ready(pool_ready),
39     .pool_data(pool_data),
40     .pool_valid(pool_valid),
41     .pool_last(pool_last),
42     .pool_frame_done(pool_frame_done)
43 );
44
45 endmodule

```

### A.9.9 Fully-Connected Top

*Implements the final classifier datapath and output logits.*

code/input/RTL/fc/FC.v

```

1  module FC #(
2      parameter PIX_W      = 8,
3      parameter KER_W      = 8,
4      parameter K          = 288,
5      parameter ACC_W      = 32,
6      parameter OUT_CHANNELS = 10
7  )(
8      input  wire  clk,
9      input  wire  rst_n,
10
11     // cfg handshake for eff_bias loading (OUT_CHANNELS-word packet)
12     input  wire      cfg_valid,
13     output wire      cfg_ready,
14     input  wire [31:0] cfg_data,
15     input  wire      cfg_last,
16     output wire      param_load_done,
17
18     // MAC data path (packed vectors, LSB = channel 0)
19     input  wire      mul_en,
20     input  wire [OUT_CHANNELS*PIX_W-1:0] pixel_vec,
21     input  wire [OUT_CHANNELS*KER_W-1:0] kernel_vec,
22     output wire [OUT_CHANNELS*ACC_W-1:0] acc_vec,
23     output wire [OUT_CHANNELS-1:0] mul_done_vec,
24     output wire      all_done
25 );
26
27 // Bias loader -> MAC wires
28 wire signed [31:0] bias_in_w;
29 wire [OUT_CHANNELS-1:0] load_bias_vec;
30
31 fc_bias_loader #(
32     .OUT_CHANNELS(OUT_CHANNELS)

```

```

33 ) u_bias_loader (
34     .clk(clk),
35     .rst_n(rst_n),
36     .cfg_valid(cfg_valid),
37     .cfg_ready(cfg_ready),
38     .cfg_data(cfg_data),
39     .cfg_last(cfg_last),
40     .param_load_done(param_load_done),
41     .bias_in(bias_in_w),
42     .load_bias_vec(load_bias_vec)
43 );
44
45 // Per-channel MAC array
46 genvar gi;
47 generate
48     for (gi = 0; gi < OUT_CHANNELS; gi = gi + 1) begin : g_mac
49         wire signed [PIX_W-1:0] pix_i = pixel_vec [gi*PIX_W +: PIX_W];
50         wire signed [KER_W-1:0] ker_i = kernel_vec[gi*KER_W +: KER_W];
51         wire signed [ACC_W-1:0] acc_i;
52         wire
53             done_i;
54
55         mac #(
56             .PIX_W(PIX_W),
57             .KER_W(KER_W),
58             .K(K),
59             .ACC_W(ACC_W)
60         ) u_mac (
61             .clk(clk),
62             .rst_n(rst_n),
63             .mul_en(mul_en),
64             .pixel(pix_i),
65             .kernel(ker_i),
66             .load_bias(load_bias_vec[gi]),
67             .bias_in(bias_in_w[ACC_W-1:0]),
68             .acc(acc_i),
69             .mul_done(done_i)
70         );
71
72         assign acc_vec[gi*ACC_W +: ACC_W] = acc_i;
73         assign mul_done_vec[gi] = done_i;
74     end
75 endgenerate
76
77 assign all_done = &mul_done_vec;
78 endmodule

```

## A.9.10 FC Bias Loader

*Loads and presents FC bias terms for classifier accumulation.*

```

1 module fc_bias_loader #(
2     parameter OUT_CHANNELS = 10
3 )
4     input wire      clk,
5     input wire      rst_n,
6
7     // 4-wire cfg handshake: OUT_CHANNELS-word packet, one bias per word.
8     input wire      cfg_valid,
9     output wire     cfg_ready,
10    input wire [31:0] cfg_data,
11    input wire      cfg_last,
12
13    output wire      param_load_done,
14    output reg signed [31:0] bias_in,
15    output reg [OUT_CHANNELS-1:0] load_bias_vec
16 );
17
18 localparam integer CFG_WORDS = OUT_CHANNELS;
19 localparam integer IDX_W      = (OUT_CHANNELS <= 1) ? 1 : $clog2(OUT_CHANNELS);
20
21 localparam [1:0] PL_IDLE      = 2'd0,
22                 PL_DISTRIBUTE = 2'd1;
23
24 reg signed [31:0] bias_reg [0:OUT_CHANNELS-1];
25 reg [IDX_W:0] rx_idx; // one extra bit to compare against CFG_WORDS
26 reg [IDX_W:0] pl_idx;
27 reg [1:0] pl_state;
28 reg preload_active;
29 reg load_done_r;
30
31 wire cfg_fire = cfg_valid && cfg_ready;
32
33 assign cfg_ready      = !preload_active && (rx_idx < CFG_WORDS);
34 assign param_load_done = load_done_r;
35
36 always @(posedge clk or negedge rst_n) begin : bias_fsm
37     integer i;
38     if (!rst_n) begin
39         for (i = 0; i < OUT_CHANNELS; i = i + 1)
40             bias_reg[i] <= 32'sd0;
41         bias_in      <= 32'sd0;
42         load_bias_vec <= {OUT_CHANNELS{1'b0}};
43         rx_idx       <= {(IDX_W+1){1'b0}};
44         pl_idx       <= {(IDX_W+1){1'b0}};
45         pl_state     <= PL_IDLE;
46         preload_active <= 1'b0;
47         load_done_r  <= 1'b0;
48     end else begin
49         // Default: clear one-hot pulse
50         load_bias_vec <= {OUT_CHANNELS{1'b0}};
51
52         // ---- RX phase: accept OUT_CHANNELS cfg words ----

```

```

53   if (cfg_fire) begin
54       load_done_r <= 1'b0;
55       bias_reg[rx_idx[IDX_W-1:0]] <= $signed(cfg_data);
56
57       if (rx_idx == CFG_WORDS - 1) begin
58           rx_idx          <= {(IDX_W+1){1'b0}};
59           preload_active <= 1'b1;
60           pl_state       <= PL_DISTRIBUTE;
61           pl_idx         <= {(IDX_W+1){1'b0}};
62       end else begin
63           rx_idx <= rx_idx + 1'b1;
64       end
65   end
66
67   // ---- Preload phase: one channel per cycle ----
68   if (preload_active && pl_state == PL_DISTRIBUTE) begin
69       bias_in          <= bias_reg[pl_idx[IDX_W-1:0]];
70       load_bias_vec[pl_idx[IDX_W-1:0]] <= 1'b1;
71
72       if (pl_idx == OUT_CHANNELS - 1) begin
73           pl_state      <= PL_IDLE;
74           preload_active <= 1'b0;
75           load_done_r   <= 1'b1;
76           pl_idx        <= {(IDX_W+1){1'b0}};
77       end else begin
78           pl_idx <= pl_idx + 1'b1;
79       end
80   end
81
82   // Early cfg_last resets RX if packet is short
83   if (cfg_fire && cfg_last && rx_idx != CFG_WORDS - 1) begin
84       rx_idx          <= {(IDX_W+1){1'b0}};
85       load_done_r <= 1'b0;
86   end
87 end
88 end
89
90 endmodule

```

### A.9.11 FC Data Adapter

*Adapts SRAM\_B feature data and FC weight data into the classifier MAC schedule.*

