# Efficient Lightweight Compression Alongside Fast Scans
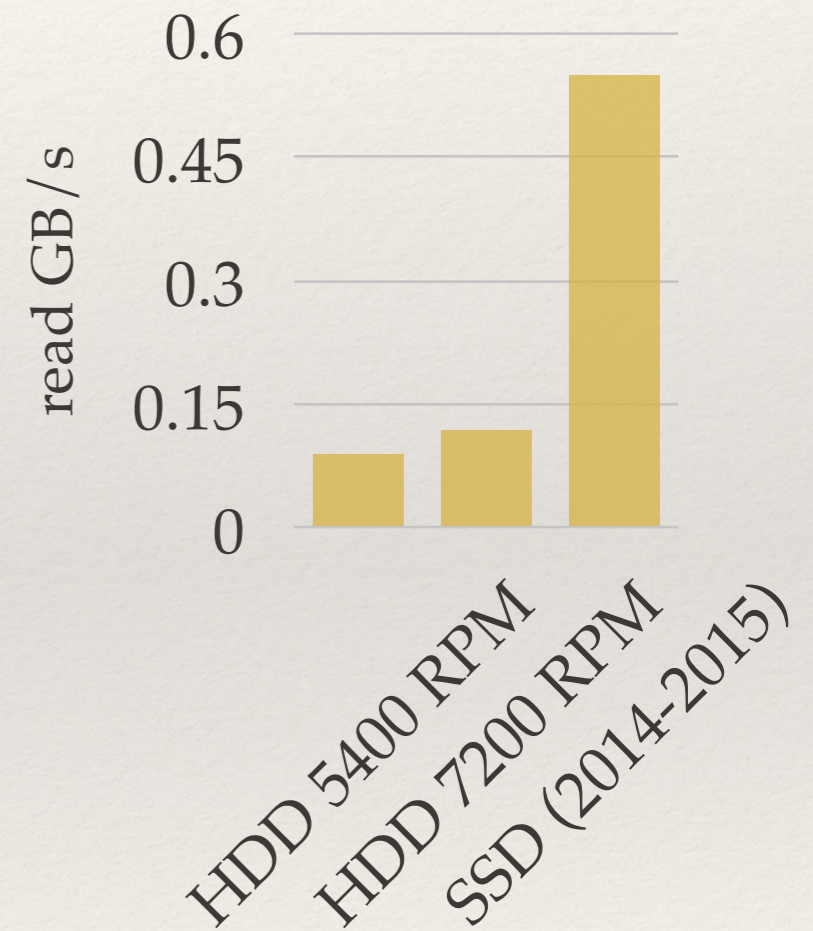
**Orestis Polychroniou**

**Kenneth A. Ross**

Columbia University
IN THE CITY OF NEW YORK

*DaMoN 2015, Melbourne, Victoria, Australia*

# Databases & Compression

❖ Process data on *disk*

 ❖ Nearly *unlimited* capacity

 ❖ Affects query *optimization*

  ❖ Minimize # of blocks fetched

  ❖ Minimize # of random block accesses

 ❖ Compress to improve disk speed

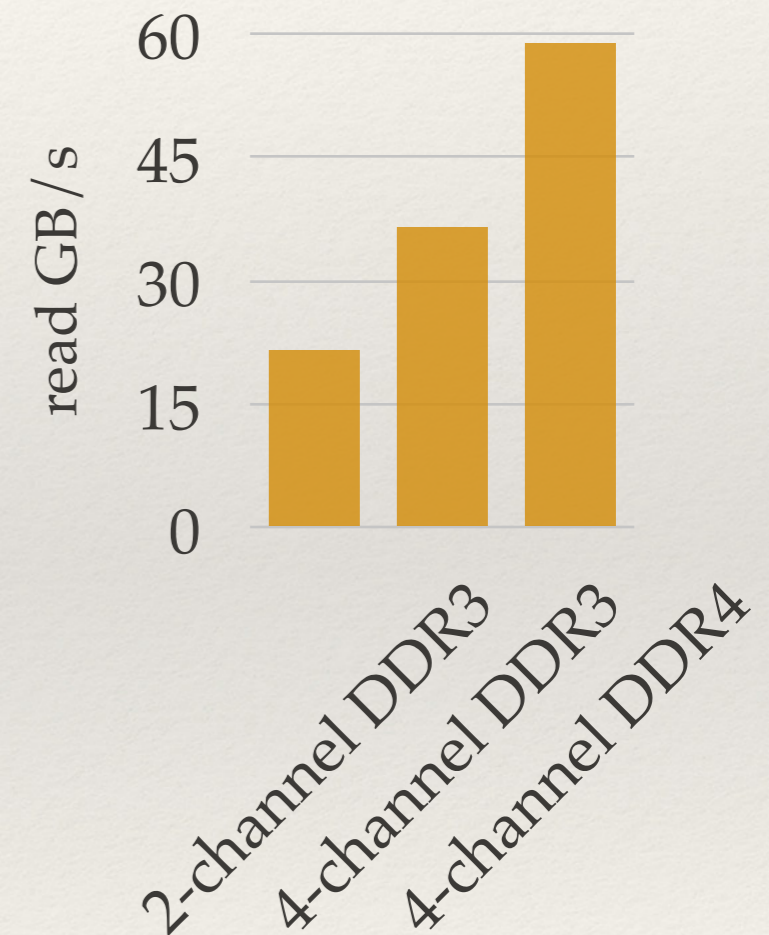  ❖ Focused on *compression rate* since disks are "slow"
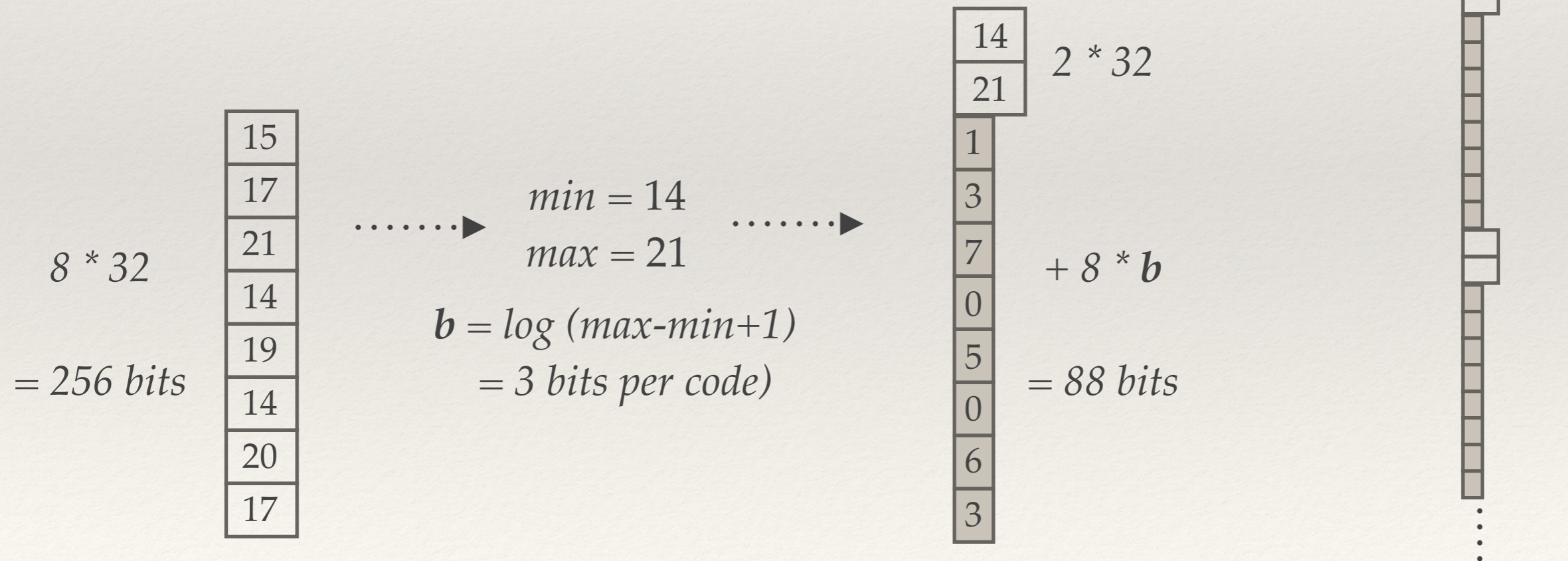
# Databases & Compression

