# COLUMBIA UNIVERSITY
## IN THE CITY OF NEW YORK

**PhD** Thesis Defense

# Analytical Query Execution Optimized for all Layers of Modern Hardware

*Orestis Polychroniou*

Ορέστης Πολυχρονίου

# Big Data

- **Volume** and **value** of (big) data
  - 20 **zettabytes** of data by 2020
  - $125 **billion** in 2015
  - $40 **billion** for databases

- Relational Analytics
  - Business Intelligence
  - Decision Support
  - > $10 **billion** market

# Database Systems

- Disk-based / Traditional DBMS
    - Data on (hard) **disk**
    - Query execution **disk-bound**
    - **Not** very distributed (e.g. Oracle)

*new hardware* **!**

- In-Memory / "**Modern**" DBMS
    - Data (mostly) in **RAM**
    - Query execution **memory-bound**
    - **Very** distributed (e.g. cloud)

# Impact of Hardware

- Traditional —> Modern DBMS

  - Driven by **hardware** advances !

- Hardware advances affecting databases

  - Large main memory capacity

  - Complex multi-core processors

  - Scalable memory hierarchy (including fast networks)

- How can we achieve **high** performance in a **modern** database ?

  - Database system **specialization**

  - **Adapting** to the hardware dynamics

# Modern Database Specialization

* Transactional DBMS

    * Focus on **transactions**

    * **Update** a **few tuples** per transaction

    * **Row**-store

* **Analytical** DBMS ·······························▶

    * Focus on **queries** for **analysis**

    * **Read** a **few columns** from **many tuples** per query

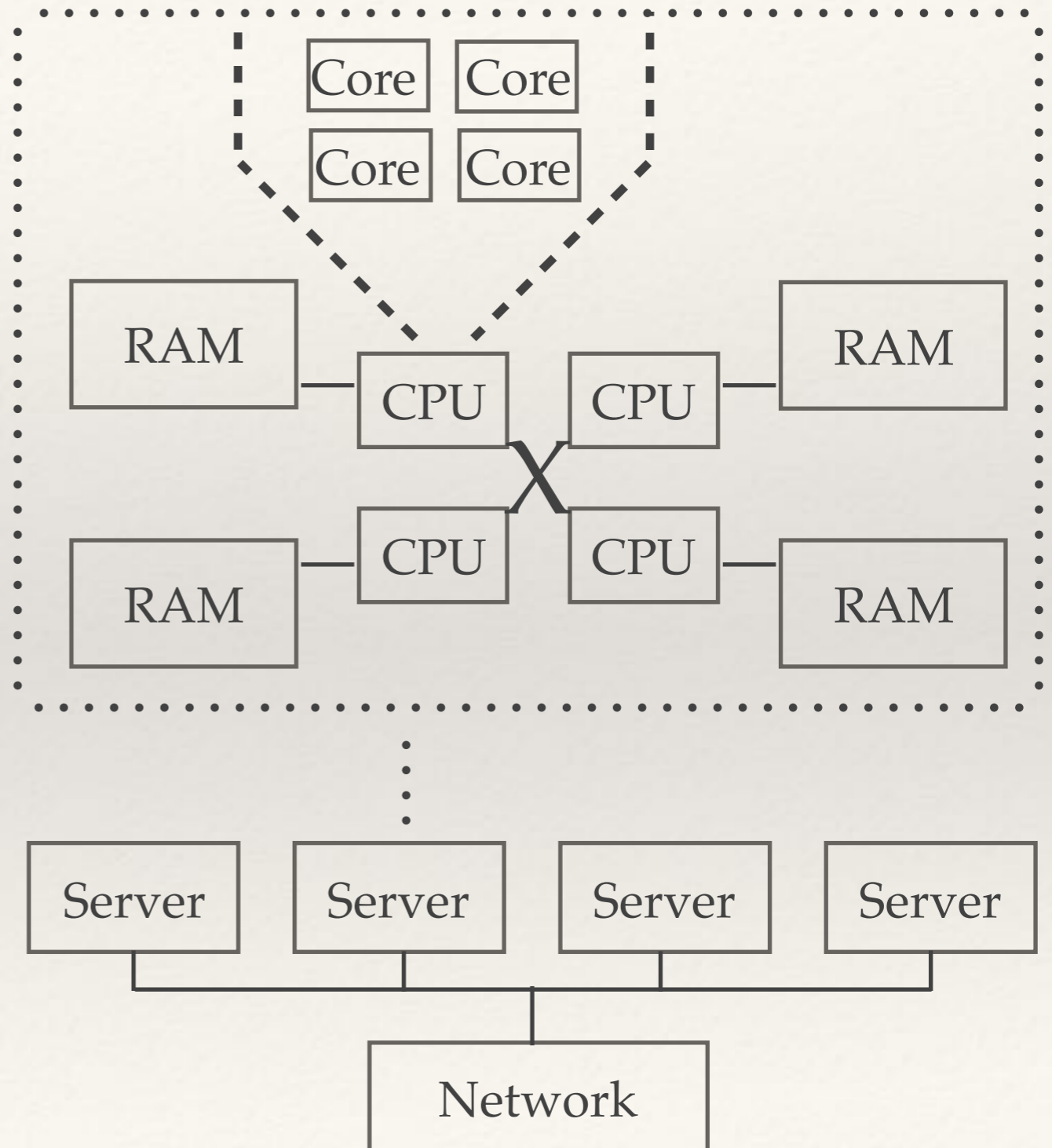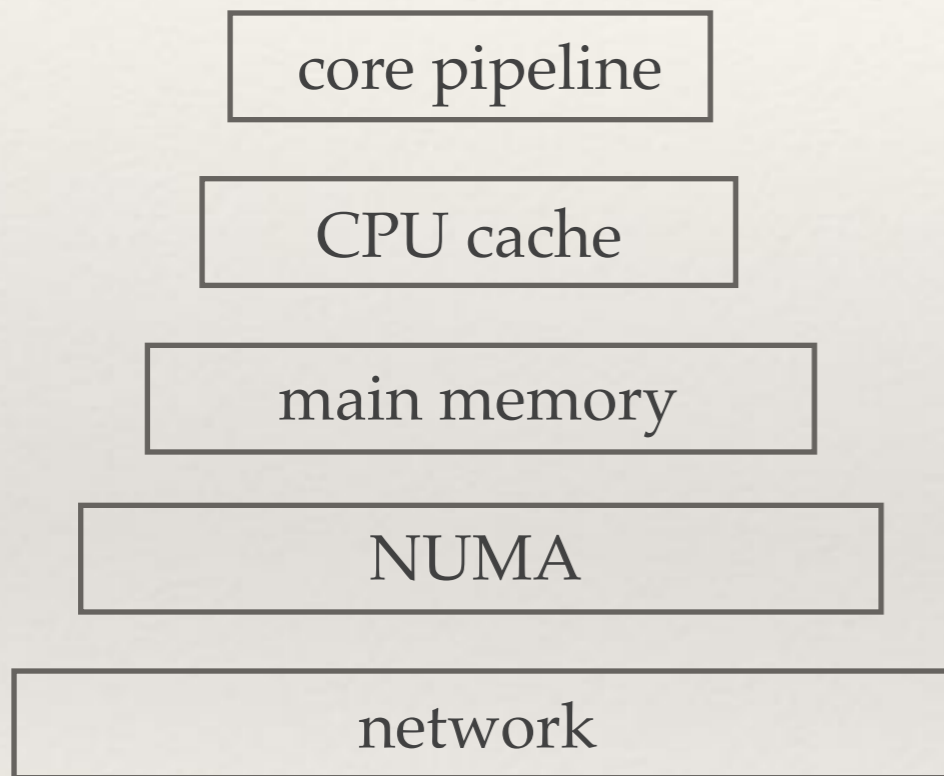    * **Column**-store

* Others (e.g. scientific, graph, …)