```
code/input/RTL/fc/fc_data_adapter.v
```

```

1 // FC Data Adapter: synchronize SRAM_FCW weight stream and SRAM_B pixel stream
2 // into FC pixel_vec/kernel_vec + mul_en.
3 //
4 // SRAM_FCW (weight): 288 words x 80b, one word per MAC cycle.
5 // Packing (LSB = channel 0): {k9, k8, k7, k6, k5, k4, k3, k2, k1, k0}
6 // SRAM_B (data):      72 words x 32b, one word per 4 MAC cycles (byte-serial).

```

```

7 //
8 // The same pixel byte is broadcast to all OUT_CHANNELS MACs in the same cycle;
9 // per-channel differentiation lives entirely in the 80-bit weight word.
10
11 module fc_data_adapter #(
12     parameter OUT_CHANNELS = 10,
13     parameter PIX_W       = 8,
14     parameter KER_W       = 8
15 )
16     input wire      clk,
17     input wire      rst_n,
18
19     // FC weight stream (from SRAM_FCW, 288 words x 80b)
20     input wire      wt_valid,
21     input wire [OUT_CHANNELS*KER_W-1:0] wt_data,
22     input wire      wt_last,
23     output wire     wt_ready,
24
25     // FC data stream (from SRAM_B, 72 packed words)
26     input wire      data_valid,
27     input wire [31:0] data_data,
28     input wire      data_last,
29     output wire     data_ready,
30
31     // FC data path (packed vectors, LSB = channel 0)
32     output wire     mul_en,
33     output wire [OUT_CHANNELS*PIX_W-1:0] pixel_vec,
34     output wire [OUT_CHANNELS*KER_W-1:0] kernel_vec,
35
36     // FC completion
37     input wire      all_done,
38     output wire     fc_done
39 );
40     localparam integer WT_W = OUT_CHANNELS*KER_W;
41
42     // --- Holding registers ---
43     reg [WT_W-1:0] wt_reg;
44     reg           wt_held;
45     reg [31:0]    data_reg;
46     reg           data_held;
47     reg [1:0]    byte_sel;
48
49     wire wt_avail    = wt_held;
50     wire pixel_avail = data_held;
51
52     assign mul_en    = wt_avail && pixel_avail && !all_done;
53     assign wt_ready  = (wt_held && mul_en) || (!wt_held && !wt_valid);
54     assign data_ready = (data_held && mul_en && (byte_sel == 2'd3))
55                         || (!data_held && !data_valid);
56
57     // Kernel vector: directly drive from held 80-bit weight word
58     assign kernel_vec = wt_reg;
59
60     // Pixel byte selection (SRAM_B packs high byte first, matching current behavior)

```

```

61 wire [7:0] pixel_byte = (byte_sel == 2'd0) ? data_reg[31:24] :
62                       (byte_sel == 2'd1) ? data_reg[23:16] :
63                       (byte_sel == 2'd2) ? data_reg[15: 8] :
64                       data_reg[ 7: 0];
65
66 // Broadcast one pixel byte to all OUT_CHANNELS slots
67 genvar gi;
68 generate
69     for (gi = 0; gi < OUT_CHANNELS; gi = gi + 1) begin : g_pix
70         assign pixel_vec[gi*PIX_W +: PIX_W] = $signed(pixel_byte);
71     end
72 endgenerate
73
74 assign fc_done = all_done;
75
76 always @(posedge clk or negedge rst_n) begin
77     if (!rst_n) begin
78         wt_reg    <= {WT_W{1'b0}};
79         wt_held   <= 1'b0;
80         data_reg  <= 32'd0;
81         data_held <= 1'b0;
82         byte_sel  <= 2'd0;
83     end else begin
84         if (!wt_held) begin
85             if (wt_valid) begin
86                 wt_reg <= wt_data;
87                 wt_held <= 1'b1;
88             end
89             end else if (mul_en) begin
90                 wt_held <= 1'b0;
91             end
92
93             if (!data_held) begin
94                 if (data_valid) begin
95                     data_reg <= data_data;
96                     data_held <= 1'b1;
97                     byte_sel <= 2'd0;
98                 end
99                 end else if (mul_en) begin
100                     if (byte_sel == 2'd3) begin
101                         byte_sel <= 2'd0;
102                         data_held <= 1'b0;
103                     end else begin
104                         byte_sel <= byte_sel + 2'd1;
105                     end
106                 end
107             end
108         end
109     end
110 endmodule

```

## A.9.12 FC MAC

*Implements the multiply-accumulate primitive used by the fully-connected classifier.*

code/input/RTL/fc/mac.v

```
1 module mac #(
2   parameter PIX_W = 8,
3   parameter KER_W = 8,
4   parameter K = 288,
5   parameter ACC_W = 32
6 ) (
7   input wire clk,
8   input wire rst_n,
9   input wire mul_en,
10  input wire signed [PIX_W-1:0] pixel,
11  input wire signed [KER_W-1:0] kernel,
12  // Bias preload
13  input wire          load_bias,
14  input wire signed [ACC_W-1:0] bias_in,
15  output reg signed [ACC_W-1:0] acc,
16  output reg mul_done
17 );
18
19 wire signed [PIX_W+KER_W-1:0] product;
20 wire signed [ACC_W-1:0] product_ext;
21
22 reg [$clog2(K)-1:0] cnt;
23 reg signed [ACC_W-1:0] bias_reg;
24
25 assign product = $signed(pixel) * $signed(kernel);
26 assign product_ext = {{(ACC_W-(PIX_W+KER_W)){product[PIX_W+KER_W-1]}}, product};
27
28 // Bias register: latch on load_bias pulse
29 always @(posedge clk or negedge rst_n) begin
30   if (!rst_n)
31     bias_reg <= {ACC_W{1'b0}};
32   else if (load_bias)
33     bias_reg <= bias_in;
34 end
35
36 // MAC accumulator
37 always @(posedge clk or negedge rst_n) begin
38   if (!rst_n) begin
39     acc <= {ACC_W{1'b0}};
40     mul_done <= 1'b0;
41     cnt <= {$clog2(K){1'b0}};
42   end else if (load_bias) begin
43     // Reset for new inference: clear accumulator state so mul_en can fire
44     acc <= {ACC_W{1'b0}};
45     mul_done <= 1'b0;
46     cnt <= {$clog2(K){1'b0}};
47   end else if (mul_en) begin
48     if (cnt == 0) begin
```

```

49     acc    <= bias_reg + product_ext;
50     mul_done <= 1'b0;
51     cnt    <= cnt + 1'b1;
52     end else if (cnt == K - 1) begin
53         acc    <= acc + product_ext;
54         mul_done <= 1'b1;
55         cnt    <= {$clog2(K){1'b0}};
56     end else begin
57         acc    <= acc + product_ext;
58         mul_done <= 1'b0;
59         cnt    <= cnt + 1'b1;
60     end
61 end
62 end
63
64 endmodule

```

## A.10 SRAM Controller RTL

### A.10.1 SRAM Address Generator

*Generates layer-dependent SRAM addresses for configuration, weights, activations, and writeback.*

code/input/RTL/SRAM/Addr\_Gen.v