- ❖ Process data on *disk*
    - ❖ Nearly *unlimited* capacity
    - ❖ Affects query *optimization*
        - ❖ Minimize # of blocks fetched
        - ❖ Minimize # of random block accesses
    - ❖ Compress to improve disk speed
        - ❖ Focused on *compression rate* since disks are "slow"

- ❖ Process data on *RAM*
    - ❖ Always *limited* capacity
    - ❖ Affects query *optimization* & query *execution*
        - ❖ Minimize # of accesses  (e.g. column stores & late materialization)
        - ❖ Minimize # of *random* (out of CPU cache) accesses   (e.g. partitioned join)
    - ❖ Compress to improve RAM speed & *avoid* disk
        - ❖ Focused on *(de-)compression efficiency* as RAM is "fast"

Chart: read GB/s vs memory type

| Memory type | read GB/s (approx) |
| --- | --- |
| 2-channel DDR3 | 22 |
| 4-channel DDR3 | 37 |
| 4-channel DDR4 | 59 |

y-axis: read GB/s (0, 15, 30, 45, 60)

# Lightweight Compression

- ❖ Compression schemes
  - ❖ Entropy compression
    - ❖ Group *nearby* similar values
    - ❖ e.g. run-length-encoding, frame-of-reference

*8 * 32*

*= 256 bits*

| 15 |
|----|
| 17 |
| 21 |
| 14 |
| 19 |
| 14 |
| 20 |
| 17 |

$min = 14$
$max = 21$

$b = log\ (max\text{-}min+1)$
$= 3\ bits\ per\ code)$

| 14 |
|----|
| 21 |
| 1 |
| 3 |
| 7 |
| 0 |
| 5 |
| 0 |
| 6 |
| 3 |

*2 * 32*
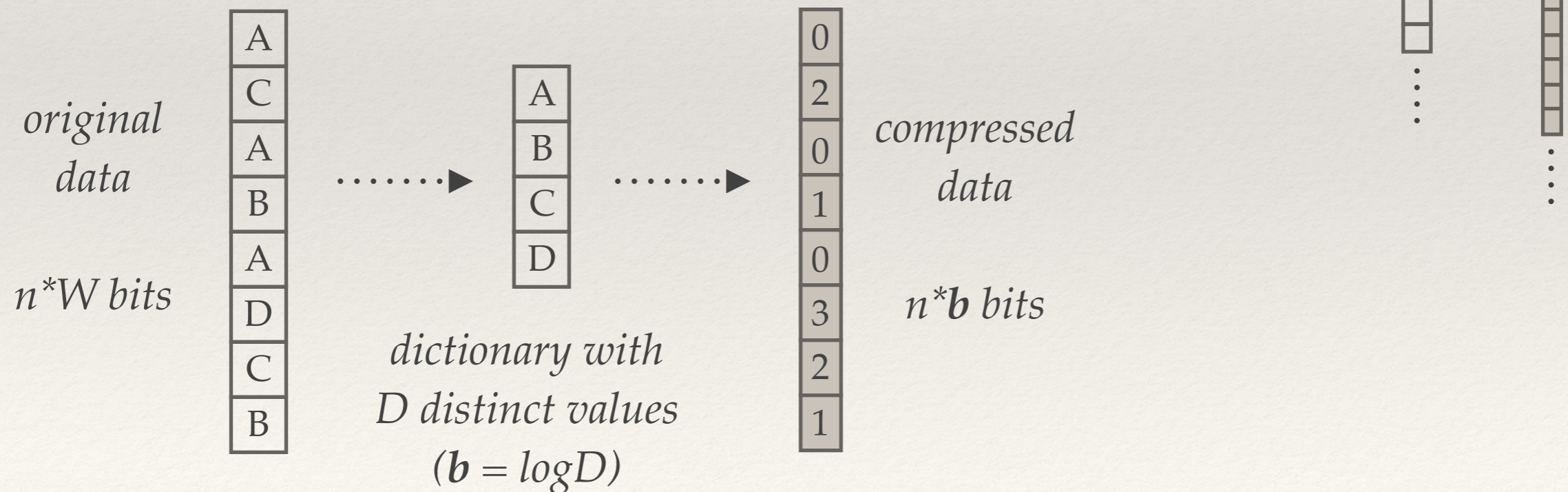
*+ 8 * b*

*= 88 bits*

# Lightweight Compression

- ❖ Compression schemes
  - ❖ Entropy compression
    - ❖ Group *nearby* similar values
    - ❖ e.g. run-length-encoding, frame-of-reference
  - ❖ Symbol compression
    - ❖ Assign a *symbol* to each distinct value
    - ❖ e.g. dictionary compression

*original data*

*n\*W bits*

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

┈┈┈➤

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

*dictionary with D distinct values*
*(b = logD)*

┈┈┈➤

| 0 |
|---|
| 2 |
| 0 |
| 1 |
| 0 |
| 3 |
| 2 |
| 1 |

*compressed data*

*n\*b bits*

+

# Lightweight Compression

- Compression schemes
  - Entropy compression
    - Group *nearby* similar values
    - e.g.  run-length-encoding, frame-of-reference
  - Symbol compression
    - Assign a *symbol* to each distinct value
    - e.g.  dictionary compression
  - Frequency (symbol) compression
    - Compress *frequent* symbols with less bits
    - e.g.  Huffman coding (slow), multiple dictionaries (fast)
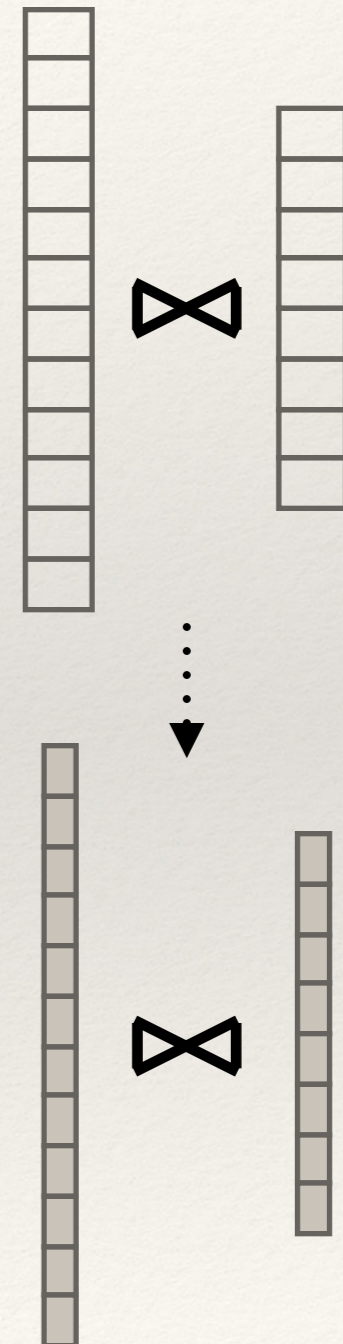
# Lightweight Compression

- ❖ Compression schemes
  - ❖ Entropy compression
    - ❖ Group *nearby* similar values
    - ❖ e.g. run-length-encoding, frame-of-reference
  - ❖ Symbol compression
    - ❖ Assign a *symbol* to each distinct value
    - ❖ e.g. dictionary compression
  - ❖ Frequency (symbol) compression
    - ❖ Compress *frequent* symbols with less bits
    - ❖ e.g. Huffman coding (slow), multiple dictionaries (fast)

- ❖ DBMS integration
  - ❖ Decompress *during* execution
    - ❖ In CPU cache (non-integrated) or in registers (integrated)

# Lightweight Compression

❖ Compression schemes

   ❖ Entropy compression

      ❖ Group *nearby* similar values

      ❖ e.g. run-length-encoding, frame-of-reference

   ❖ Symbol compression

      ❖ Assign a *symbol* to each distinct value

      ❖ e.g. dictionary compression

   ❖ Frequency (symbol) compression

      ❖ Compress *frequent* symbols with less bits

      ❖ e.g. Huffman coding (slow), multiple dictionaries (fast)