# Research Statement

❖ How can we **increase performance** in a **modern analytical database ?**

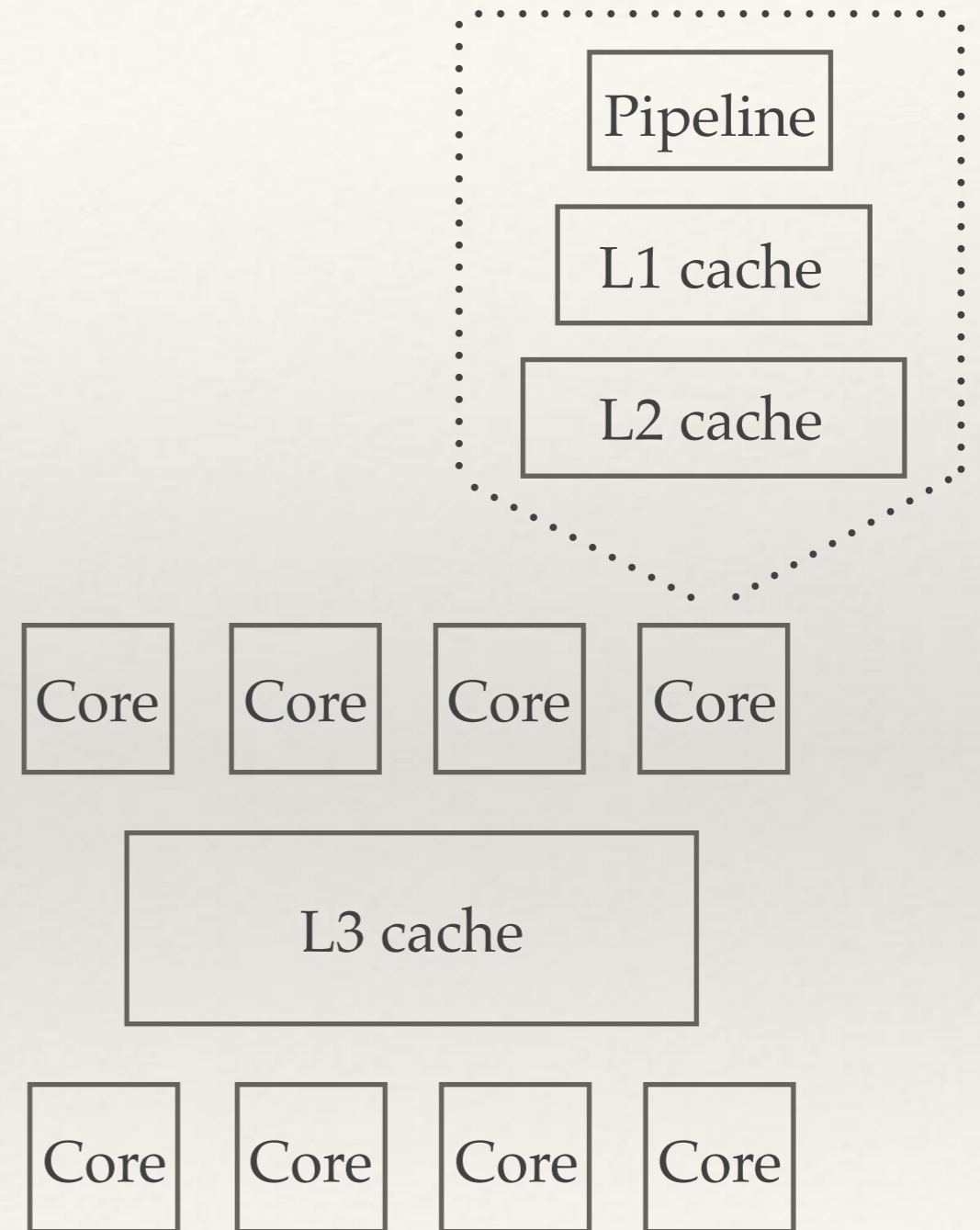　❖ By optimizing **analytical query execution** for each **hardware layer !**

❖ **Why** does it work **?**

　❖ Hardware has **always** driven database **design** & **implementation** (disks originally)

　❖ Hardware becomes more **complex** making hardware-oblivious designs **ineffective**

　❖ **Our** solutions are **hardware-conscious** and utilize complex modern hardware **features**

　❖ Hardware becomes increasingly **parallel** to efficiently process **larger** datasets

　❖ **Our** solutions push the boundaries of parallelism (via data parallelism, many-cores, …)

# Layers of Modern Hardware

core pipeline

CPU cache

main memory

NUMA

network

Core Core
Core Core

RAM — CPU CPU — RAM

RAM — CPU CPU — RAM

Server    Server    Server    Server

Network

# Modern Mainstream CPUs

- Thread parallelism
  - Multiple cores
  - Multiple threads per core

- Instruction level parallelism
  - Out-of-order execution
  - Super-scalar pipeline

- Data parallelism
  - SIMD vectorization

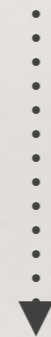Pipeline

L1 cache

L2 cache

Core  Core  Core  Core

L3 cache

Core  Core  Core  Core

# SIMD Vectorization

❖ As compiler optimization

  ❖ Works for **simple** loops only

  ❖ Insufficient for database operators

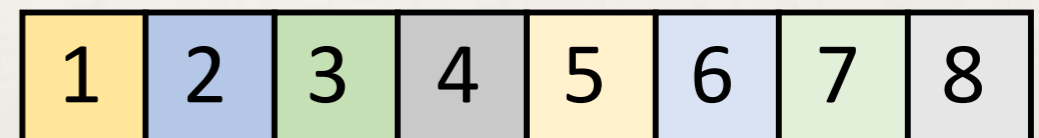```
for (i = 0; i < n; ++i) {
    c[i] = a[i] + b[i];
}
```
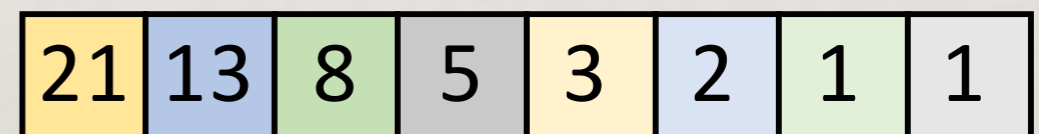
*scalar code*

*SIMD code*