```

1  module Addr_Gen #(
2      parameter ADDR_WIDTH = 11
3  )(
4      input  clk,
5      input  rst_n,
6
7      input  start,
8      input  enable,
9      input  [ADDR_WIDTH-1:0] base_addr,
10     input  [ADDR_WIDTH-1:0] length,
11
12     output [ADDR_WIDTH-1:0] addr,
13     output done
14 );
15
16     // -----
17     // Internal Registers
18     // -----
19     reg [ADDR_WIDTH-1:0] counter;
20     reg                 running;
21
22     // -----
23     // Running Control
24     // -----
25     always @(posedge clk or negedge rst_n) begin

```

```

26     if (!rst_n)
27         running <= 1'b0;
28     else if (start)
29         running <= 1'b1;
30     else if (done)
31         running <= 1'b0;
32 end
33
34 // -----
35 // Counter
36 // -----
37 always @(posedge clk or negedge rst_n) begin
38     if (!rst_n)
39         counter <= {ADDR_WIDTH{1'b0}};
40     else if (start)
41         counter <= {ADDR_WIDTH{1'b0}};
42     else if (running && enable && !done)
43         counter <= counter + 1'b1;
44 end
45
46 // -----
47 // Combinational Outputs
48 // -----
49 assign addr = base_addr + counter;
50
51 assign done = running && enable &&
52             (counter == (length - 1'b1));
53
54 endmodule

```

## A.10.2 FC Weight Preload Packer

*Packs incoming FC weight preload words into the SRAM\_FCW layout.*

```
code/input/RTL/SRAM/fcw_preload_packer.v
```

```

1  `timescale 1ns/1ps
2
3  // fcw_preload_packer: pack three 32-bit host words into one 80-bit FCW write.
4  //
5  // Host preload stream carries FC weights at 32b/word. Each FCW slot holds
6  // 10 kernels x 8b = 80b. The packer consumes 3 host words per FCW slot:
7  //   word0 bits[31:0]  -> fcw_word[31:0]   = {k3, k2, k1, k0}
8  //   word1 bits[31:0]  -> fcw_word[63:32]  = {k7, k6, k5, k4}
9  //   word2 bits[15:0]  -> fcw_word[79:64]  = {k9, k8}   (upper 16b ignored)
10 //
11 // On the third beat the packer emits a single wea pulse at the current
12 // FCW address, then advances the address. 'start' resets the address
13 // counter and phase counter.
14
15 module fcw_preload_packer # (

```

```

16     parameter AW = 9,
17     parameter DW = 80
18 ) (
19     input          clk,
20     input          rst_n,
21
22     input          start,          // resets address/phase counters
23
24     // 32b host preload stream
25     input          up_valid,
26     input    [31:0] up_data,
27     output         up_ready,
28
29     // 80b write port to sram_FCW_wrapper
30     output         wea,
31     output    [AW-1:0]  addra,
32     output    [DW-1:0]  dina
33 );
34 reg [1:0]    phase;
35 reg [AW-1:0] wr_cnt;
36 reg [31:0]  word0_r;
37 reg [31:0]  word1_r;
38
39 wire fire = up_valid && up_ready;
40
41 assign up_ready = 1'b1;
42
43 // Combinationally emit wea on the cycle the third host word arrives.
44 assign wea    = fire && (phase == 2'd2);
45 assign addra = wr_cnt;
46 assign dina  = {up_data[15:0], word1_r, word0_r};
47
48 always @(posedge clk or negedge rst_n) begin
49     if (!rst_n) begin
50         phase <= 2'd0;
51         wr_cnt <= {AW{1'b0}};
52         word0_r <= 32'd0;
53         word1_r <= 32'd0;
54     end else if (start) begin
55         phase <= 2'd0;
56         wr_cnt <= {AW{1'b0}};
57     end else if (fire) begin
58         case (phase)
59             2'd0: begin
60                 word0_r <= up_data;
61                 phase <= 2'd1;
62             end
63             2'd1: begin
64                 word1_r <= up_data;
65                 phase <= 2'd2;
66             end
67             2'd2: begin
68                 phase <= 2'd0;
69                 wr_cnt <= wr_cnt + {{(AW-1){1'b0}}, 1'b1};

```

```

70         end
71         default: phase <= 2'd0;
72     endcase
73     end
74 end
75
76 endmodule

```

### A.10.3 SRAM A Controller

*Controls SRAM\_A reads and writes for preload data, convolution input, and intermediate feature maps.*

code/input/RTL/SRAM/sram\_A\_controller.v

```

1  module sram_A_controller #(
2      parameter ADDR_WIDTH = 10
3  )(
4      input [2:0] layer_sel,
5      input [1:0] data_sel,
6      input      pass_id,
7
8      output reg [ADDR_WIDTH-1:0] base_addr,
9      output reg [ADDR_WIDTH-1:0] length,
10     output reg      read_mode,
11     output reg      write_mode
12 );
13
14 localparam [2:0] LAYER_PRELOAD = 3'd0;
15 localparam [2:0] LAYER_L1      = 3'd1;
16 localparam [2:0] LAYER_L2      = 3'd2;
17 localparam [2:0] LAYER_L3      = 3'd3;
18 localparam [2:0] LAYER_FC      = 3'd4;
19
20 localparam [1:0] SEL_CFG        = 2'd0;
21 localparam [1:0] SEL_WT        = 2'd1;
22 localparam [1:0] SEL_DATA      = 2'd2;
23 localparam [1:0] SEL_FCW       = 2'd3; // FC weights in dedicated 80b SRAM_FCW
24
25 localparam [ADDR_WIDTH-1:0] PRELOAD_CFG_BASE = 10'h000;
26 localparam [ADDR_WIDTH-1:0] PRELOAD_CFG_LEN  = 10'd45; // conv L1/L2/L3 cfg only; FC bias moved to
    dedicated slot
27 localparam [ADDR_WIDTH-1:0] PRELOAD_WT_BASE  = 10'h030;
28 localparam [ADDR_WIDTH-1:0] PRELOAD_WT_LEN   = 10'd225; // conv L1/L2/L3 weights only; FC weights moved
    to SRAM_FCW
29
30 // FC bias preload reuses SRAM_A (32b words) but at its own base outside
31 // the conv region. Preload writes 10 bias words to the free zone that used
32 // to hold the 3-class FC weight table.
33 localparam [ADDR_WIDTH-1:0] PRELOAD_FC_CFG_BASE = 10'h111;

```