❖ DBMS integration

   ❖ Decompress *during* execution

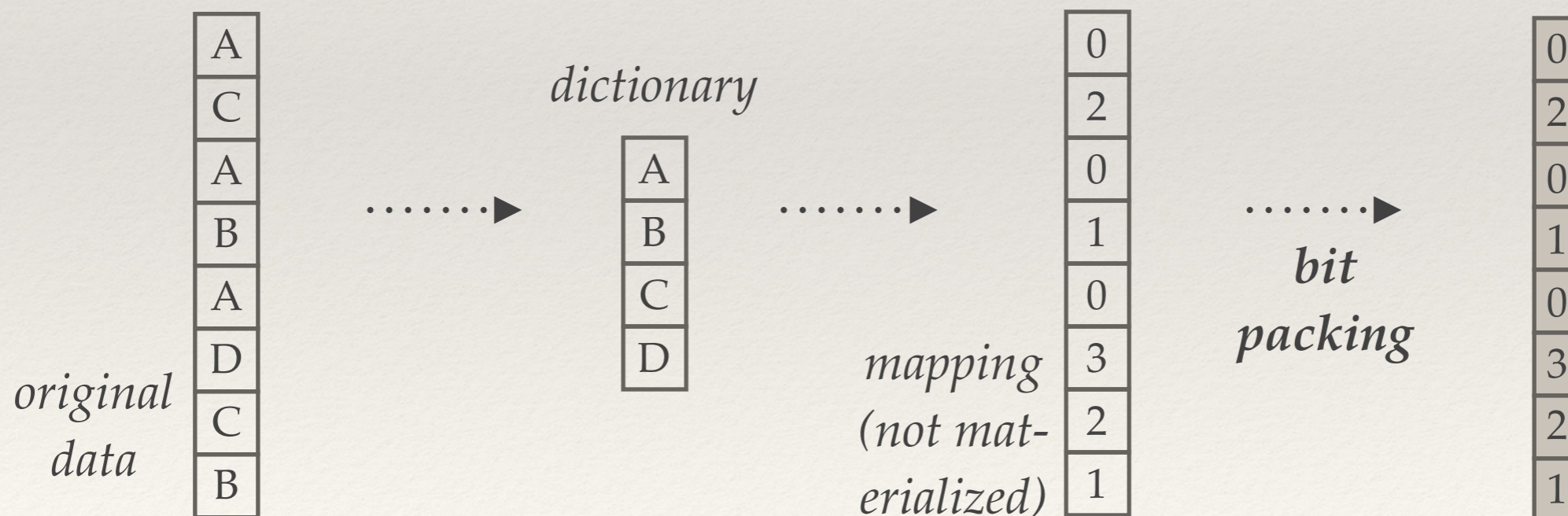      ❖ In CPU cache (non-integrated) or in registers (integrated)

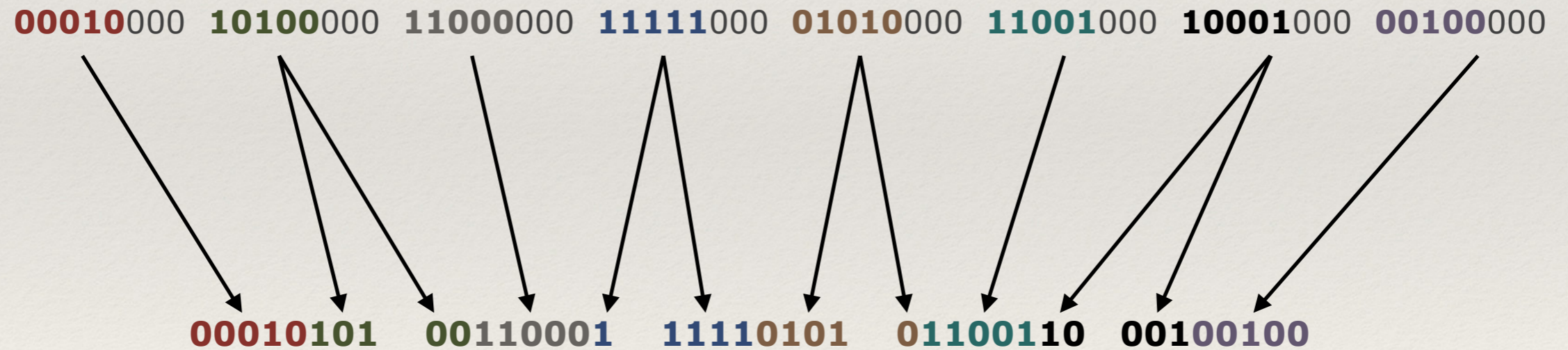   ❖ Process *compressed* data without decompressing

# Bit Packing

- ❖ Definition
  - ❖ Input code width is hardware-supported
    - ❖ 8-bit, 16-bit, 32-bit, 64-bit
  - ❖ Output code width $b$ must be (almost) constant
    - ❖ Either constant across the *entire* input
    - ❖ Or constant for the next *group* of items  (e.g. frame-of-reference)

| original data | | dictionary | | mapping (not mat-erialized) | | bit packing |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| A | | | | 0 | | 0 |
| C | | | | 2 | | 2 |
| A | | A | | 0 | | 0 |
| B | ·····▶ | B | ·····▶ | 1 | ·····▶ | 1 |
| A | | C | | 0 | | 0 |
| D | | D | | 3 | | 3 |
| C | | | | 2 | | 2 |
| B | | | | 1 | | 1 |

# Bit Packing

- Layouts
  - *Horizontal* bit packing
    - Bits per code are *contiguous*

**00010**000 **10100**000 **11000000** **11111**000 **01010**000 **11001**000 **10001**000 **00100**000

**00010101** **00110001** **11110101** **0110011**0 **001**00100

# Bit Packing

- Layouts
  - *Horizontal* bit packing
    - Bits per code are *contiguous*
  - *Vertical* bit packing
    - Bits of codes are *interleaved*

# Bit Packing

* Layouts

  * *Horizontal* bit packing

    * Bits per code are *contiguous*

  * *Vertical* bit packing

    * Bits of codes are *interleaved*



$b = 5$
$k = 4$

$b = 5$
$k = 8$

# Outline

- ❖ Operations
  - ❖ Packing
  - ❖ Unpacking
  - ❖ Scanning

# Outline

- Operations
  - Packing
  - Unpacking
  - Scanning

- Horizontal layouts
  - Fully packed
    - Fast unpacking & scanning
  - Word aligned
    - Faster scanning

# Outline

- ❖ Operations
  - ❖ Packing
  - ❖ Unpacking
  - ❖ Scanning

- ❖ Horizontal layouts
  - ❖ Fully packed
    - ❖ Fast unpacking & scanning
  - ❖ Word aligned
    - ❖ Faster scanning

- ❖ Vertical layout
  - ❖ Known traits
    - ❖ Fastest scanning
  - ❖ **New traits**
    - ❖ **Fast packing & unpacking**

# Horizontal Layout

- ❖ *Fully* packed
    - ❖ No space wasted
        - ❖ Codes can *span* across 2 packed words

# Horizontal Layout

❖ ***Fully*** packed

    ❖ No space wasted

        ❖ Codes can *span* across 2 packed words

    ❖ Packing

        ❖ Process 1 unpacked *code* per iteration

        ❖ Branch to store output packed *word*

    ❖ Unpacking

        ❖ Process 1 output *code* per iteration

        ❖ Branch to load input packed *word*

# Horizontal Layout

- *Fully* packed
  - No space wasted
    - Codes can *span* across 2 packed words
  - Packing
    - Process 1 unpacked *code* per iteration
    - Branch to store output packed *word*
  - Unpacking
    - Process 1 output *code* per iteration
    - Branch to load input packed *word*
    - Can be written in *SIMD* !

Based on paper by
T. Willhalm et al.
@ VLDB 2009
(& improved using
latest SIMD ISA)