```
for (i = 0; i < n; i += 16) {
    __m512i x = _mm512_load_si512(&a[i]);
    __m512i y = _mm512_load_si512(&b[i]);
    __m512i z = _mm512_add_epi32(x, y);
    _mm512_store_si512(&c[i], z);
}
```

*8-way SIMD addition*

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |

**+**

| 21 | 13 | 8 | 5 | 3 | 2 | 1 | 1 |

**=**

| 22 | 15 | 11 | 9 | 8 | 8 | 8 | 9 |

# Bottlenecks of Query Execution

❖ Network-bound

  ❖ Distributed joins with minimal network traffic (Part 1)

core pipeline

CPU cache

main memory

NUMA

network
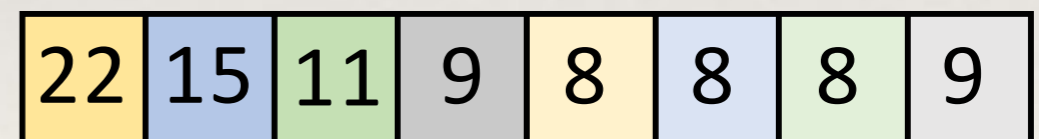
# Bottlenecks of Query Execution

❖ Network-bound

  ❖ Distributed joins with minimal network traffic (Part 1)

❖ Memory-bound (random access)

  ❖ (Cache-RAM-NUMA)-aware partitioning (Part 2)

| core pipeline |
| CPU cache |
| main memory |
| NUMA |
| network |

# Bottlenecks of Query Execution

❖ Network-bound

   ❖ Distributed joins with minimal network traffic (Part 1)

❖ Memory-bound (random access)

   ❖ (Cache-RAM-NUMA)-aware partitioning (Part 2)

❖ Memory-bound (sequential access)

   ❖ Lightweight in-memory compression (Part 3)

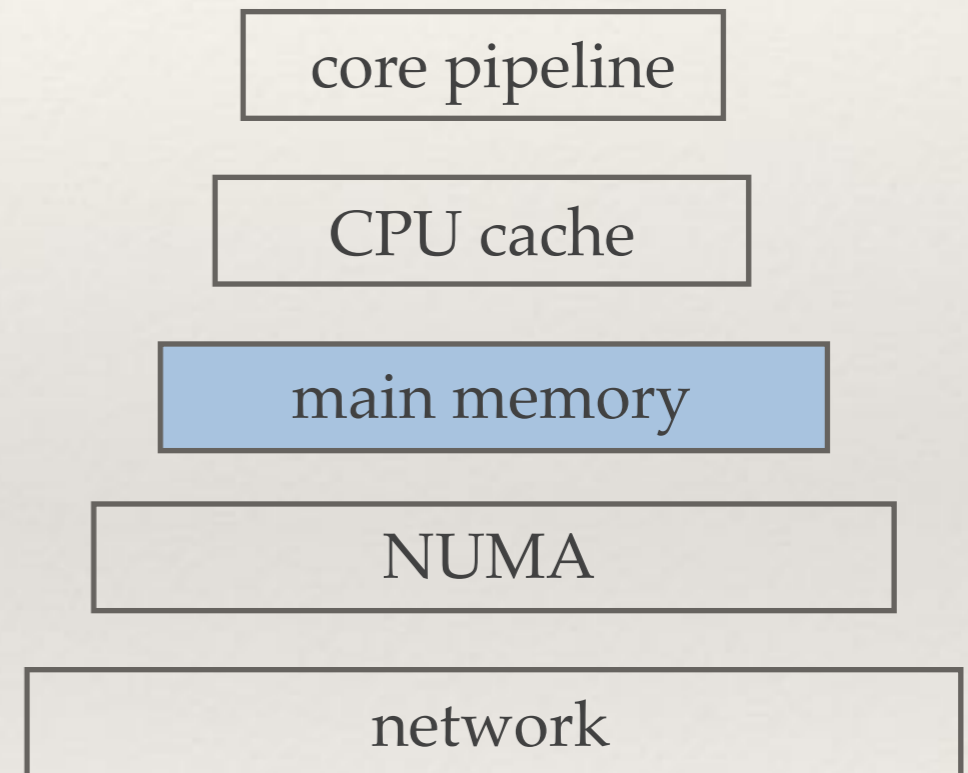| core pipeline |
| CPU cache |
| main memory |
| NUMA |
| network |

# Bottlenecks of Query Execution

- Network-bound
  - Distributed joins with minimal network traffic (Part 1)

- Memory-bound (random access)
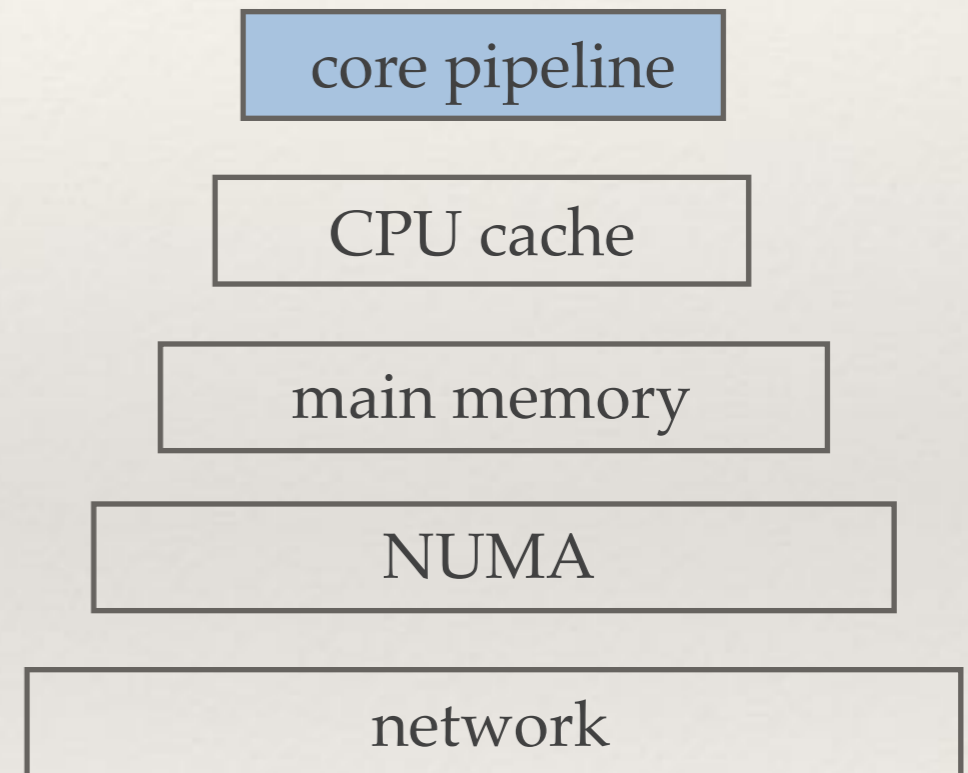  - (Cache-RAM-NUMA)-aware partitioning (Part 2)

- Memory-bound (sequential access)
  - Lightweight in-memory compression (Part 3)

- Compute-bound
  - Advanced SIMD vectorization techniques (Part 4)

| core pipeline |
| CPU cache |
| main memory |
| NUMA |
| network |

# Part 1: Network-bound



core pipeline

CPU cache

main memory

NUMA

network

Core  Core
Core  Core

RAM  CPU  CPU  RAM

RAM  CPU  CPU  RAM

Server  Server  Server  Server

Network

# Network << RAM

- Network is slower than RAM

  - 4-channel DDR4 x4 = ~200 GB/s

  - FDR InfiniBand 4X = ~5 GB/s

- Optimize for network traffic

  - Joins are **dominant**

  - Disk hash join ~ network hash join

  - **Our contribution**: track join

Bandwidth (GB/sec)

| | 200 | 160 | 120 | 80 | 40 | 0 |
|---|---|---|---|---|---|---|

2-channel DDR4 x1 | 4-channel DDR4 x1 | 4-channel DDR4 x4

5 | 4 | 3 | 2 | 1 | 0

10 Gbit Ethernet | QDR InfiniBand 4X | FDR InfiniBand 4X

# Previous Work: Broadcast Join



Network Nodes

1    2    3    4

Input Relations

R

S

- Good if one table is small

- Bad for large number of nodes

# Previous Work: Hash Join

Network Nodes



- ❖ Good if both tables are large
- ❖ Bad if one table is small

# Track Join: Minimize Network Traffic

❖ Basic idea

    ❖ Logical decomposition into **Cartesian product** joins

    ❖ **Optimize the transfers** for each Cartesian product

❖ Basic steps

    ❖ **Track** tuple locations per unique join key

    ❖ Generate **optimal** transfer schedule per key

    ❖ Transfer data and execute join

❖ Multiple variants

    ❖ 2-phase, 3-phase, 4-phase

# Track Join



Network Nodes

1    2    3    4

Input Relations

**R**

**S**

*Hash partition **unique** keys*

❖ Tracking phase

   ❖ Like a hash join of **keys only**

# Track Join (2-phase)

Network Nodes

1  2  3  4

**R**

*key=X
nodes=[1,3,4]*

Input Relations

*R tuples
(key=X)*

*R tuples
(key=X)*

*R tuples
(key=X)*

**S**

1  2  3  4

❖ Selective broadcast

   ❖ On locations that have **at least one** tuple

# Hash Join & Track Join

❖ Hash Join (network cost = 10)



❖ 2-phase Track Join (network cost = 12)



❖ 3-phase Track Join (network cost = 8)



❖ 4-phase Track Join (network cost = 6)

# Real Workload 1-1

- ❖ Real workload 1
  - ❖ 1–1 join
  - ❖ Pre-existing **locality**



*track join variants*

Legend:
- S Tuples
- R Tuples
- Drive the schedule
- Tracking phase

Y-axis: Network Traffic (GB), 0 to 20

X-axis: HJ, 2TJ-R, 2TJ-S, 3TJ, 4TJ

# Real Workload M-N

❖ Real workload 2

   ❖ M—N join i.e., output = **~5X** inputs

   ❖ Pre-existing **locality**

*broadcasts large table*

**Network Traffic (GB)** — chart, categories: HJ, 2TJ-R, 2TJ-S, 3TJ, 4TJ; y-axis 0 to 40.

Legend:
- S Tuples
- R Tuples
- Drive the schedule
- Tracking phase

*insignificant tracking cost*

# No locality? Use 4TJ

❖ Real workload 2

  ❖ M—N join i.e., output = **~5X** inputs

  ❖ Shuffle to **remove** locality



Network Traffic (GB)

40
30
20
10
0

HJ    2TJ-R    2TJ-S    3TJ    4TJ

S Tuples
R Tuples
Drive the schedule
Tracking phase

*still better than hash join*

# CPU+Network Time

❖ Non-pipelined implementation

  ❖ 4 servers **x** 2 CPUs/sever **x** 4 cores/CPU **x** 2 threads/core

  ❖ 10 Gbit Ethernet projected from 1 Gbit (Columbia CLIC lab)

■ Network   ■ CPU

*slightly higher CPU time*

*much lower network time*

HJ   2TJ   3TJ   4TJ

# Part 2: Memory-bound (random)

# Why partitioning?

- Random accesses << sequential accesses
  - Cache misses
  - TLB misses

- **Where** to use partitioning
  - **Sorting**
  - Joins
  - Group-by aggregation
  - Materialization

# Variants of Partitioning

# Variants of Partitioning

# Previous Work: Partitioning small arrays

❖ Compute **histogram**

  ❖ Contiguous arrays

  ❖ Prefix sum of histogram

❖ Shuffle the data

  ❖ **Copy** from input to output

input          output

# Previous Work: Partitioning large arrays

- If size of array >> size of cache:

  - TLB thrashing

  - Cache conflicts

  - Cache pollution

- Use **buffering**

  - Store tuples in **cache-resident** buffers

  - **Write-combine** full buffers to output

- Parallel

  - **Interleave** histograms during prefix sum

input

buffers

output

# Partitioning large arrays in place

- ❖ Transfer data in cache lines
  - ❖ **Amortize** out-of-cache accesses
  - ❖ RAM <—> CPU cache

- ❖ "Work" on the cached buffers
  - ❖ Similar to in-cache ("swap cycles")
  - ❖ Data transferred across buffers

- ❖ Recycle buffers when done
  - ❖ Flush buffer when filled
  - ❖ Refill buffer with next data

input & output

buffers

# Partitioning large arrays

# Parallel partitioning in-place

❖ Swap tuples in-place

   ❖ Using atomics

   ❖ Extreme **synchronization** cost

❖ Swap blocks of tuples in-place

   ❖ Partition to **list of blocks** in-place

   ❖ Swap **blocks** of tuples

   ❖ **Amortize** synchronization cost

*partition in blocks*

*swap blocks*

# Partitioning Function

- ❖ Radix
  - ❖ Trivial