```

34 localparam [ADDR_WIDTH-1:0] PRELOAD_FC_CFG_LEN = 10'd10;
35
36 // FC weight preload streams through SRAM_A write port (host 32b) into the
37 // SRAM_FCW packer: 288 FCW slots x 3 host words = 864 beats.
38 localparam [ADDR_WIDTH-1:0] PRELOAD_FCW_LEN = 10'd864;
39 // DATA preload is no longer stored in SRAM_A.
40 localparam [ADDR_WIDTH-1:0] PRELOAD_DATA_BASE = 10'h231;
41 localparam [ADDR_WIDTH-1:0] PRELOAD_DATA_LEN = 10'd0;
42
43 localparam [ADDR_WIDTH-1:0] L1_CFG_BASE = 10'h000;
44 localparam [ADDR_WIDTH-1:0] L1_CFG_LEN = 10'd9;
45 localparam [ADDR_WIDTH-1:0] L1_WT_BASE = 10'h030;
46 localparam [ADDR_WIDTH-1:0] L1_WT_LEN = 10'd9;
47 // L1 input data is no longer read from SRAM_A.
48 localparam [ADDR_WIDTH-1:0] L1_DATA_BASE = 10'h231;
49 localparam [ADDR_WIDTH-1:0] L1_DATA_LEN = 10'd0;
50
51 localparam [ADDR_WIDTH-1:0] L2_PASS0_CFG_BASE = 10'h009;
52 localparam [ADDR_WIDTH-1:0] L2_PASS0_CFG_LEN = 10'd9;
53 localparam [ADDR_WIDTH-1:0] L2_PASS1_CFG_BASE = 10'h012;
54 localparam [ADDR_WIDTH-1:0] L2_PASS1_CFG_LEN = 10'd9;
55 localparam [ADDR_WIDTH-1:0] L2_PASS0_WT_BASE = 10'h039;
56 localparam [ADDR_WIDTH-1:0] L2_PASS0_WT_LEN = 10'd36;
57 localparam [ADDR_WIDTH-1:0] L2_PASS1_WT_BASE = 10'h05D;
58 localparam [ADDR_WIDTH-1:0] L2_PASS1_WT_LEN = 10'd36;
59 localparam [ADDR_WIDTH-1:0] L2_PASS0_OUT_BASE = 10'h231;
60 localparam [ADDR_WIDTH-1:0] L2_PASS0_OUT_LEN = 10'd196;
61 // L2 pass1 uses same base as pass0; stride-2 interleave done in top_sram_A
62 localparam [ADDR_WIDTH-1:0] L2_PASS1_OUT_BASE = 10'h231;
63 localparam [ADDR_WIDTH-1:0] L2_PASS1_OUT_LEN = 10'd196;
64
65 localparam [ADDR_WIDTH-1:0] L3_PASS0_CFG_BASE = 10'h01B;
66 localparam [ADDR_WIDTH-1:0] L3_PASS0_CFG_LEN = 10'd9;
67 localparam [ADDR_WIDTH-1:0] L3_PASS1_CFG_BASE = 10'h024;
68 localparam [ADDR_WIDTH-1:0] L3_PASS1_CFG_LEN = 10'd9;
69 localparam [ADDR_WIDTH-1:0] L3_PASS0_WT_BASE = 10'h081;
70 localparam [ADDR_WIDTH-1:0] L3_PASS0_WT_LEN = 10'd72;
71 localparam [ADDR_WIDTH-1:0] L3_PASS1_WT_BASE = 10'h0C9;
72 localparam [ADDR_WIDTH-1:0] L3_PASS1_WT_LEN = 10'd72;
73 localparam [ADDR_WIDTH-1:0] L3_DATA_BASE = 10'h231;
74 localparam [ADDR_WIDTH-1:0] L3_DATA_LEN = 10'd392;
75
76 // FC bias lives at 0x111 (the region freed by moving FC weights to SRAM_FCW).
77 // FC weights no longer live in SRAM_A at all -- LAYER_FC+SEL_WT is a no-op
78 // here; top_sram_A routes that transaction to the SRAM_FCW peer instead.
79 localparam [ADDR_WIDTH-1:0] FC_CFG_BASE = 10'h111;
80 localparam [ADDR_WIDTH-1:0] FC_CFG_LEN = 10'd10;
81
82 always @(*) begin
83     base_addr = {ADDR_WIDTH{1'b0}};
84     length = {ADDR_WIDTH{1'b0}};
85     read_mode = 1'b0;
86     write_mode = 1'b0;
87

```

```

88     case (layer_sel)
89         LAYER_PRELOAD: begin
90             write_mode = 1'b1;
91             case (data_sel)
92                 SEL_CFG: begin
93                     base_addr = PRELOAD_CFG_BASE;
94                     length    = PRELOAD_CFG_LEN;
95                 end
96                 SEL_WT: begin
97                     base_addr = PRELOAD_WT_BASE;
98                     length    = PRELOAD_WT_LEN;
99                 end
100                SEL_DATA: begin
101                    // Preload DATA no longer uses SRAM_A.
102                    write_mode = 1'b0;
103                end
104                SEL_FCW: begin
105                    // FC weight preload: top_sram_A redirects the write stream
106                    // to the SRAM_FCW packer. length counts host 32b beats.
107                    base_addr = {ADDR_WIDTH{1'b0}};
108                    length    = PRELOAD_FCW_LEN;
109                end
110                default: begin
111                    write_mode = 1'b0;
112                end
113            endcase
114        end
115
116        LAYER_L1: begin
117            read_mode = 1'b1;
118            case (data_sel)
119                SEL_CFG: begin
120                    base_addr = L1_CFG_BASE;
121                    length    = L1_CFG_LEN;
122                end
123                SEL_WT: begin
124                    base_addr = L1_WT_BASE;
125                    length    = L1_WT_LEN;
126                end
127                SEL_DATA: begin
128                    // L1 DATA no longer reads from SRAM_A.
129                    read_mode = 1'b0;
130                end
131                default: begin
132                    read_mode = 1'b0;
133                end
134            endcase
135        end
136
137        LAYER_L2: begin
138            case (data_sel)
139                SEL_CFG: begin
140                    read_mode = 1'b1;
141                    base_addr = pass_id ? L2_PASS1_CFG_BASE : L2_PASS0_CFG_BASE;

```

```

142         length    = pass_id ? L2_PASS1_CFG_LEN  : L2_PASSO_CFG_LEN;
143     end
144     SEL_WT: begin
145         read_mode = 1'b1;
146         base_addr = pass_id ? L2_PASS1_WT_BASE : L2_PASSO_WT_BASE;
147         length    = pass_id ? L2_PASS1_WT_LEN  : L2_PASSO_WT_LEN;
148     end
149     SEL_DATA: begin
150         write_mode = 1'b1;
151         base_addr = pass_id ? L2_PASS1_OUT_BASE : L2_PASSO_OUT_BASE;
152         length    = pass_id ? L2_PASS1_OUT_LEN : L2_PASSO_OUT_LEN;
153     end
154     default: begin
155     end
156 endcase
157 end
158
159 LAYER_L3: begin
160     read_mode = 1'b1;
161     case (data_sel)
162     SEL_CFG: begin
163         base_addr = pass_id ? L3_PASS1_CFG_BASE : L3_PASSO_CFG_BASE;
164         length    = pass_id ? L3_PASS1_CFG_LEN  : L3_PASSO_CFG_LEN;
165     end
166     SEL_WT: begin
167         base_addr = pass_id ? L3_PASS1_WT_BASE : L3_PASSO_WT_BASE;
168         length    = pass_id ? L3_PASS1_WT_LEN  : L3_PASSO_WT_LEN;
169     end
170     SEL_DATA: begin
171         base_addr = L3_DATA_BASE;
172         length    = L3_DATA_LEN;
173     end
174     default: begin
175         read_mode = 1'b0;
176     end
177 endcase
178 end
179
180 LAYER_FC: begin
181     case (data_sel)
182     SEL_CFG: begin
183         // FC bias lives in the freed old-FC-WT region of SRAM_A.
184         // Enable BOTH read and write so the same transaction can
185         // preload (write via pool_valid from host load_data) or
186         // fetch at inference time (read via data_ready from FC).
187         // Only one of read_en/write_en actually fires per cycle
188         // because pool_valid and data_ready are mutually exclusive
189         // in their respective modes.
190         read_mode = 1'b1;
191         write_mode = 1'b1;
192         base_addr = FC_CFG_BASE;
193         length    = FC_CFG_LEN;
194     end
195     SEL_WT: begin

```