*LSB*          *MSB*

| 00010101 | 00110001 | 11110101 | 01100110 |

*8-bit —> 4-bit*

| 0001 | 0101 | 0011 | 0001 | 1111 | 0101 | 0110 | 0110 |

**shuffle**

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

*4-bit —> 8-bit*

| 00010101 | 01010011 | 00110001 | 00011111 |

**shift**   <<   <<   <<   <<

| 00010101 | 10100110 | 11000100 | 11111000 |

**mask**   &   &   &   &

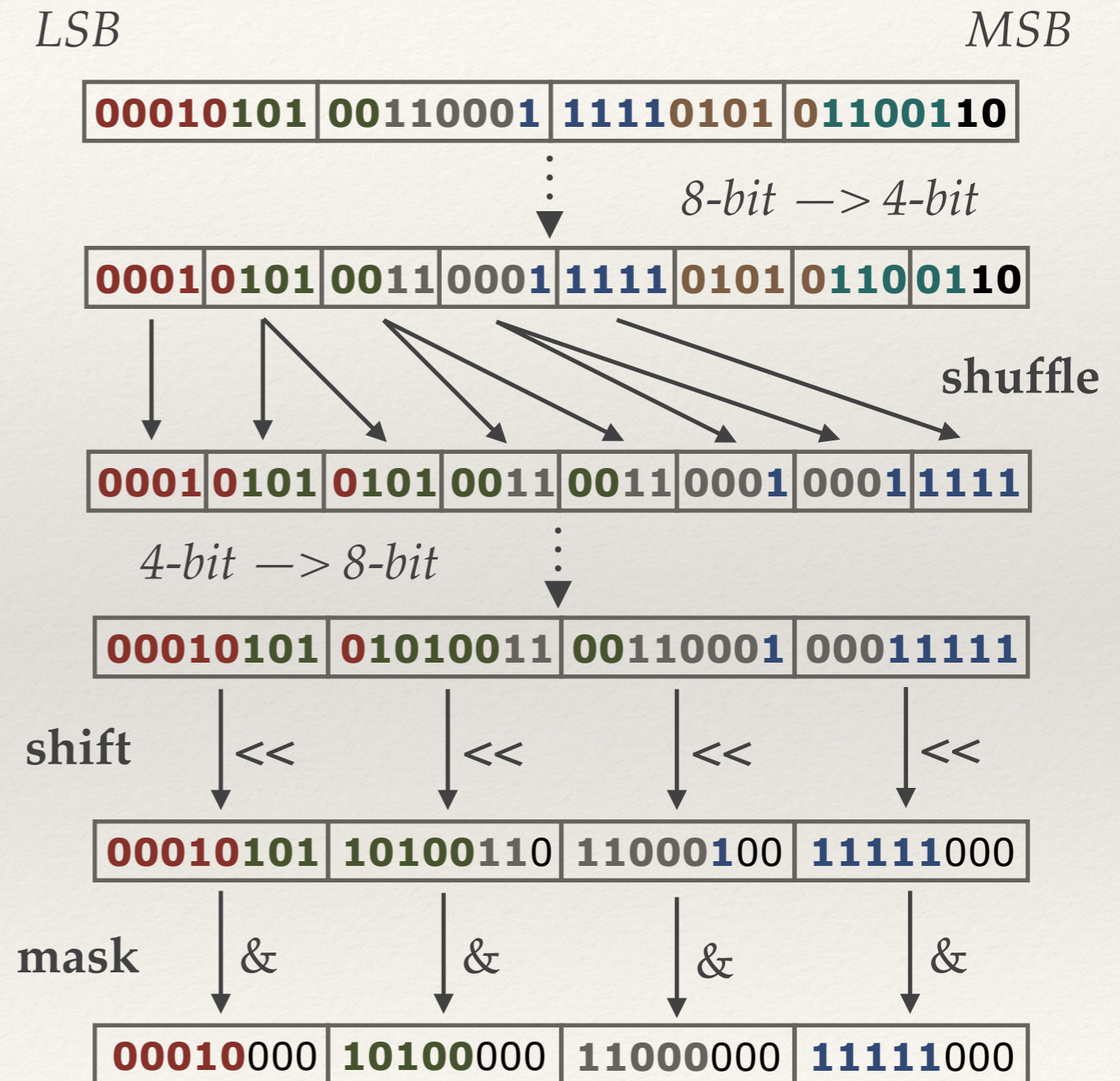| 00010000 | 10100000 | 11000000 | 11111000 |

# Horizontal Layout

- *Fully* packed
  - No space wasted
    - Codes can *span* across 2 packed words
  - Packing
    - Process 1 unpacked *code* per iteration
    - Branch to store output packed *word*
  - Unpacking
    - Process 1 output *code* per iteration
    - Branch to load input packed *word*
    - Can be written in *SIMD* !



*up to 7X improvement from SIMD*

# Horizontal Layout
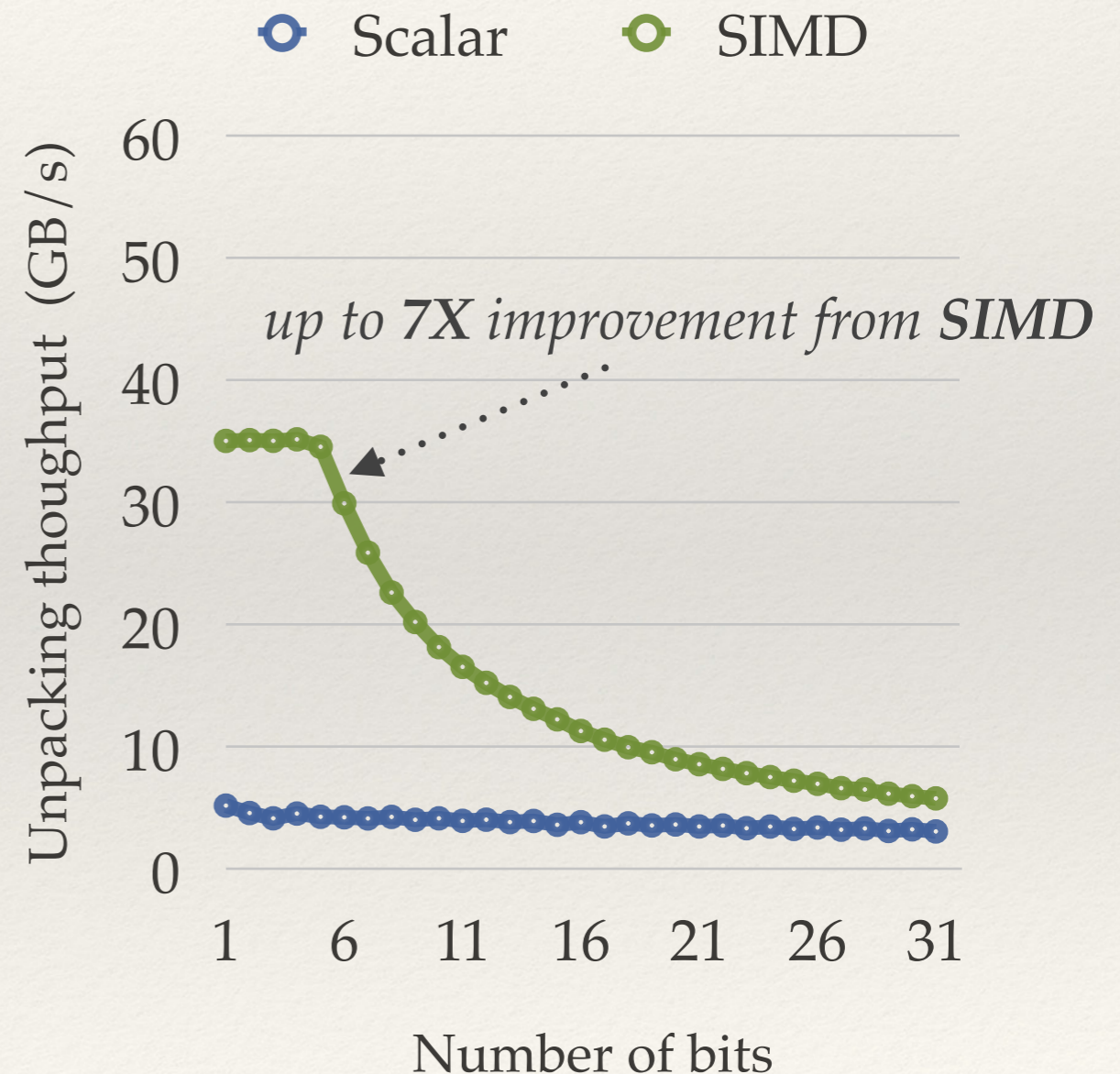
❖ *Fully* packed

    ❖ No space wasted

        ❖ Codes can *span* across 2 packed words

    ❖ Packing

        ❖ Process 1 unpacked *code* per iteration

        ❖ Branch to store output packed *word*

    ❖ Unpacking

        ❖ Process 1 output *code* per iteration

        ❖ Branch to load input packed *word*

        ❖ Can be written in *SIMD* !

    ❖ Scanning

        ❖ *Unpack* the codes in CPU registers

        ❖ *Evaluate* selective predicates and append to bitmap

        ❖ Must unpack first thus bounded by *O(n)*

# Horizontal Layout

❖ *Fully* packed

  ❖ No space wasted

    ❖ Codes can *span* across 2 packed words

  ❖ Packing

    ❖ Process 1 unpacked *code* per iteration

    ❖ Branch to store output packed *word*

  ❖ Unpacking

    ❖ Process 1 output *code* per iteration

    ❖ Branch to load input packed *word*

    ❖ Can be written in *SIMD* !