- ❖ Hash
  - ❖ Depends on hash function

- ❖ Range
  - ❖ **Slow** with binary search
  - ❖ **Fast** with range index



Legend:
- range (SIMD index)
- range (binary search)
- radix
- hash

*RAM **bandwidth** for hash and radix*

*5X by range index*

Y-axis: Billion tuples per second (0, 5, 10, 15, 20, 25, 30, 35)

X-axis: Partitioning fanout (128, 200, 256, 360, 512, 1000, 1024, 1800, 2048)

# Large-scale Sorting

- Stable LSB radix-sort
  - Parallel radix partitioning (not in-place)

- In-place MSB radix-sort
  - Parallel in-place radix partitioning
  - In-place radix partitioning

- Comparison-sort (CMP)
  - Parallel range partitioning (not in-place)
  - SIMD comb-sort in the cache



LSB    MSB    CMP

Billion tuples per second

0.8
0.6
0.4
0.2
0

1    2.5    5    10    25    50

Billion tuples
(32-bit key & payload)

# NUMA Awareness

- Optimize for NUMA
  - Use **local** RAM per CPU
  - **Minimize** NUMA transfers

- Transfers per sorting variant
  - LSB: up to 1 transfer
  - MSB: up to 2 transfers
  - CMP: up to 1 transfer

# Part 3: Memory-bound (sequential)

core pipeline

CPU cache

main memory

NUMA

network

Core Core
Core Core

RAM — CPU CPU — RAM

CPU CPU

RAM — RAM

Server  Server  Server  Server

Network

# Compression in Databases

❖ Why compress?

  ❖ Make dataset RAM resident

  ❖ Process data **faster** than RAM bandwidth

❖ **Dictionary** encoding

  ❖ Process **without** decompressing



*mapping*

*dictionary*

*original data*

*bit packing*

# Bit Packing Layouts

❖ Horizontal Bit Packing

**00010**000 **10100**000 **11000**000 **11111**000 **01010**000 **11001**000 **10001**000 **00100**000

$b = 5$

**00010101** **00110001** **11110101** **01100110** **00100100**

❖ Vertical Bit Packing

**00010**000 **10100**000 **11000**000 **11111**000 **01010**000 **11001**000 **10001**000 **00100**000 · · · ·

· · · · · · · ·

$b = 5$
$k = 8$

**01110110** **00111100** **01010001** **10011000** **00010110** · · · · ·

# Previous Work: Scan Horizontal

*LSB*  *MSB*

| 00010101 | 00110001 | 11110101 | 0110011**0** |
|----------|----------|----------|----------|

*load*   *8-bit —> 4-bit*

| 0001 | 0101 | 0011 | 0001 | 1111 | 0101 | 0110 | 01**10** |
|------|------|------|------|------|------|------|------|

*shuffle*

| 0001 | 0101 | 0101 | 0011 | 0011 | 0001 | 0001 | 1111 |
|------|------|------|------|------|------|------|------|

*4-bit —> 8-bit*

| 00010101 | 01010011 | 00110001 | 00011111 |
|----------|----------|----------|----------|

*shift*   << << << <<

| 00010101 | 10100110 | 11000100 | 11111000 |
|----------|----------|----------|----------|

*mask*   & & & &

| 00010000 | 10100000 | 11000000 | 11111000 |
|----------|----------|----------|----------|

❖ Fully packed

  ❖ No bits wasted

  ❖ Unpack **before** evaluating predicates

  ❖ Unpack in **SIMD**

# Previous Work: Scan Horizontal

select ... where column < C ...

- ❖ Word aligned

  - ❖ Scan **without** unpacking

  - ❖ Using scalar code

  - ❖ Bits **wasted**

  - ❖ Parallel bit extraction

**01**

*set constant C*

**01**0 **01**0 00

*invert code*

**01**0 **10**0 00 ^ 110 110 00 = **10**0 **01**0 00

*add constant C*

**10**0 **01**0 00 + **01**0 **01**0 00 = 110 001 00

*extract bits*

110 001 00 —> **01**

# Previous Work: Scan Horizontal



because we **unpack**

**wasted** space

— Fully packed (scalar)
— Fully packed (SIMD)
— Word aligned (scalar)

Scan thoughput (GB/s)

Number of bits

# Previous Work: Scan Vertical

select ... where column < C ...

❖ Scan **without** unpacking

    ❖ Bit-wise operations

    ❖ Both scalar & SIMD

    ❖ Can skip bits and **stop early**

*result |= constant & ~data*

| | | |
|---|---|---|
| 00010110 | 00000000 | ___0_00_ |
| 10011000 | 11111111 | _110_001 |
| 01010001 | 11111111 | 11101001 |
| 00111100 | 00000000 | *result* |
| 01110110 | 00000000 | |
| *data* | *constant* | |

# Previous Work: Scan Vertical



*vertical scan without unpacking*

**2X** *throughput for large k
(k: number of interleaved codes)*

Vertical (k = 64)
Horizontal (full)
Horizontal (word)
Vertical (k = 8192)

Scanning thoughput (GB/s)

Number of bits

# Pack Vertical Layout

- ❖ Scalar
  - ❖ **Extract** 1 bit per instruction
  - ❖ $O(n * b)$

- ❖ SIMD
  - ❖ Keep codes in SIMD registers
  - ❖ Maximize **bits per SIMD instruction**

| 000**00010** | 000**10100** | 000**11000** | 000**11111** |
|---|---|---|---|

| 000**01010** | 000**11001** | 000**10001** | 000**00100** |
|---|---|---|---|

*pack*

| **0010** | **0100** | **1000** | **1111** | **1010** | **1001** | **0001** | **0100** |
|---|---|---|---|---|---|---|---|

**00010110**

| **0001** | **0010** | **0100** | **0111** | **0101** | **0100** | **0000** | **0010** |
|---|---|---|---|---|---|---|---|

*extract & shift*

**10011000**

| 000**0** | 000**1** | 00**10** | 00**11** | 00**10** | 00**10** | 00**00** | 000**1** |
|---|---|---|---|---|---|---|---|

**01010001**

| 0000**0** | 0000**0** | 000**1** | 000**1** | 000**1** | 000**1** | 0000**0** | 0000**0** |
|---|---|---|---|---|---|---|---|

**00111100**

*shift & pack*

| 0000**0** | 0000**1** | 000**1** | 000**1** | 0000**0** | 000**1** | 000**1** | 0000**0** |
|---|---|---|---|---|---|---|---|

46

# Pack Vertical Layout

❖ Scalar

  ❖ **Extract** 1 bit per instruction

  ❖ O(n * b)

❖ SIMD

  ❖ Keep codes in SIMD registers

  ❖ Maximize **bits per SIMD instruction**

SIMD  Scalar

*up to* **27X** *improvement* **!**

Packing thoughput (GB/s)

Number of bits

# Unpack Vertical Layout

❖ Scalar

  ❖ **Insert** 1 bit per instruction

  ❖ O(n * b)

❖ SIMD

  ❖ Keep codes in SIMD registers

  ❖ Maximize **bits per SIMD instruction**



Legend: — SIMD   — Scalar

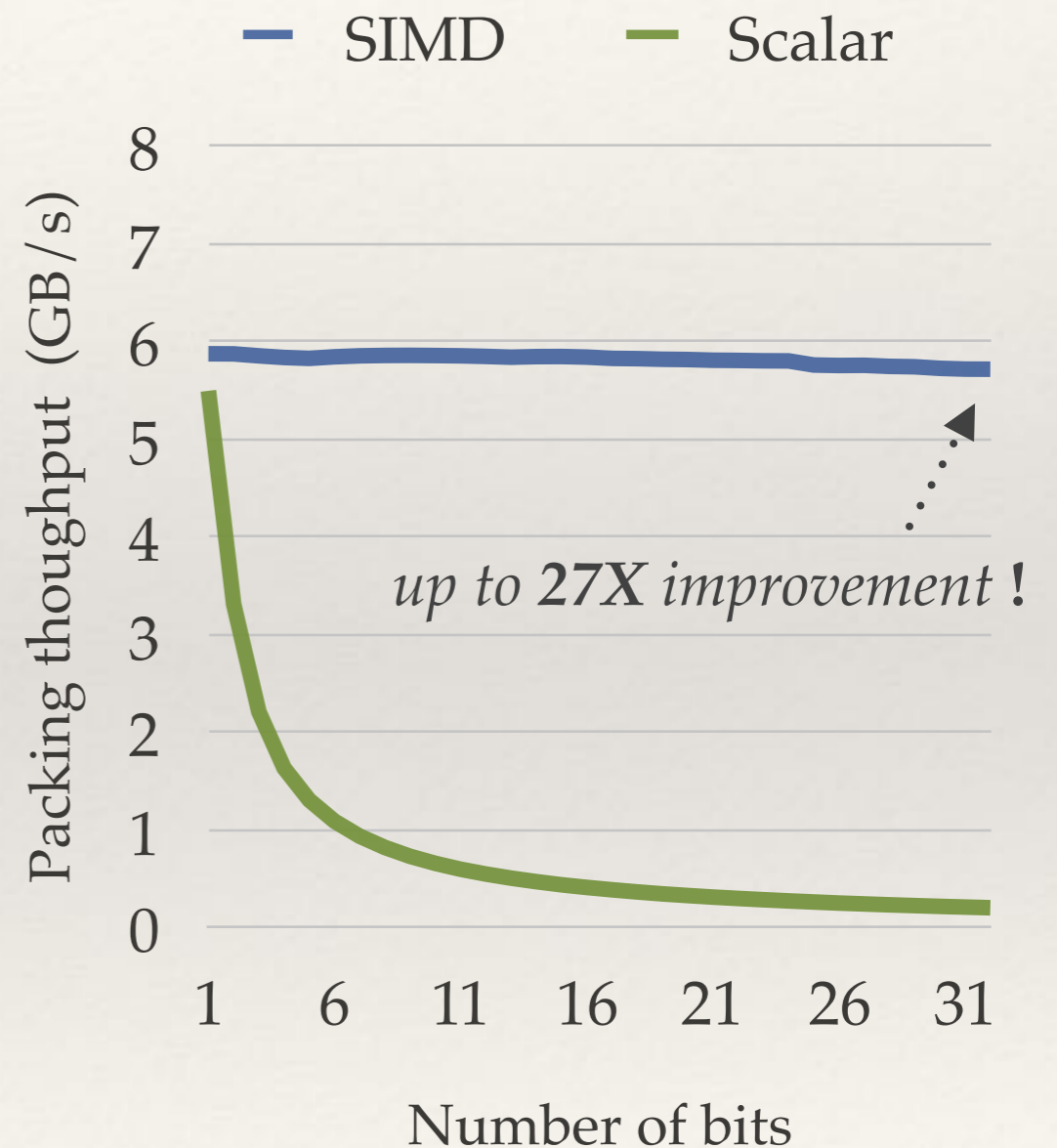Y-axis: Unpacking thoughput (GB/s) [logarithmic scale]: 0.1, 1, 10, 100

X-axis: Number of bits: 1, 6, 11, 16, 21, 26, 31

*11–20X improvement* **!**

# What if not memory-bound?

# Part 4: Compute-bound

core pipeline

CPU cache

main memory

NUMA

network

Core    Core
Core    Core

RAM    CPU    CPU    RAM

RAM    CPU    CPU    RAM

Server    Server    Server    Server

Network

# Many-Core (MIC) Platforms

- Mainstream CPUs

  - Aggressively **out-of-order**

  - Massively **super-scalar**

  - ~20 cores ($$$$)

- **Many-core** co-processors

  - 1st generation

  - In-order

  - Not super-scalar

  - 16 GB of fast RAM

  - ~60 cores