```

196             // FC weights live in SRAM_FCW, not here.
197             // top_sram_A routes this transaction to the FCW peer and
198             // leaves the 32b SRAM_A idle; length=0 so Addr_Gen never
199             // asserts done on this path.
200             end
201             default: begin
202             end
203         endcase
204     end
205
206     default: begin
207     end
208 endcase
209 end
210
211 endmodule

```

#### A.10.4 SRAM A Wrapper

*Wraps the SRAM\_A model with active-low chip/write enables and read-valid timing.*

code/input/RTL/SRAM/sram\_A\_wrapper.v

```

1  `timescale 1ns/1ps
2  `define DELAY #1
3
4  module sram_A_wrapper # (
5      parameter AW = 10
6  ) (
7      input clk,
8      input rst_n,
9      input write_en,
10     input read_en,
11     input [AW-1:0] addr,
12     input [31:0] write_data,
13
14     output [31:0] read_data,
15     output reg read_valid
16 );
17     wire do_write;
18     wire do_read;
19
20     wire CEN;
21     wire WEN;
22     wire [AW-1:0] addr_delayed;
23     wire [31:0] write_data_delayed;
24
25     assign do_write = write_en & ~read_en;
26     assign do_read = read_en & ~write_en;
27
28     assign `DELAY CEN = ~(do_write | do_read);

```

```

29     assign 'DELAY WEN = ~do_write;
30     assign 'DELAY addr_delayed = addr;
31     assign 'DELAY write_data_delayed = write_data;
32
33     always @(posedge clk or negedge rst_n) begin
34         if (!rst_n)
35             read_valid <= 1'b0;
36         else
37             read_valid <= do_read;
38     end
39
40     sram_A_model sram_A_inst (
41         .CLK(clk),
42         .CEN(CEN),
43         .WEN(WEN),
44         .A(addr_delayed),
45         .D(write_data_delayed),
46         .EMA(3'b000),
47         .RETN(1'b1),
48         .Q(read_data)
49     );
50
51
52 endmodule

```

### A.10.5 SRAM B Controller

*Controls SRAM\_B read/write traffic for pooled convolution outputs and downstream layer inputs.*

code/input/RTL/SRAM/sram\_B\_controller.v

```

1  module sram_B_controller #(
2      parameter ADDR_WIDTH = 10
3  )(
4      input  [2:0] layer_sel,
5      input  [1:0] data_sel,
6      input          pass_id,
7
8      output reg [ADDR_WIDTH-1:0] base_addr,
9      output reg [ADDR_WIDTH-1:0] length,
10     output reg          read_mode,
11     output reg          write_mode
12 );
13
14 localparam [2:0] LAYER_PRELOAD = 3'd0;
15 localparam [2:0] LAYER_L1     = 3'd1;
16 localparam [2:0] LAYER_L2     = 3'd2;
17 localparam [2:0] LAYER_L3     = 3'd3;
18 localparam [2:0] LAYER_FC     = 3'd4;
19
20 localparam [1:0] SEL_CFG      = 2'd0;

```

```

21 localparam [1:0] SEL_WT      = 2'd1;
22 localparam [1:0] SEL_DATA   = 2'd2;
23
24 // SRAM_B only stores runtime feature maps in the current top-level memory plan.
25 // There is no preload transaction on SRAM_B.
26 // It is reused in time:
27 // - L1 pooled output shares the same 31x31x4 buffer used as L2 input.
28 // - L3/FC data traffic is routed by top_sram_B into a dedicated FC-order
29 //   buffer, so the FC stage no longer needs interleaved reads out of SRAM_B.
30
31 localparam [ADDR_WIDTH-1:0] L1_OUT_BASE      = 10'h000;
32 localparam [ADDR_WIDTH-1:0] L1_OUT_LEN      = 10'd961; // 31*31*4 bytes / 4
33
34 localparam [ADDR_WIDTH-1:0] L2_DATA_BASE     = 10'h000;
35 localparam [ADDR_WIDTH-1:0] L2_DATA_LEN     = 10'd961; // 31*31*4 bytes / 4
36
37 localparam [ADDR_WIDTH-1:0] L3_PASS0_OUT_BASE = 10'h000;
38 localparam [ADDR_WIDTH-1:0] L3_PASS0_OUT_LEN = 10'd36; // 6*6*4 bytes / 4
39 localparam [ADDR_WIDTH-1:0] L3_PASS1_OUT_BASE = 10'h024;
40 localparam [ADDR_WIDTH-1:0] L3_PASS1_OUT_LEN = 10'd36; // 6*6*4 bytes / 4
41
42 localparam [ADDR_WIDTH-1:0] FC_DATA_BASE     = 10'h000;
43 localparam [ADDR_WIDTH-1:0] FC_DATA_LEN     = 10'd72; // 6*6*8 bytes / 4
44
45 always @(*) begin
46     base_addr = {ADDR_WIDTH{1'b0}};
47     length    = {ADDR_WIDTH{1'b0}};
48     read_mode = 1'b0;
49     write_mode = 1'b0;
50
51     case (layer_sel)
52         LAYER_L1: begin
53             if (data_sel == SEL_DATA) begin
54                 write_mode = 1'b1;
55                 base_addr  = L1_OUT_BASE;
56                 length     = L1_OUT_LEN;
57             end
58         end
59
60         LAYER_L2: begin
61             if (data_sel == SEL_DATA) begin
62                 read_mode = 1'b1;
63                 base_addr  = L2_DATA_BASE;
64                 length     = L2_DATA_LEN;
65             end
66         end
67
68         LAYER_L3: begin
69             if (data_sel == SEL_DATA) begin
70                 write_mode = 1'b1;
71                 base_addr  = pass_id ? L3_PASS1_OUT_BASE : L3_PASS0_OUT_BASE;
72                 length     = pass_id ? L3_PASS1_OUT_LEN  : L3_PASS0_OUT_LEN;
73             end
74         end

```

```

75
76     LAYER_FC: begin
77         if (data_sel == SEL_DATA) begin
78             read_mode = 1'b1;
79             base_addr = FC_DATA_BASE;
80             length     = FC_DATA_LEN;
81         end
82     end
83
84     default: begin
85     end
86 endcase
87 end
88
89 endmodule

```

### A.10.6 SRAM B Wrapper

*Wraps the SRAM\_B model with the project-level SRAM control convention.*

```
code/input/RTL/SRAM/sram_B_wrapper.v
```

```

1  `timescale 1ns/1ps
2  `define DELAY #1
3
4  module sram_B_wrapper # (
5      parameter AW = 11
6  ) (
7      input  clk,
8      input  rst_n,
9      input  write_en,
10     input  read_en,
11     input  [AW-1:0] addr,
12     input  [31:0] write_data,
13
14     output [31:0] read_data,
15     output reg read_valid
16 );
17     wire do_write;
18     wire do_read;
19
20     wire CEN;
21     wire WEN;
22     wire [AW-1:0] addr_delayed;
23     wire [31:0] write_data_delayed;
24
25     assign do_write = write_en & ~read_en;
26     assign do_read  = read_en  & ~write_en;
27
28     assign `DELAY CEN = ~(do_write | do_read);
29     assign `DELAY WEN = ~do_write;

```