  ❖ Scanning

    ❖ *Unpack* the codes in CPU registers

    ❖ *Evaluate* selective predicates and append to bitmap

    ❖ Must unpack first thus bounded by *O(n)*

    ❖ Can also be written in *SIMD* via SIMD unpacking

*select … where column < C …*

| **00010101** | **00110001** | **11110101** | **01100110** |

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

**compare with C**

| **01100**000 | **01100**000 | **01100**000 | **01100**000 |

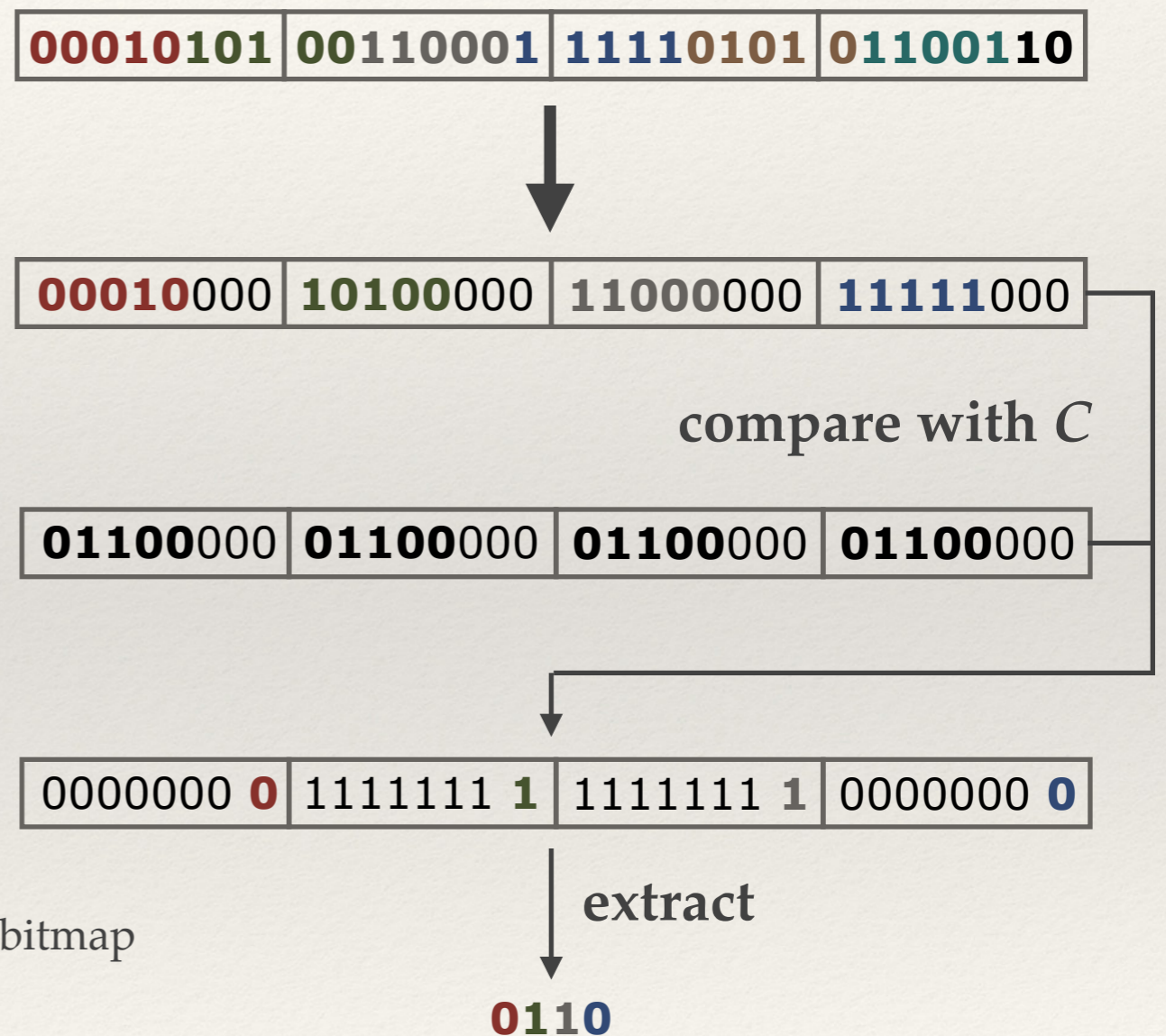| 0000000 **0** | 1111111 **1** | 1111111 **1** | 0000000 **0** |

**extract**

**0110**

# Horizontal Layout
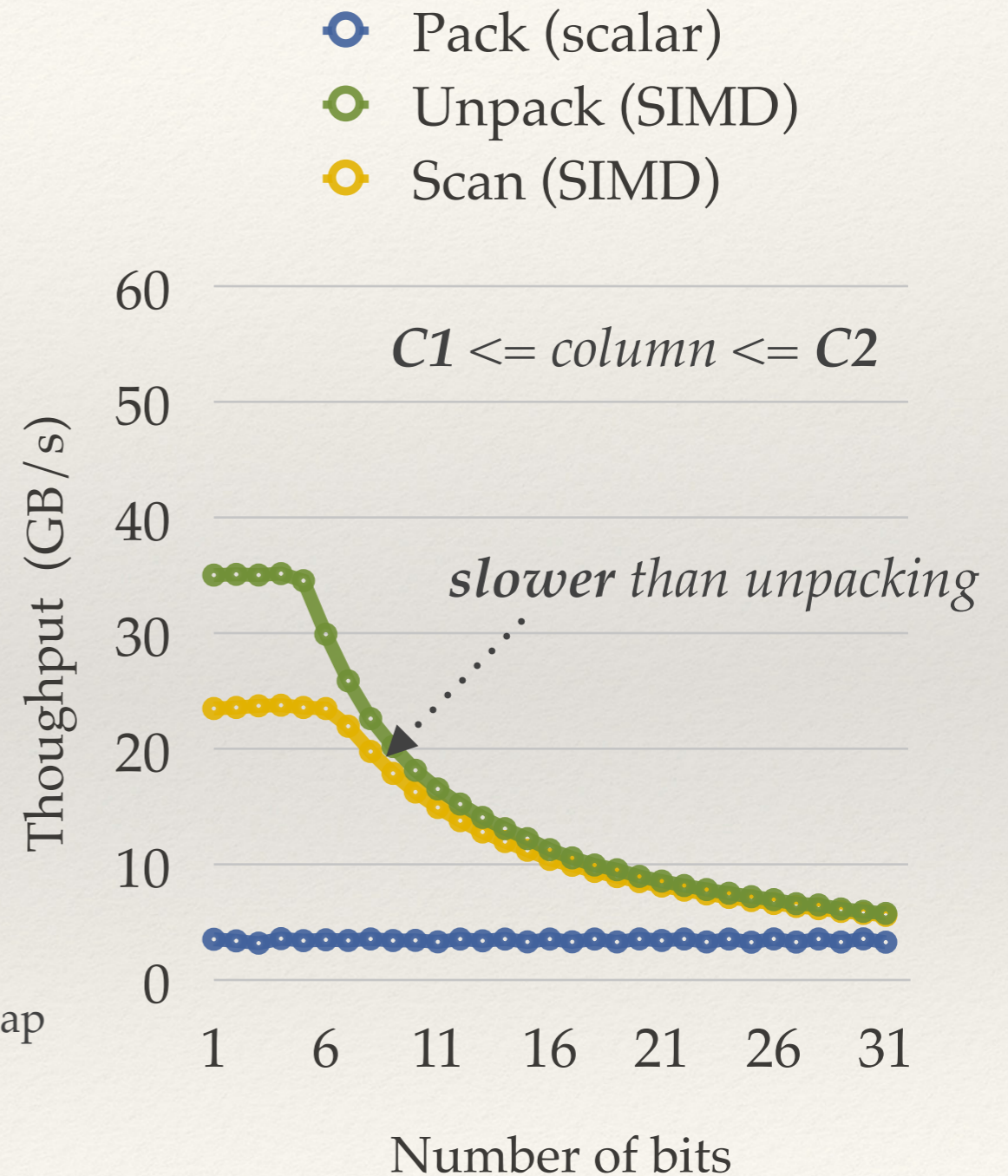
❖ *Fully* packed

    ❖ No space wasted

        ❖ Codes can *span* across 2 packed words

    ❖ Packing

        ❖ Process 1 unpacked *code* per iteration

        ❖ Branch to store output packed *word*

    ❖ Unpacking

        ❖ Process 1 output *code* per iteration

        ❖ Branch to load input packed *word*

        ❖ Can be written in *SIMD* !