~ *60* cores

# Advanced SIMD Vectorization

- Baseline operator
  - O(f(n)) complexity in scalar code

- **Fully** vectorized
  - O(f(n) / **W**) complexity in SIMD code
  - Excluding **random** memory accesses

- **Reusable** vectorization techniques
  - Reuse fundamental operations

*8-way **SIMD** permutation*

| 1 | 1 | 2 | 3 | 5 | 8 | 13 | 21 |

| 5 | 4 | 6 | 1 | 7 | 3 | 2 | 0 |

| 8 | 5 | 13 | 1 | 21 | 3 | 2 | 1 |

# Fundamental Operations

(input) vector [ A | B | C | D ]

array [ X | Y | | | | ] ....

❖ Selective load

[ 0 1 1 0 ] mask

(output) vector [ A | X | Y | D ]

vector [ A | B | C | D ]

[ 0 1 1 0 ] mask

❖ Selective store

array [ B | C | | | | ] ....

# Fundamental Operations

❖ (Selective) gather

| index vector | 5 | 0 | 2 | 7 |
| --- | --- | --- | --- | --- |

| array | A | B | C | D | E | F | G | H |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |

| value vector | F | B | C | H |
| --- | --- | --- | --- | --- |

❖ (Selective) scatter

| value vector | A | B | C | D |
| --- | --- | --- | --- | --- |

| index vector | 5 | 0 | 2 | 7 |
| --- | --- | --- | --- | --- |

| array | B | | C | | | A | | D |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |

# Vectorized Selection Scans

select … where column > C …

❖ Scalar

   ❖ Branching

   ❖ Branchless

❖ Vectorized

   ❖ Using **selective stores**

   ❖ Store tuples (early)

   ❖ Store rids & dereference (late)

input payloads | v1 | v2 | v3 | v4 | v5 | v6 | v7 | ⋯⋯

input keys | k1 | k2 | k3 | k4 | k5 | k6 | k7 | ⋯⋯

*1) evaluate predicate(s)*

| k1 | k2 | k3 | k4 | ⋯ | C | C | C | C |

*2) selectively store qualifiers*

output keys | k0 | k2 | k3 | | ⋯⋯

output payloads | v0 | v2 | v3 | | ⋯⋯

# Vectorized Selection Scans



Throughput (billion tuples / second) vs Selectivity (%)

"late" is good for low selectivity

big speedup overall

- Scalar (branching)
- Scalar (branchless)
- Vectorized (early)
- Vectorized (late)

# Previous Work: Vectorized Hash Probing

❖ Scalar

   ❖ 1 input key at a time

   ❖ 1 table key per input key

❖ **Horizontal** vectorization

   ❖ 1 input keys at a time

   ❖ W table keys per input key

linear probing
**bucketized**
hash table

hash
index

$h$

input key

$k$

$k$ $k$ $k$ $k$

keys

payloads

# Vectorized Hash Probing

**Vertical** vectorization

- W input keys at a time
- 1 table keys per input key

input keys

hash indexes

| k1 | | h1 |
| k2 | | h2 |
| k3 | | h3 |
| k4 | | h4 |

*gather buckets*

linear probing hash table

keys — payloads

| k2 | |
| | |
| k5 | |
| | |
| k3 | |
| | |
| k7 | |

*matching **lanes***

*non-matching **lanes***

# Vectorized Hash Probing

**Vertical** vectorization

- W input keys at a time

- 1 table keys per input key

input keys

| k1 |
| k5 |
| k6 |
| k4 |

linear probing offsets

| 0 |
| 1 |
| 1 |
| 0 |

*selectively load* *input keys*

hash indexes

| h1+1 |
| h5 |
| h6 |
| h4+1 |

linear probing hash table

| k2 | |
| | |
| k5 | |
| | |
| k3 | |
| | |
| k7 | |

*matching **lanes** replaced*

*non-matching **lanes** kept*

59

# Vectorized Hash Table Probing



*out of the cache*

*much faster in the cache* !

Probing throughput (billion tuples / second)

- Scalar
- Vectorized (horizontal)
- Vectorized (vertical)

# Vectorized Data Shuffling

- ❖ Scalar
  - ❖ Move 1 tuple at a time

- ❖ Vectorized
  - ❖ **Scatter** tuples to output
  - ❖ **Serialize** conflicts

*1) **gather** offsets*

keys

| k1 | k2 | k3 | k4 |
|----|----|----|----|

partition ids

| h1 | h2 | h3 | h4 |
|----|----|----|----|

offset array

|  |  | l1 |  | l3 |  | l2 |  |  |  |
|--|--|----|--|----|--|----|--|--|--|

*conflicting **lanes***    *non-conflicting **lanes***

partition offsets

| l1 | l2 | l3 | l4 |
|----|----|----|----|

*2) **scatter** tuples*

output keys

| k1 |  |  |  | | k4 |  |  |  | | k3 |  |  |  |

output payloads

| v1 |  |  |  | | v4 |  |  |  | | v3 |  |  |  |

# Vectorized Data Shuffling

**Scalar**

- Move 1 tuple at a time

**Vectorized**

- **Scatter** tuples to output
- **Serialize** conflicts

*1) **gather** offsets*

keys

| k1 | k2 | k3 | k4 |

partition ids

| h1 | h2 | h3 | h4 |

offset array

| | | l1 | | | l3 | | | l2 | | | |

*conflicting **lanes***   *non-conflicting **lanes***

serialization offsets

| 0 | 0 | 0 | **1** |

+

partition offsets

| l1 | l2 | l3 | l4 |

*2) **serialize** conflicts*

*3) **scatter** tuples*

output keys

... | k1 | | | | ... | k2 | k4 | | | ... | k3 | | | | ...

... | v1 | | | | ... | v2 | v4 | | | ... | v3 | | | | ...
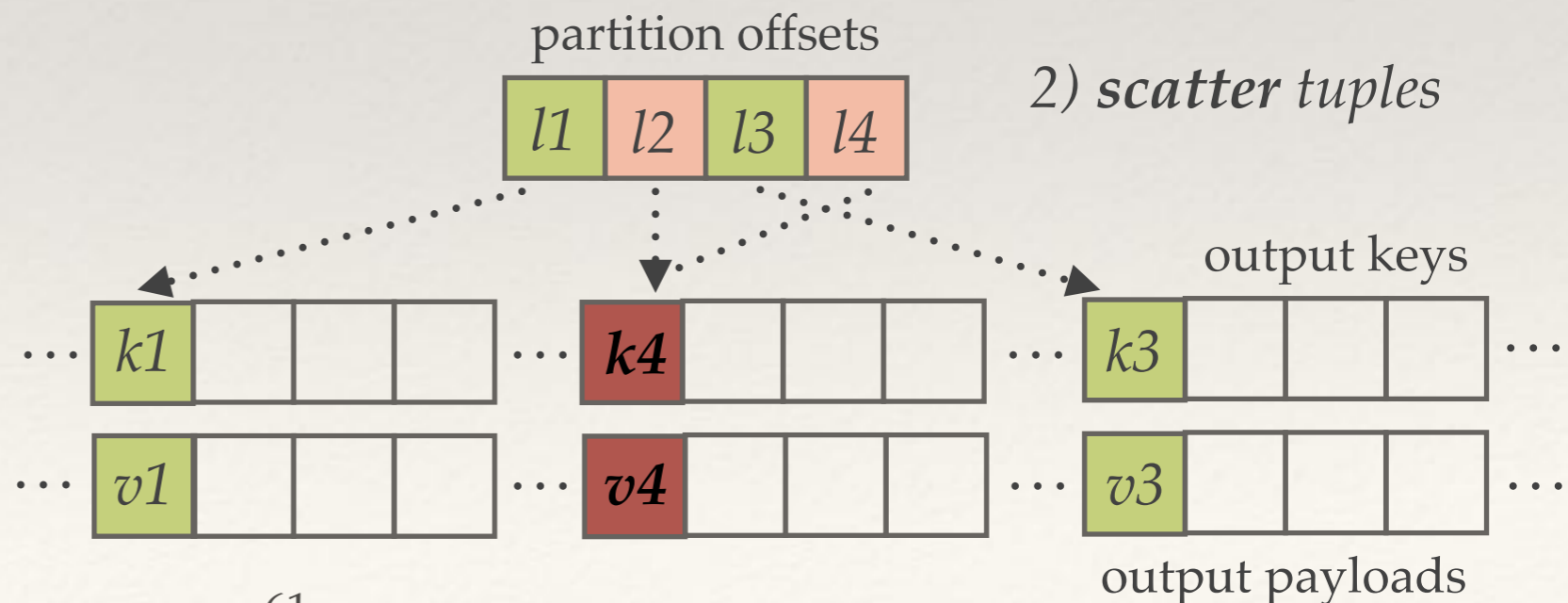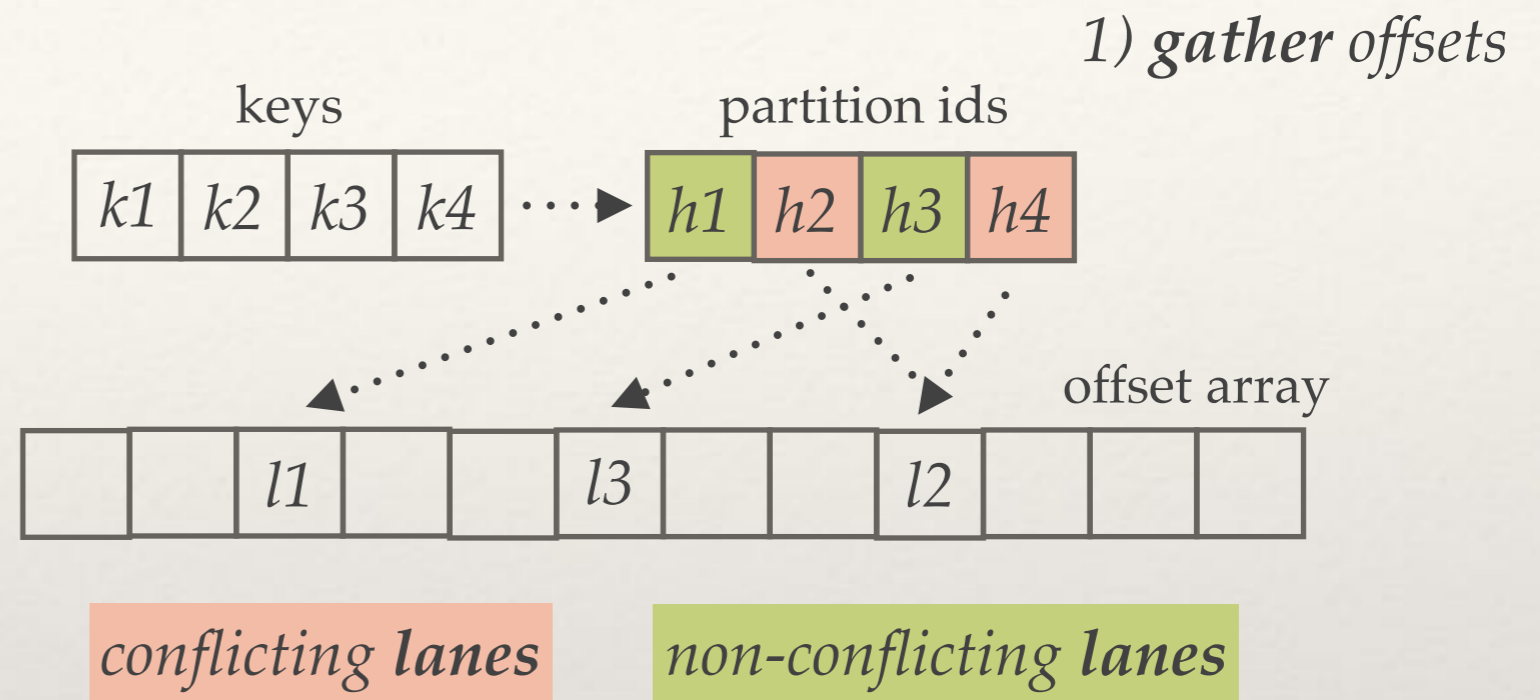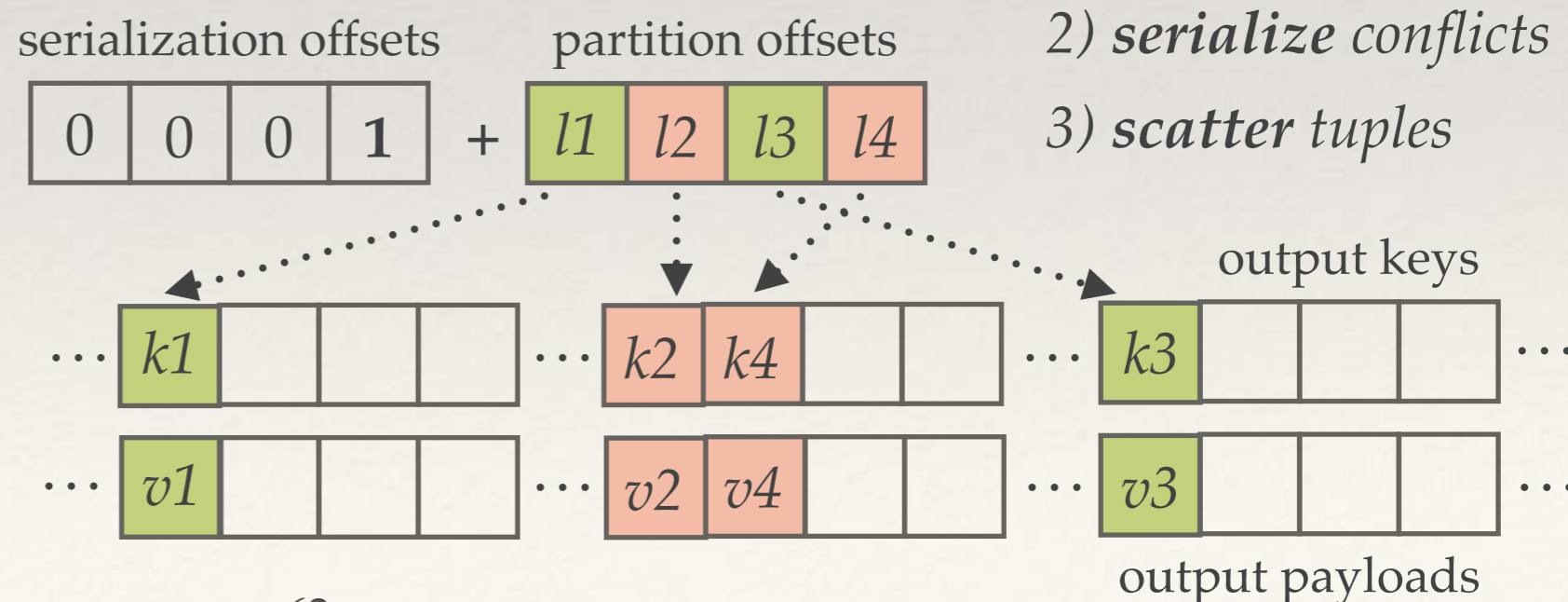
output payloads

62

# Vectorized Buffered Data Shuffling

- ❖ Scalar
  - ❖ Move 1 tuple at a time

- ❖ Vectorized
  - ❖ **Scatter** tuples to **buffers**
  - ❖ **Serialize** conflicts

*vertical* logic

input

*1-tuple scatters*

buffers

*W-tuple stream stores*

output

*horizontal* logic

# Vectorized Partitioning



Shuffling throughput (billion tuples / second) vs Fanout (log)

*vectorization* ***orthogonal*** *to buffering*

- Scalar unbuffered
- Scalar buffered
- Vector unbuffered
- Vector buffered

# Vectorized Operators

❖ Selection scans

❖ Partitioning

 ❖ Histogram

 ❖ Data shuffling

❖ Hash table building & probing

 ❖ Linear probing

 ❖ Double hashing

 ❖ Cuckoo hashing

❖ Bloom filter probing

❖ Regular expression matching

❖ Sorting

 ❖ LSB radix-sort

❖ Hash joins

 ❖ Non-partitioned

 ❖ Partitioned

# Hash Join Performance

❖ Hash **join** 200 with 200 million tuples  (2X 32-bit key & payload)



❖ Being **cache-conscious** matters !

# Part 5: An Engine for Many-Cores

core pipeline

CPU cache

main memory

NUMA

network

Core  Core
Core  Core

RAM — CPU   CPU — RAM

CPU   CPU

RAM — CPU   CPU — RAM

Server   Server   Server   Server

Network

# Why many-core CPUs?

- More **complex** cores
  - Super-scalar out-of-order cores
  - Core size: 1st-gen << 2nd-gen << mainstream

- Additional layer of on-chip *MCDRAM*
  - **~4X** higher **bandwidth** than DDR4 DRAM
  - **Larger** than the caches (16 GB)

- Advanced **SIMD**: AVX-512
  - **Same** as upcoming mainstream CPUs

core pipeline

CPU cache

on-chip memory

off-chip memory

NUMA

network

# Baseline: Code Generation

- Code generation

  - Generate code **per query** at runtime

  - **Pipelined** operators

  - **Specialized** data structures