```

30     assign 'DELAY addr_delayed = addr;
31     assign 'DELAY write_data_delayed = write_data;
32
33     always @(posedge clk or negedge rst_n) begin
34         if (!rst_n)
35             read_valid <= 1'b0;
36         else
37             read_valid <= do_read;
38     end
39
40     sram_B_model sram_B_inst (
41         .CLK(clk),
42         .CEN(CEN),
43         .WEN(WEN),
44         .A(addr_delayed),
45         .D(write_data_delayed),
46         .EMA(3'b000),
47         .RETN(1'b1),
48         .Q(read_data)
49     );
50
51
52 endmodule

```

## A.10.7 SRAM FCW Wrapper

*Wraps the wide FC weight SRAM used by the classifier.*

code/input/RTL/SRAM/sram\_FCW\_wrapper.v

```

1  'timescale 1ns/1ps
2
3  // sram_FCW_wrapper: dedicated FC-weight memory.
4  //   Depth  = 288 words (FC input K = 6*6*8)
5  //   Width  = 80 bits  (10 output channels x 8-bit INT8 kernel)
6  //
7  // Written so Quartus (Cyclone V M10K) and Xilinx Vivado both infer
8  // BRAM. Keys to inference:
9  //   (1) memory read/write in ONE tiny always block, no reset.
10 //   (2) no reset of 'douta' inside the RAM block -- the M10K output
11 //       register cannot be reset and still map to BRAM. Any reset of
12 //       'douta' forces the tool into FFs.
13 //   (3) 'douta_valid' tracking is in a SEPARATE always block so it
14 //       does not contaminate the RAM template.
15 //   (4) 80-bit width > M10K's 40-bit ceiling -> Quartus automatically
16 //       concatenates two M10K blocks (fine as long as (1)-(3) hold).
17
18 module sram_FCW_wrapper # (
19     parameter AW = 9,          // 2^9 = 512 >= 288
20     parameter DW = 80
21 ) (

```

```

22     input          clka,
23     input          rsta,      // sync active-high reset (valid flag only)
24     input          ena,
25     input          wea,
26     input [AW-1:0] addra,
27     input [DW-1:0] dina,
28
29     output reg [DW-1:0] douta,
30     output reg      douta_valid
31 );
32 // -----
33 // RAM block: plain sync write / sync read, no reset (BRAM template)
34 // -----
35 (* ramstyle = "M1OK" *)
36 reg [DW-1:0] mem [0:(1<<AW)-1];
37
38 always @(posedge clka) begin
39     if (ena) begin
40         if (wea)
41             mem[addra] <= dina;
42         else
43             douta <= mem[addra];
44     end
45 end
46
47 // -----
48 // Valid flag: separate always block so it does NOT inhibit
49 // BRAM inference of 'mem' / 'douta'.
50 // -----
51 always @(posedge clka) begin
52     if (rsta)
53         douta_valid <= 1'b0;
54     else if (ena && !wea)
55         douta_valid <= 1'b1;
56     else
57         douta_valid <= 1'b0;
58 end
59
60 endmodule

```

## A.10.8 Top SRAM A Integration

*Integrates SRAM\_A controller, address generation, wrapper, and stream-facing control.*

```
code/input/RTL/SRAM/top_sram_A.v
```

```

1 module top_sram_A(
2     input clk,
3     input rst_n,
4
5     input [2:0] layer_sel,

```

```

6     input [1:0] data_sel,
7     input      pass_id,
8     input      start,
9     output     busy,
10    output     done,
11
12    // Read path from SRAM_A to conv / pool pipeline (32b, shared).
13    input      data_ready,
14    output [31:0] read_data,
15    output     data_valid,
16    output     data_last,
17
18    // Write path from pool pipeline back to SRAM_A (32b, shared).
19    output     pool_ready,
20    input signed [31:0] pool_data,
21    input      pool_valid,
22    input      pool_last,
23
24    // FC weight read path (80b, dedicated SRAM_FCW).
25    input      fc_wt_ready,
26    output [79:0] fc_wt_read_data,
27    output     fc_wt_valid,
28    output     fc_wt_last
29 );
30
31 localparam ADDR_WIDTH = 10;
32 localparam FCW_ADDR_WIDTH = 9; // 512 depth covers 288 entries
33 localparam FCW_DATA_WIDTH = 80; // 10 x 8-bit kernels
34 localparam FCW_LEN = 10'd288;
35
36 // SRAM_A controller localparams -- mirror sram_a_controller encodings.
37 localparam [2:0] LAYER_PRELOAD = 3'd0;
38 localparam [2:0] LAYER_FC = 3'd4;
39 localparam [1:0] SEL_CFG = 2'd0;
40 localparam [1:0] SEL_WT = 2'd1;
41 localparam [1:0] SEL_FCW = 2'd3;
42
43 // -----
44 // Classify incoming transaction
45 // -----
46 wire fcw_preload_txn = (layer_sel == LAYER_PRELOAD) && (data_sel == SEL_FCW);
47 wire fcw_read_txn = (layer_sel == LAYER_FC) && (data_sel == SEL_WT);
48 wire fcw_txn = fcw_preload_txn | fcw_read_txn;
49
50 // =====
51 // Existing 32b SRAM_A path (conv L1/L2/L3, FC CFG, preload CFG/WT)
52 // =====
53
54 wire [ADDR_WIDTH-1:0] addr;
55 wire [ADDR_WIDTH-1:0] base_addr;
56 wire [ADDR_WIDTH-1:0] length;
57 wire      read_mode;
58 wire      write_mode;
59 wire      read_en;

```

```

60 wire      write_en;
61 wire      addr_step;
62 wire      txn_done;
63 wire      pool_ready_int;
64 wire      unused_pool_last;
65
66 reg       read_last_d;
67 reg       txn_active;
68
69 // Active-transaction flag is shared: only ONE of (32b SRAM_A, SRAM_FCW) runs
70 // at a time. fcw_txn diverts start/done to the FCW datapath below.
71 wire      a32_start = start && !fcw_txn;
72
73 assign addr_step = read_en | write_en;
74 assign read_en   = txn_active && read_mode  && data_ready;
75 assign write_en  = txn_active && write_mode && pool_valid;
76 assign pool_ready_int = txn_active && write_mode;
77 assign data_last  = read_last_d & data_valid;
78 assign unused_pool_last = pool_last;
79
80 always @(posedge clk or negedge rst_n) begin
81     if (!rst_n)
82         read_last_d <= 1'b0;
83     else
84         read_last_d <= txn_done && read_en;
85 end
86
87 always @(posedge clk or negedge rst_n) begin
88     if (!rst_n)
89         txn_active <= 1'b0;
90     else if (a32_start)
91         txn_active <= (length != {ADDR_WIDTH{1'b0}}) && (read_mode || write_mode);
92     else if (txn_done)
93         txn_active <= 1'b0;
94 end
95
96 Addr_Gen #(
97     .ADDR_WIDTH(ADDR_WIDTH)
98 ) u_addr_gen (
99     .clk(clk),
100    .rst_n(rst_n),
101    .start(a32_start),
102    .enable(addr_step),
103    .base_addr(base_addr),
104    .length(length),
105    .addr(addr),
106    .done(txn_done)
107 );
108
109 // L2 DATA write interleave: pass0 writes even offsets, pass1 writes odd offsets.
110 wire l2_data_write = (layer_sel == 3'd2) && (data_sel == 2'd2) && write_mode;
111 wire [ADDR_WIDTH-1:0] counter_val = addr - base_addr;
112 wire [ADDR_WIDTH-1:0] addr_final  = l2_data_write
113     ? (base_addr + {counter_val[ADDR_WIDTH-2:0], pass_id})

```