    ❖ Scanning

        ❖ *Unpack* the codes in CPU registers

        ❖ *Evaluate* selective predicates and append to bitmap

        ❖ Must unpack first thus bounded by *O(n)*

        ❖ Can also be written in *SIMD* via SIMD unpacking

- ○ Pack (scalar)
- ○ Unpack (SIMD)
- ○ Scan (SIMD)

*C1 <= column <= C2*

*slower* than unpacking

Thoughput (GB/s) vs Number of bits

# Horizontal Layout

- ❖ Word *aligned*
  - ❖ Waste space to get alignment
    - ❖ Pack *b' = w / (b+1)* codes per processor word
    - ❖ Extra bit per word used for *scanning*

*fully* packed

**01 10 11 00**

**01**0 **10**0 00  **11**0 **00**0 00

*word* **aligned**

*unused* high order bits per word

**01**0 **10**0 00

*unused* extra bit per code

# Horizontal Layout

- Word *aligned*
  - Waste space to get alignment
    - Pack $b' = w / (b+1)$ codes per processor word
    - Extra bit per word used for *scanning*
  - Packing
    - 1 packed word at a time
    - Nested loop to pack $b'$ codes

# Horizontal Layout

❖ Word *aligned*

  ❖ Waste space to get alignment

    ❖ Pack $b' = w / (b+1)$ codes per processor word

      ❖ Extra bit per word used for *scanning*

  ❖ Packing

    ❖ 1 packed word at a time

    ❖ Nested loop to pack *b'* codes

  ❖ Unpacking

    ❖ 1 packed word at a time

    ❖ Nested loop to unpack *b'* codes

Legend:
- Fully packed (scalar)
- Fully packed (SIMD)
- Word aligned (scalar)

*slower* than SIMD

X-axis: Number of bits (1, 6, 11, 16, 21, 26, 31)
Y-axis: Unpacking thoughput (GB/s) (0, 10, 20, 30, 40, 50, 60)
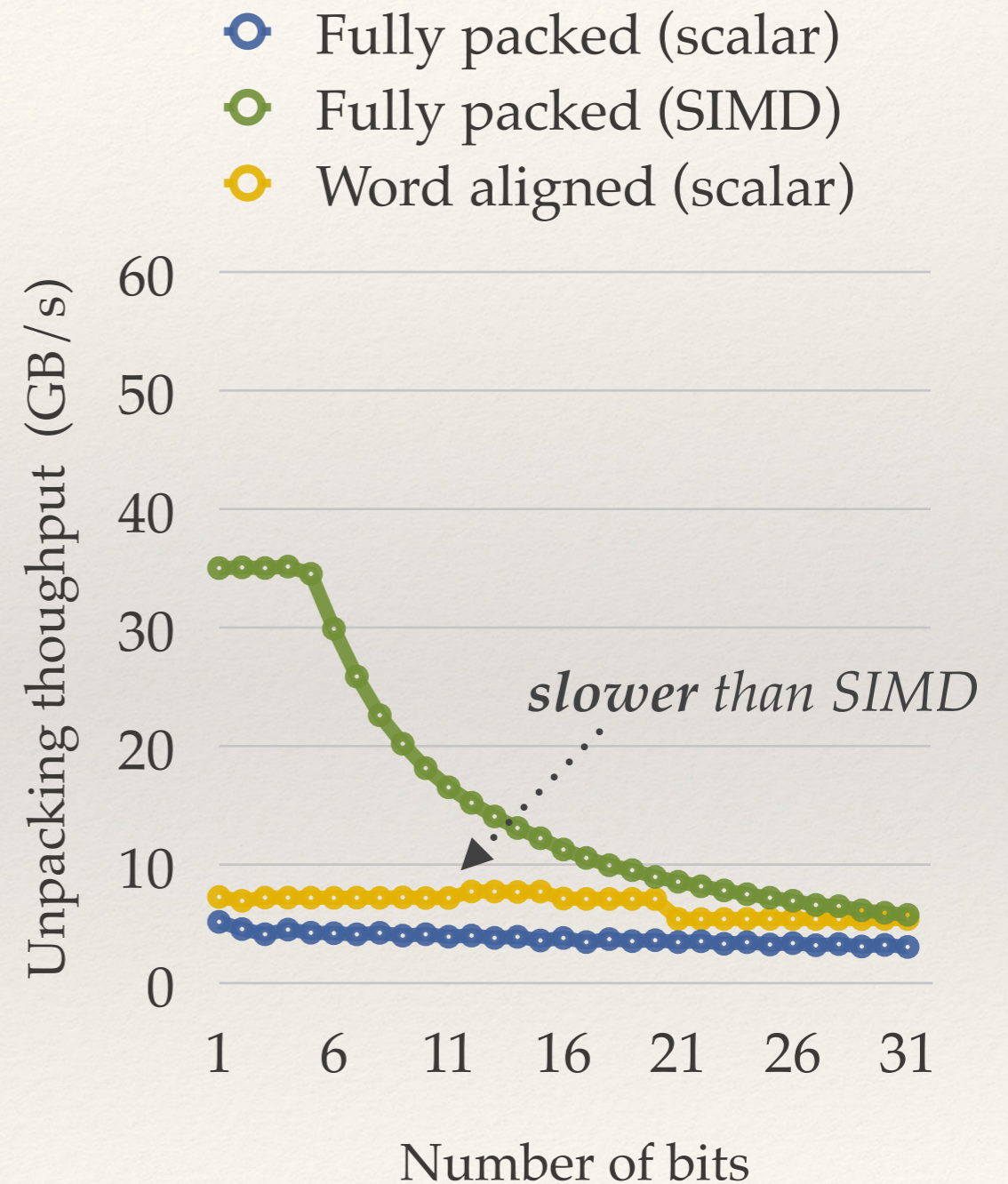
# Horizontal Layout

- Word *aligned*
  - Waste space to get alignment
    - Pack $b' = w / (b+1)$ codes per processor word
    - Extra bit per word used for *scanning*
  - Packing
    - 1 packed word at a time
    - Nested loop to pack $b'$ codes
  - Unpacking
    - 1 packed word at a time
    - Nested loop to unpack $b'$ codes
  - Scanning
    - Evaluate predicates *without* unpacking
    - Works with *simple* order predicates: <,=,>
    - Boolean result in *overflow* bit of *b*-bit arithmetic
    - Executing $< O(n)$ operations

*select … where column < C …*

**01**

**set constant C**

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

**invert code bits**

**01**0 **10**0 00 ^ **11**0 110 00 = **10**0 **01**0 00

**add constant**

**10**0 **01**0 00 + **01**0 **01**0 00 = 11**0** 001**1** 00

**extract sign**

11**0** 001**1** 00 —> **01**

## Based on paper by Leslie Lamport @ CACM 1975

# Horizontal Layout
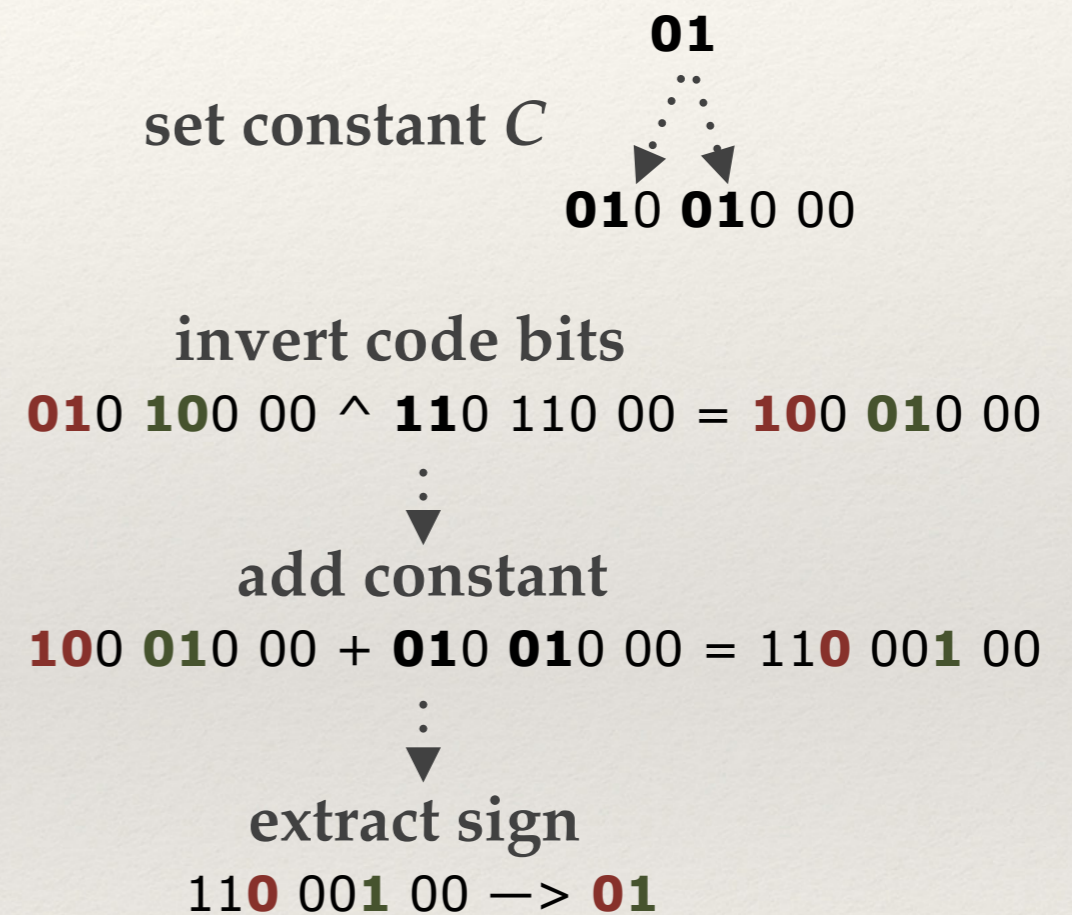
- Word *aligned*
  - Waste space to get alignment
    - Pack $b' = w / (b+1)$ codes per processor word
    - Extra bit per word used for *scanning*
  - Packing
    - 1 packed word at a time
    - Nested loop to pack $b'$ codes
  - Unpacking
    - 1 packed word at a time
    - Nested loop to unpack $b'$ codes
  - Scanning
    - Evaluate predicates *without* unpacking
    - Works with *simple* order predicates: <,=,>
    - Boolean result in *overflow* bit of $b$-bit arithmetic
    - Executing $< O(n)$ operations