```
select sum(F.val * A.val * B.val)
from F, A, B
where F.key_A = A.key
   and F.key_B = B.key
   and F.val between x0 and y0
   and A.val between x1 and y1
   and B.val between x2 and y2;
```
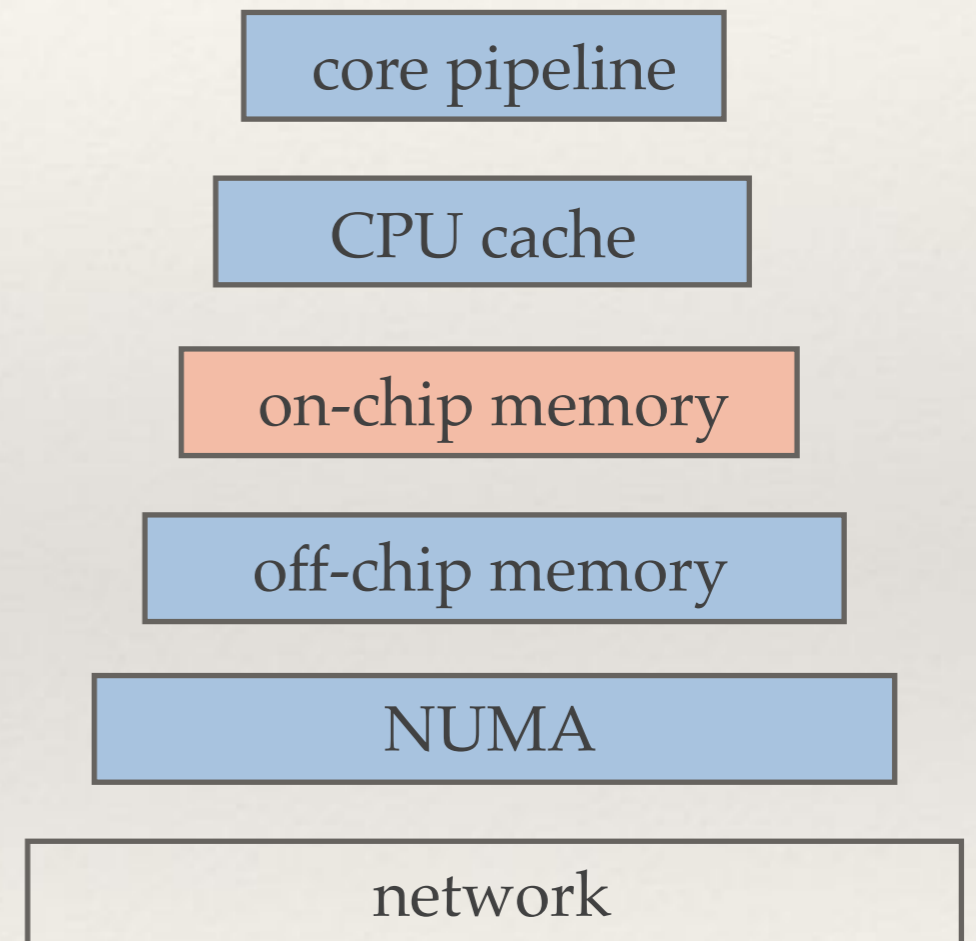
```
typedef struct {
    int set:1;
    A_key_t key;
    A_val_t val;
} A_key_val_t;