```

114     : addr;
115
116 sram_A_wrapper #(
117     .AW(ADDR_WIDTH)
118 ) u_sram_A_wrapper (
119     .clk(clk),
120     .rst_n(rst_n),
121     .write_en(write_en),
122     .read_en(read_en),
123     .addr(addr_final),
124     .write_data(pool_data),
125     .read_data(read_data),
126     .read_valid(data_valid)
127 );
128
129 sram_A_controller #(
130     .ADDR_WIDTH(ADDR_WIDTH)
131 ) u_sram_A_controller(
132     .layer_sel(layer_sel),
133     .data_sel(data_sel),
134     .pass_id(pass_id),
135     .base_addr(base_addr),
136     .length(length),
137     .read_mode(read_mode),
138     .write_mode(write_mode)
139 );
140
141 // =====
142 // SRAM_FCW dedicated 80b path (FC weight preload + inference read)
143 // =====
144
145 wire          fcw_start    = start && fcw_txn;
146 reg          fcw_active;
147 reg          fcw_rd_mode; // 1 = inference read, 0 = preload write
148
149 wire          fcw_rd_done;
150 wire [FCW_ADDR_WIDTH-1:0] fcw_rd_addr;
151 wire [FCW_ADDR_WIDTH-1:0] fcw_rd_base = {FCW_ADDR_WIDTH{1'b0}};
152 wire [FCW_ADDR_WIDTH-1:0] fcw_rd_len  = FCW_LEN[FCW_ADDR_WIDTH-1:0];
153
154 // Read-side enable: on inference, advance when adapter is ready.
155 wire fcw_read_en = fcw_active && fcw_rd_mode && fc_wt_ready;
156
157 // Read address generator
158 Addr_Gen #(
159     .ADDR_WIDTH(FCW_ADDR_WIDTH)
160 ) u_fcw_read_addr (
161     .clk(clk),
162     .rst_n(rst_n),
163     .start(fcw_start && fcw_read_txn),
164     .enable(fcw_read_en),
165     .base_addr(fcw_rd_base),
166     .length(fcw_rd_len),
167     .addr(fcw_rd_addr),

```

```

168     .done(fcw_rd_done)
169 );
170
171 // Track last flag for 80b reads: pulses in the cycle the registered data
172 // for the final FCW entry becomes valid on the wrapper output.
173 reg fcw_read_last_d;
174
175 always @(posedge clk or negedge rst_n) begin
176     if (!rst_n)
177         fcw_read_last_d <= 1'b0;
178     else
179         fcw_read_last_d <= fcw_rd_done && fcw_read_en;
180 end
181
182 // Preload path: fcw_preload_packer accepts 32b host words (from pool_data)
183 // and emits 80b writes to SRAM_FCW. Beat count is driven by length=864 from
184 // the controller, but we manage it locally here via pool_valid beats.
185 reg [9:0] fcw_pl_beat_cnt; // up to 864 beats
186 wire      fcw_pl_en = fcw_active && !fcw_rd_mode && pool_valid;
187
188 always @(posedge clk or negedge rst_n) begin
189     if (!rst_n) begin
190         fcw_pl_beat_cnt <= 10'd0;
191     end else if (fcw_start && fcw_preload_txn) begin
192         fcw_pl_beat_cnt <= 10'd0;
193     end else if (fcw_pl_en) begin
194         fcw_pl_beat_cnt <= fcw_pl_beat_cnt + 10'd1;
195     end
196 end
197
198 wire fcw_pl_done = fcw_active && !fcw_rd_mode && (fcw_pl_beat_cnt == 10'd863) && pool_valid;
199
200 // fcw_active tracks the FCW transaction
201 always @(posedge clk or negedge rst_n) begin
202     if (!rst_n)
203         fcw_active <= 1'b0;
204     else if (fcw_start)
205         fcw_active <= 1'b1;
206     else if (fcw_rd_done || fcw_pl_done)
207         fcw_active <= 1'b0;
208 end
209
210 always @(posedge clk or negedge rst_n) begin
211     if (!rst_n)
212         fcw_rd_mode <= 1'b0;
213     else if (fcw_start)
214         fcw_rd_mode <= fcw_read_txn;
215 end
216
217 // Packer instance (only driven during preload)
218 wire      packer_wea;
219 wire [FCW_ADDR_WIDTH-1:0] packer_addr;
220 wire [FCW_DATA_WIDTH-1:0] packer_dina;
221 wire      packer_up_ready;

```

```

222
223 fcw_preload_packer #(
224     .AW(FCW_ADDR_WIDTH),
225     .DW(FCW_DATA_WIDTH)
226 ) u_fcw_packer (
227     .clk(clk),
228     .rst_n(rst_n),
229     .start(fcw_start && fcw_preload_txn),
230     .up_valid(fcw_pl_en),
231     .up_data(pool_data),
232     .up_ready(packer_up_ready),
233     .wea(packer_wea),
234     .addra(packer_addr),
235     .dina(packer_dina)
236 );
237
238 // SRAM_FCW access port: mux read vs. write
239 wire          fcw_ena   = fcw_read_en | packer_wea;
240 wire          fcw_wea   = packer_wea;
241 wire [FCW_ADDR_WIDTH-1:0] fcw_addr = packer_wea ? packer_addr : fcw_rd_addr;
242 wire [FCW_DATA_WIDTH-1:0] fcw_dina = packer_dina;
243
244 wire [FCW_DATA_WIDTH-1:0] fcw_douta;
245 wire          fcw_douta_valid;
246
247 sram_FCW_wrapper #(
248     .AW(FCW_ADDR_WIDTH),
249     .DW(FCW_DATA_WIDTH)
250 ) u_sram_FCW (
251     .clka(clk),
252     .rsta(~rst_n),
253     .ena(fcw_ena),
254     .wea(fcw_wea),
255     .addra(fcw_addr),
256     .dina(fcw_dina),
257     .douta(fcw_douta),
258     .douta_valid(fcw_douta_valid)
259 );
260
261 assign fc_wt_read_data = fcw_douta;
262 assign fc_wt_valid     = fcw_douta_valid & fcw_rd_mode;
263 assign fc_wt_last     = fcw_read_last_d;
264
265 // =====
266 // Shared done / busy
267 // =====
268 assign busy          = txn_active | fcw_active;
269 // For FCW reads, use Addr_Gen.done pulse (fires on the last read enable) so
270 // done aligns with the existing 32b SRAM_A convention (done coincident with
271 // the last valid data beat being consumed).
272 wire fcw_done_pulse = fcw_active
273     ? (fcw_rd_mode ? fcw_rd_done : fcw_pl_done)
274     : 1'b0;
275 wire a32_write_done = write_en && txn_done;

```

```

276 wire a32_read_done = read_mode && data_last;
277 assign done        = fcw_done_pulse
278                   | (!fcw_active && (a32_write_done | a32_read_done));
279
280 // pool_ready: during non-FCW transactions behaves as before; during FCW
281 // preload the packer always accepts (up_ready=1) so we echo that through.
282 assign pool_ready = fcw_active
283                   ? (!fcw_rd_mode ? packer_up_ready : 1'b0)
284                   : pool_ready_int;
285
286 endmodule

```

### A.10.9 Top SRAM B Integration

*Integrates SRAM\_B controller, address generation, wrapper, and pool/input stream routing.*

code/input/RTL/SRAM/top\_sram\_B.v