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

$C1 <= column <= C2$

*faster* due to not unpacking

*& slower* due to wasted space

Scanning thoughput (GB/s)

60
50
40
30
20
10
0

1   6   11   16   21   26   31

Number of bits

# Vertical Layout

❖ *Fully* packed & word *aligned*

  ❖ Interleave bits of $k$ codes

    ❖ $k$ divides the processor word

# Vertical Layout

- *Fully* packed & word *aligned*
  - Interleave bits of *k* codes
    - *k* divides the processor word
  - Scanning
    - Evaluate *without* unpacking
    - Can *skip* words early

```
00010110   00000000   11101001   ___0_00_
10011000   11111111   100_1__0   _110_001
01010001   11111111   0___0___   11101001
00111100   00000000
01110110   00000000
```

"="   $X \&= \sim(column \wedge C)$ ......

"<"   $Y \mathrel{|}= C \& (\sim X)$ ..........................

**stop if** $X = 0$

Based on paper by
Y. Li et al.
@ SIGMOD 2013

# Vertical Layout

- *Fully* packed & word *aligned*
  - Interleave bits of *k* codes
    - *k* divides the processor word
  - Scanning
    - Evaluate *without* unpacking
    - Can *skip* words early



Vertical k = 64 (scalar)
Horizontal full (SIMD)
Horizontal word (scalar)

*fastest* in most cases

Scanning thoughput (GB/s)

Number of bits

# Vertical Layout

❖ *Fully* packed & word *aligned*
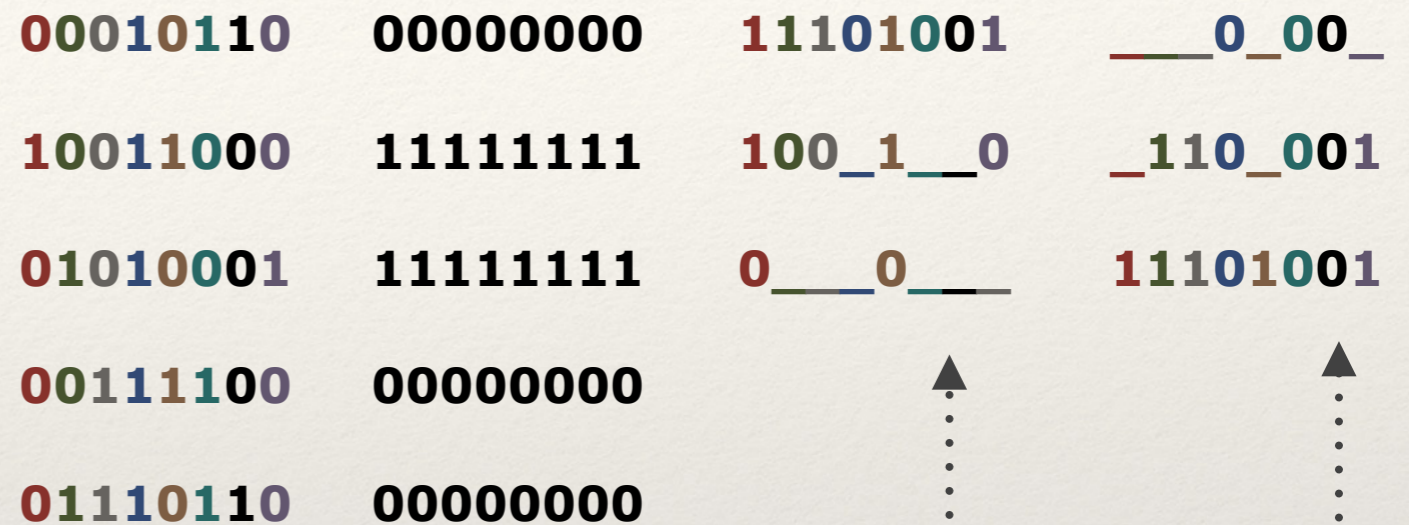
  ❖ Interleave bits of *k* codes

    ❖ *k* divides the processor word

  ❖ Scanning

    ❖ Evaluate *without* unpacking

    ❖ Can *skip* words early

    ❖ Increase *k* to minimize false (pre)fetches

$k = 64$
$b = 5$

$k = 256$
$b = 5$

# Vertical Layout

❖ *Fully* packed & word *aligned*

    ❖ Interleave bits of *k* codes

        ❖ *k* divides the processor word

    ❖ Scanning

        ❖ Evaluate *without* unpacking

        ❖ Can *skip* words early

        ❖ Increase *k* to minimize false (pre)fetches



Legend:
- Vertical k = 64 (scalar)
- Horizontal full (SIMD)
- Horizontal word (scalar)
- Vertical k = 8192 (SIMD)

*faster due to cache line* **skip**

Y-axis: Scanning thoughput (GB/s)

X-axis: Number of bits

# Vertical Layout

❖ *Fully* packed & word *aligned*

   ❖ Interleave bits of $k$ codes

      ❖ $k$ divides the processor word

   ❖ Scanning

      ❖ Evaluate *without* unpacking

      ❖ Can *skip* words early

      ❖ Increase $k$ to minimize false (pre)fetches

   ❖ Packing

      ❖ Transfer *nb* bits across registers

00100000
10001000
11001000
01010000
11111000
11000000
10100000
00010000

00000000    00000000
10000000    10000000
10000000    11000000
10000000    01100000
00000000    10110000
10000000    01011000
00000000    00101100
00000000    00010110

00000000    00000000
00000000    10000000
10000000    11000000
00000000    01100000
10000000    11110000
00000000    11111000
00000000    00111100
00000000    00011110