for (size_t i = 0; i != F_tuples; ++i) {
    if (F_val[i] >= x0 && F_val[i] <= y0) {

        size_t h1 = hash(F_key_A[i], buckets_HJT_a);
        while (HJT_a[h1].set) {
            if (HTJ_a[h1].key == F_key_A[i]) {

                size_t h2 = hash(F_key[i], buckets_HJT_b);
                while (HJT_b[h2].set) {
                    if (HJT_b[h2].key == F_key[i]) {

                        sum += F_val[i] * HJT_a[h1].val
                              * HJT_b[h2].val; }

                    if (++h2 == buckets_HJT_b) h2 = 0; }}

            if (++h1 == buckets_HJT_a) h1 = 0; }}}
```

# Baseline + SIMD Vectorization

- ❖ **Maximize data parallelism**
  - ❖ Written entirely in SIMD
  - ❖ No **register-resident** execution
  - ❖ Move data in cache-resident **buffers**

- ❖ **SIMD can hurt performance**
  - ❖ Due to **cache** & **TLB misses**
  - ❖ Fast RAM does **not** help

*mostly hash probe*
*SIMD **<20%** faster*

*mostly selection scan*
*SIMD **2.5X** faster*

Time (in seconds)

Legend:
- Scalar (blue)
- Vector (green)

X-axis:
- HBW   LBW
- 10% selectivity
- HBW   LBW
- 90% selectivity

# VIP Engine

❖ Based on "sub-operators" that …

  ❖ Process a **block** of tuples at a time

  ❖ Process one **column** at a time within that block

  ❖ Designed to be **data-parallel**

  ❖ Implemented **entirely** in SIMD

❖ **Why** is the design fast ?

  ❖ **Specialized** sub-operators can be extremely **optimized**

  ❖ Block at a time execution reduces **materialization** & **interpretation** cost

  ❖ Use **cache-conscious** execution to utilize both **SIMD** and **fast RAM**

# Sub-operators: An Example

- Hash a composite key <A, B> (of types X, Y)

  - Hash one **block** at a time

  - Hash one **column** at a time per block

  - Call **hash_X()** on column **A** of type **X** for a block of tuples

  - Call **hash_Y()** on column **B** of type **Y** for a block of tuples

  - Keep working set (block of hash values) **cache-resident**

  - **Amortize** interpretation cost

  *32-bit integer prototype*

    **void hash_int32**(**const int32_t**\* data, **uint32_t**\* hash, **size_t** tuples);

# Selection Scans in VIP

- ❖ Based on sub-operators
    - ❖ Combine results using **bitmaps**
    - ❖ **Skip** tuples already determined
    - ❖ Process **W** items in SIMD

- ❖ Built-in **compression**
    - ❖ **Horizontal** dictionary compression
    - ❖ **Skip** tuples if determined
    - ❖ Decompress in 5 SIMD instructions

**select** * **from** T
**where** x = 9
**and** y > 1;

x = [1, 9, 8, 9]

y = [_, 0, _, 7]

*ignore or **skip***

**select** * **from** T
**where** x = 9
**or** y > 1;

x = [1, 9, 8, 9]

y = [2, _, 1, _]

# Selection Scans in VIP

❖ From TPC-H Q19  (SF = 1000)

  ❖ Selection on **part** table

  ❖ Neither CNF nor DNF

  ❖ **0.24**% selectivity

  ❖ **Skip** is essential here



*baseline is compute-bound*

*VIP is not*

Throughput (in billions of tuples per second)

20

15

10

5

0

HBW    LBW
uncompressed

HBW    LBW
compressed

Baseline    VIP

# Hash Joins in VIP

- ❖ Partition

  - ❖ **Inner** table must fit in the cache

- ❖ Hash join using **hash values**

  - ❖ **Specialized** data types & code

  - ❖ Generate rids lists

- ❖ Evaluate predicates

  - ❖ Use rid lists to access columns

  - ❖ Also evaluate non-equality predicates

  - ❖ Resolve hash **conflicts**

```
typedef struct {
    uint32_t hash;
    int32_t rid;
} join_bucket_t;

void build_hashes(
    const uint32_t* hashes,
    join_bucket_t* hash_table, [...]);