```

1  module top_sram_B(
2      input clk,
3      input rst_n,
4
5      input [2:0] layer_sel,
6      input [1:0] data_sel,
7      input      pass_id,
8      input      start,
9      output     busy,
10     output     done,
11
12     // Read path from SRAM_B to conv / FC pipeline.
13     input      data_ready,
14     output [31:0] read_data,
15     output     data_valid,
16     output     data_last,
17
18     // Write path from pool pipeline back to SRAM_B.
19     output     pool_ready,
20     input signed [31:0] pool_data,
21     input      pool_valid,
22     input      pool_last
23 );
24
25 localparam [2:0] LAYER_L3 = 3'd3;
26 localparam [2:0] LAYER_FC = 3'd4;
27 localparam [1:0] SEL_DATA = 2'd2;
28 localparam integer FC_BUF_AW = 7; // 2^7 = 128 >= 72 words
29 localparam [5:0] FC_PASS_LEN = 6'd36;
30 localparam [6:0] FC_FULL_LEN = 7'd72;
31
32 wire [9:0] addr;
33 wire [9:0] base_addr;

```

```

34 wire [9:0] length;
35 wire      read_mode;
36 wire      write_mode;
37 wire      read_en;
38 wire      write_en;
39 wire      addr_step;
40 wire      txn_done;
41 wire      pool_ready_int;
42 wire      unused_pool_last;
43 wire [31:0] main_read_data;
44 wire      main_data_valid;
45 wire      main_data_last;
46
47 reg      read_last_d;
48 reg      txn_active;
49
50 // L3 output -> FC input reorder buffer.
51 // Instead of storing pass0[0..35], pass1[0..35] and fixing the order on the
52 // FC read side, write pass0 to even addresses and pass1 to odd addresses.
53 // This trades a little extra BRAM for a simpler FC read path.
54 wire fc_buf_write_mirror = (layer_sel == LAYER_L3) && (data_sel == SEL_DATA) && write_mode;
55 wire fc_buf_read_txn     = (layer_sel == LAYER_FC) && (data_sel == SEL_DATA) && read_mode;
56 wire main_start         = start && !fc_buf_read_txn;
57
58 assign addr_step        = read_en | write_en;
59 assign unused_pool_last = pool_last;
60 assign read_en          = txn_active && read_mode && data_ready;
61 assign write_en         = txn_active && write_mode && pool_valid;
62 assign pool_ready_int  = txn_active && write_mode;
63 assign main_data_last  = read_last_d & main_data_valid;
64
65 always @(posedge clk or negedge rst_n) begin
66     if (!rst_n)
67         read_last_d <= 1'b0;
68     else
69         read_last_d <= txn_done && read_en;
70 end
71
72 // Assumes the FSM holds layer_sel / data_sel / pass_id stable until txn_done.
73 always @(posedge clk or negedge rst_n) begin
74     if (!rst_n)
75         txn_active <= 1'b0;
76     else if (main_start)
77         txn_active <= (length != 10'd0) && (read_mode || write_mode);
78     else if (txn_done)
79         txn_active <= 1'b0;
80 end
81
82 Addr_Gen #(
83     .ADDR_WIDTH(10)
84 ) u_addr_gen (
85     .clk(clk),
86     .rst_n(rst_n),
87     .start(main_start),

```

```

88     .enable(addr_step),
89     .base_addr(base_addr),
90     .length(length),
91     .addr(addr),
92     .done(txn_done)
93 );
94
95 sram_B_wrapper #(
96     .AW(10)
97 ) u_sram_B_wrapper (
98     .clk(clk),
99     .rst_n(rst_n),
100    .write_en(write_en),
101    .read_en(read_en),
102    .addr(addr),
103    .write_data(pool_data),
104    .read_data(main_read_data),
105    .read_valid(main_data_valid)
106 );
107
108 sram_B_controller u_sram_B_controller(
109     .layer_sel(layer_sel),
110     .data_sel(data_sel),
111     .pass_id(pass_id),
112     .base_addr(base_addr),
113     .length(length),
114     .read_mode(read_mode),
115     .write_mode(write_mode)
116 );
117
118 // Dedicated FC-order buffer. L3 still writes the original SRAM_B path; we
119 // mirror those writes into a second BRAM organized for FC sequential reads.
120 // pass0 writes even slots, pass1 writes odd slots, so FC can read 0..71
121 // with no address interleave logic.
122 (* ramstyle = "M10K" *)
123 reg [31:0] fc_input_mem [0:(1<<FC_BUF_AW)-1];
124
125 reg      fc_buf_rd_active;
126 reg [6:0] fc_buf_rd_count;
127 reg [31:0] fc_buf_read_data_q;
128 reg      fc_buf_data_valid_q;
129 reg      fc_buf_read_last_d;
130
131 wire [9:0] counter_val = addr - base_addr;
132 wire      fc_buf_write_en = write_en && fc_buf_write_mirror;
133 wire      fc_buf_read_en  = fc_buf_rd_active && data_ready;
134 wire [6:0] fc_buf_wr_addr  = {counter_val[5:0], 1'b0} + pass_id;
135 wire      fc_buf_rd_done  = fc_buf_read_en  && (fc_buf_rd_count == (FC_FULL_LEN - 1'b1));
136 wire      fc_buf_data_last = fc_buf_data_valid_q && fc_buf_read_last_d;
137
138 always @(posedge clk or negedge rst_n) begin
139     if (!rst_n) begin
140         fc_buf_rd_active  <= 1'b0;
141         fc_buf_rd_count  <= 7'd0;

```

```

142     fc_buf_read_data_q <= 32'd0;
143     fc_buf_data_valid_q <= 1'b0;
144     fc_buf_read_last_d <= 1'b0;
145 end else begin
146     fc_buf_data_valid_q <= 1'b0;
147     fc_buf_read_last_d <= 1'b0;
148
149     if (fc_buf_write_en)
150         fc_input_mem[fc_buf_wr_addr] <= pool_data;
151
152     if (start && fc_buf_read_txn) begin
153         fc_buf_rd_active <= 1'b1;
154         fc_buf_rd_count <= 7'd0;
155     end else if (fc_buf_read_en) begin
156         fc_buf_read_data_q <= fc_input_mem[fc_buf_rd_count];
157         fc_buf_data_valid_q <= 1'b1;
158         fc_buf_read_last_d <= fc_buf_rd_done;
159         if (fc_buf_rd_done)
160             fc_buf_rd_active <= 1'b0;
161         else
162             fc_buf_rd_count <= fc_buf_rd_count + 7'd1;
163     end
164 end
165 end
166
167 assign busy      = fc_buf_read_txn ? fc_buf_rd_active : txn_active;
168 assign done      = fc_buf_read_txn ? fc_buf_data_last
169                 : (read_mode ? main_data_last : (write_en && txn_done));
170 assign pool_ready = pool_ready_int;
171 assign read_data  = fc_buf_read_txn ? fc_buf_read_data_q : main_read_data;
172 assign data_valid = fc_buf_read_txn ? fc_buf_data_valid_q : main_data_valid;
173 assign data_last  = fc_buf_read_txn ? fc_buf_data_last : main_data_last;
174
175 endmodule

```