void probe_hashes(
    const uint32_t* hashes,
    const join_bucket_t* hash_table,
    int32_t* inner_rids,
    int32_t* outer_rids, [...]);
```
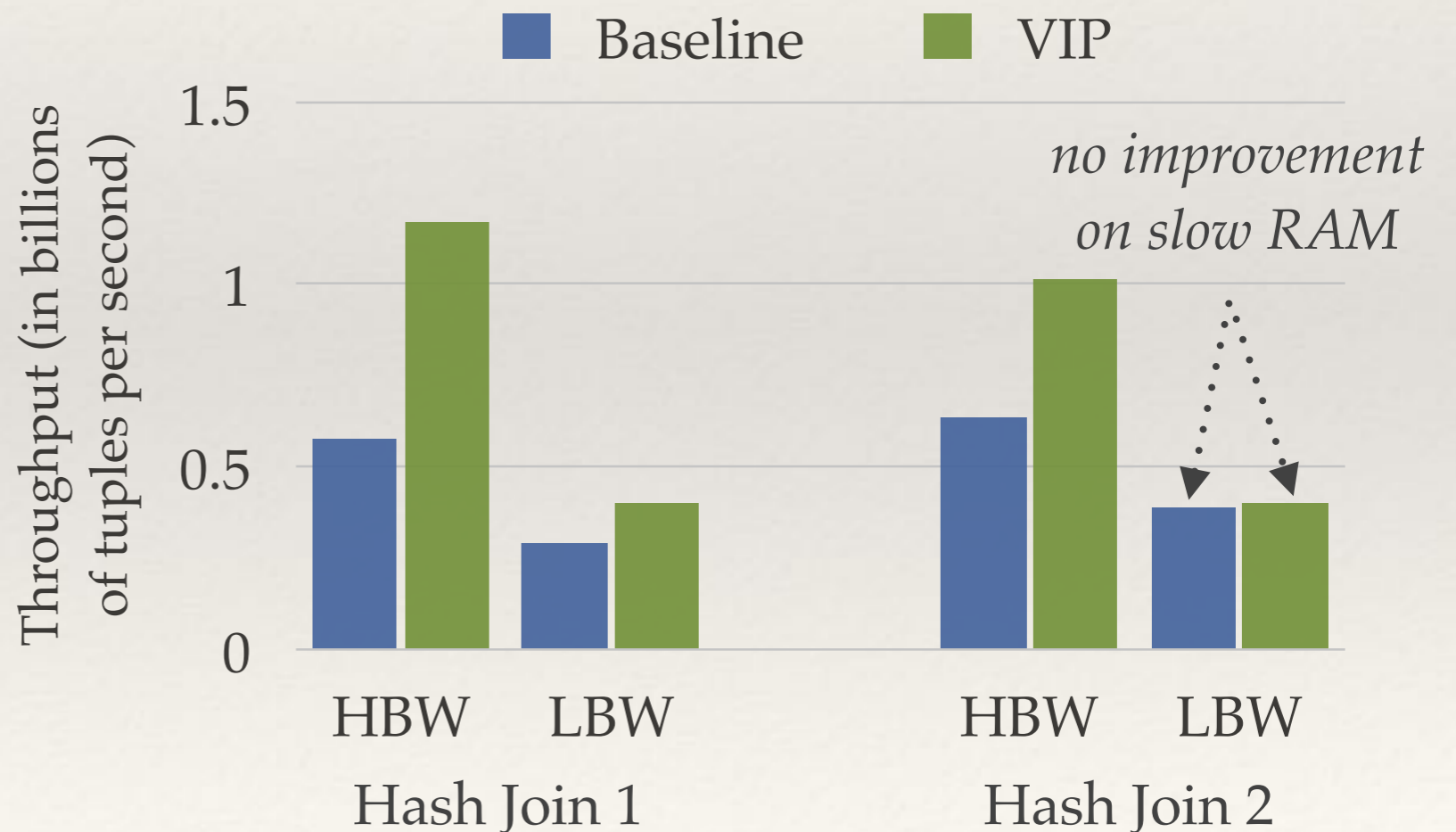
# Hash Joins in VIP

**select** l_partkey, l_suppkey, o_custkey
**from** lineitem, orders
**where** l_orderkey = o_orderkey;

**select** l_orderkey, l_partkey, l_suppkey
**from** lineitem, partsupp
**where** l_partkey = ps_partkey
   **and** l_suppkey = ps_suppkey;
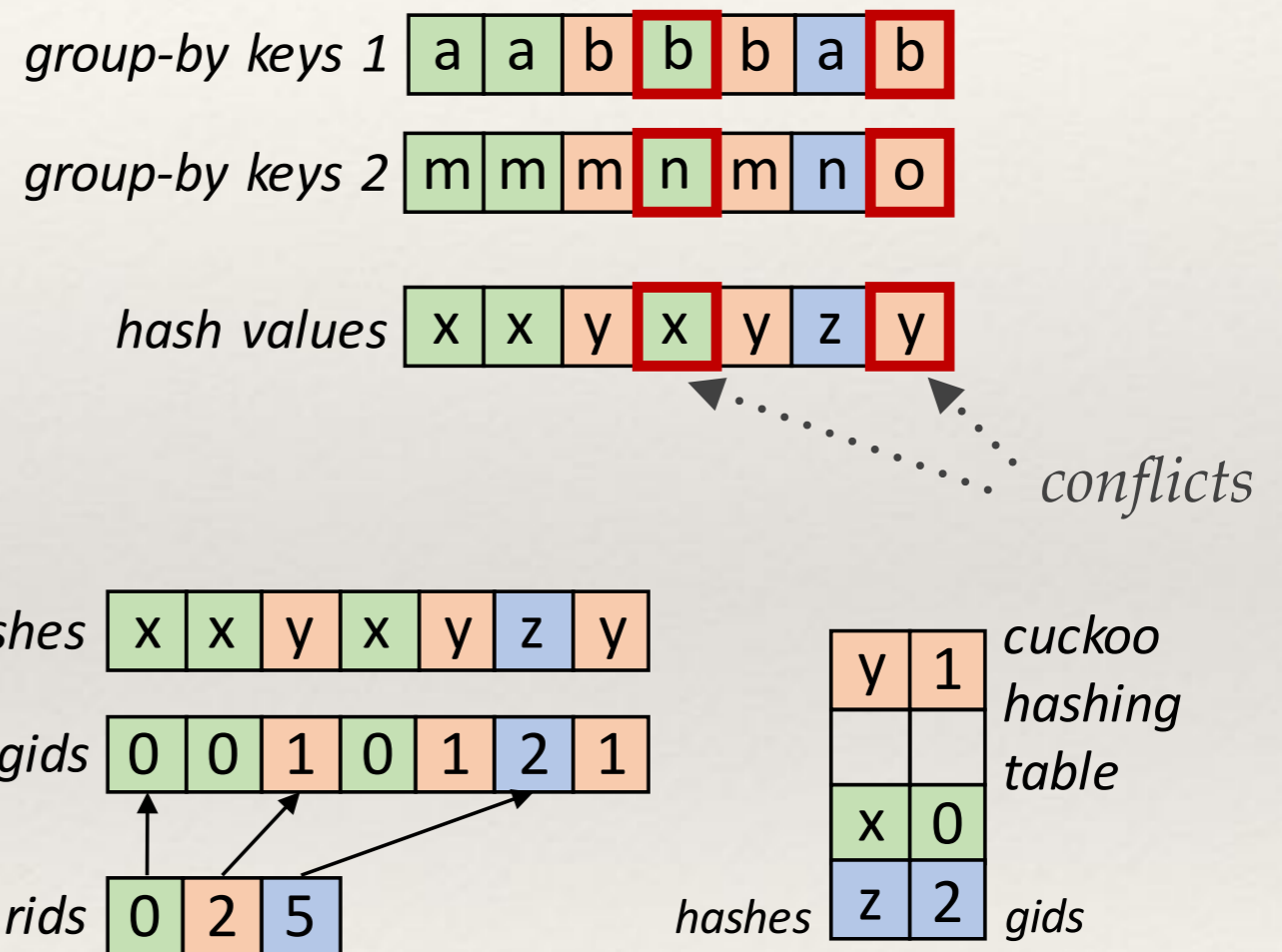
❖ From TPC-H  (SF = 30)

  ❖ **Largest** base tables

  ❖ **Core** joins of TPC-H

*no improvement on slow RAM*



*Throughput (in billions of tuples per second)*

HBW    LBW

Hash Join 1

HBW    LBW

Hash Join 2

Baseline    VIP

# Group-by Aggregation in VIP

- ❖ Partition
  - ❖ Estimate number of groups
  - ❖ **Output** groups must fit in the cache

*group-by keys 1* | a | a | b | b | b | a | b |

*group-by keys 2* | m | m | m | n | m | n | o |

- ❖ Map hashes to group-ids
  - ❖ Using **specialized** sub-operator
  - ❖ Fix hash **conflicts**

*hash values* | x | x | y | x | y | z | y |

*conflicts*

- ❖ Compute **expressions**
  - ❖ Store in **cache-resident** buffers

*hashes* | x | x | y | x | y | z | y |

*gids* | 0 | 0 | 1 | 0 | 1 | 2 | 1 |

*rids* | 0 | 2 | 5 |

*cuckoo hashing table*

| y | 1 |
| | |
| x | 0 |
| z | 2 |

*hashes*    *gids*

- ❖ Update aggregates via group-ids
  - ❖ Keep partial aggregates **in the cache**

77

# Group-by Aggregation in VIP

❖ From TPC-H Q1  (SF = 100)

❖ **2** group-by attributes

❖ **4** payload attributes

❖ **8** aggregate functions

❖ Reuse buffers for **sub-expressions**

**sum**(l_extendedprice),
**sum**(l_extendedprice * (1 - l_discount),
**sum**(l_extendedprice * (1 - l_discount) * (1 + l_tax)

*same speedup on HBW/LBW* **!**

Throughput (in billions of tuples per second)

| | Baseline | VIP |
|---|---|---|

# Future Work

❖ Track Join

  ❖ Overlap CPU & network computation to reduce **end-to-end** time

  ❖ Combine with **scheduling** algorithms for network transfers

❖ Compression

  ❖ **Multiple** dictionaries or more complex schemes (e.g. Huffman encoding)

  ❖ **Dynamic** dictionary encoding (e.g. add & update dictionary values)

❖ Vectorization

  ❖ Evaluate new hardware **platforms** with better SIMD (e.g. AVX-512)

  ❖ Design better **hardware** for database (e.g. better SIMD instructions)

❖ VIP engine

  ❖ **Pipeline** operators when cache misses cannot occur

  ❖ Evaluate **materialization** strategies & build operators in VIP

# Published Papers

- SIMD-Accelerated Regular Expression Matching

    - At DaMoN '16 with Eva Sitaridi, Kenneth A. Ross

- Rethinking SIMD Vectorization for In-Memory Databases

    - At SIGMOD '15 with Arun Raghavan, Kenneth A. Ross

- Efficient Lightweight Compression Alongside Fast Scans

    - At DaMoN '15 with Kenneth A. Ross

- Energy Analysis of Hardware and Software Range Partitioning

    - At TOCS with Lisa Wu, Raymond J. Barker, Martha A. Kim, Kenneth A. Ross

- A Comprehensive Study of Main-Memory Partitioning and its Application to Large-Scalar Comparison- and Radix-Sort

    - At SIGMOD '14 with Kenneth A. Ross

- Track Join: Distributed Joins with Minimal Network Traffic

    - At SIGMOD '14 with Rajkumar Sen, Kenneth A. Ross

- Vectorized Bloom Filters for Advanced SIMD Processors

    - At DaMoN '14 with Kenneth A. Ross

- High Throughput Heavy Hitter Aggregation for Modern SIMD Processors

    - At DaMoN '14 with Kenneth A. Ross

# Acknowledgments

❖ My advisor Ken

❖ My PhD thesis committee

   ❖ Martha, Luis, Eugene, and Stratos

❖ Friends and colleagues from the DB group

   ❖ Fotis, John, Eva, Pablo, and Wangda

❖ Colleagues from Oracle and Amazon

   ❖ Arun, Eric, Ippokratis, Michalis, and Raj

❖ More friends from Columbia & New York

# Thank you!