# Vertical Layout

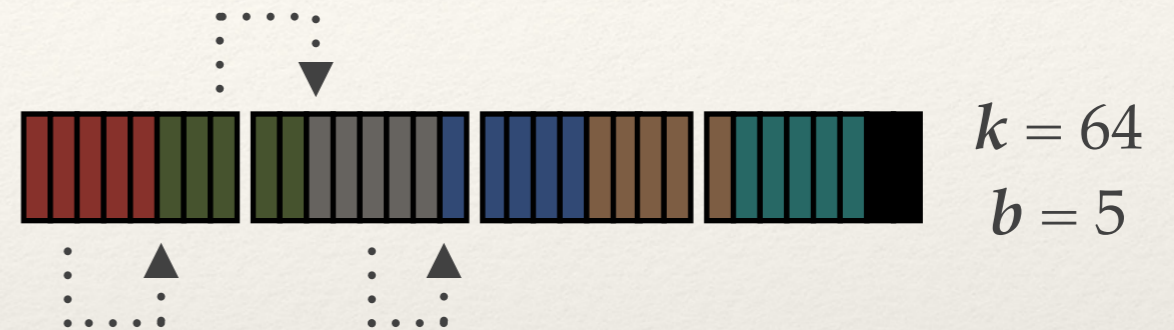- *Fully* packed & word *aligned*
  - Interleave bits of *k* codes
    - *k* divides the processor word
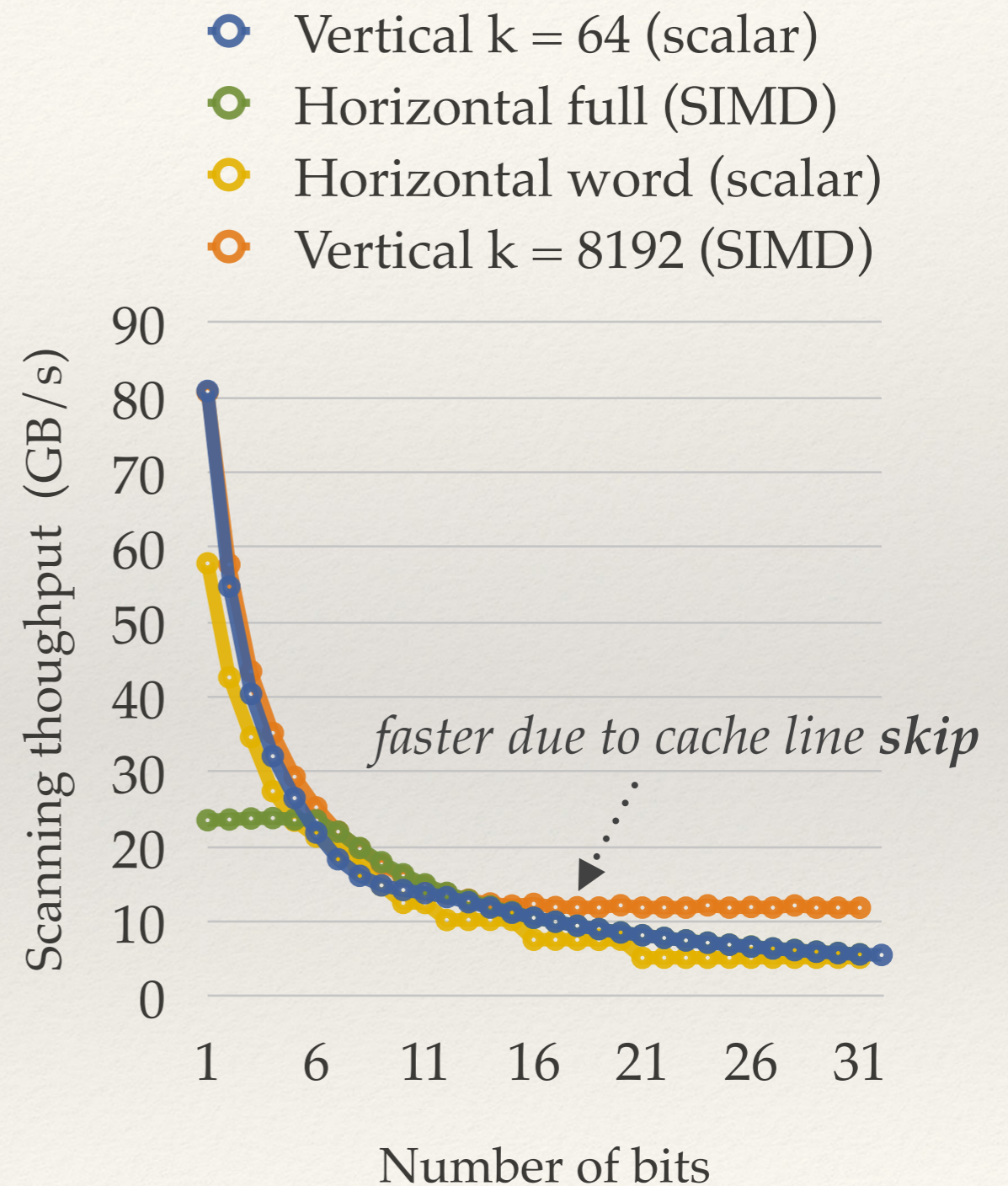  - Scanning
    - Evaluate *without* unpacking
    - Can *skip* words early
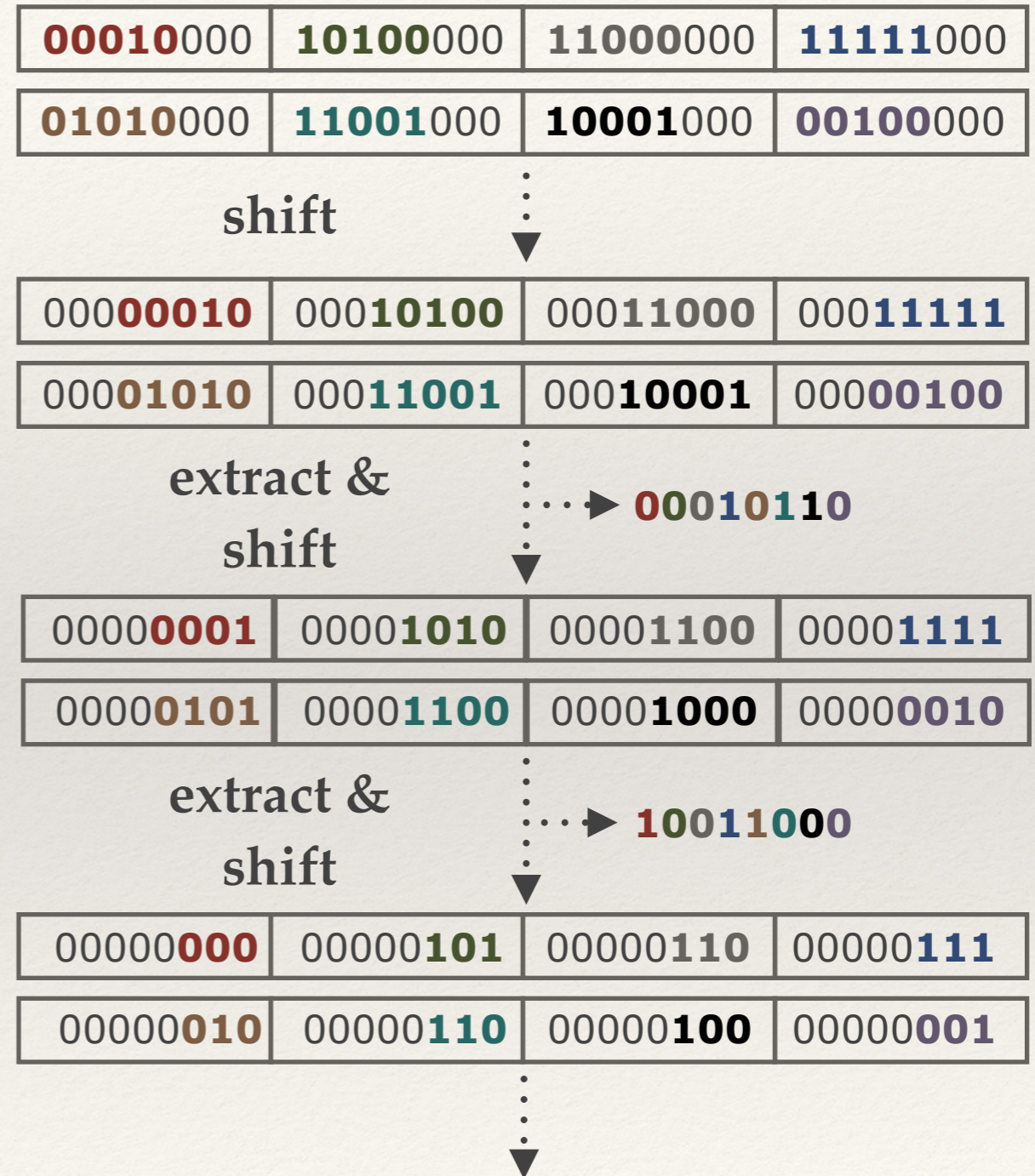    - Increase *k* to minimize false (pre)fetches
  - Packing
    - Transfer *nb* bits across registers
    - Can be written in *SIMD* !

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

**shift**

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

**extract & shift** → **00010110**

| | | | |
|---|---|---|---|
| 00000**0001** | 00000**1010** | 00000**1100** | 00000**1111** |
| 00000**0101** | 00000**1100** | 00000**1000** | 00000**0010** |

**extract & shift** → **10011000**

| | | | |
|---|---|---|---|
| 00000**000** | 00000**101** | 00000**110** | 00000**111** |
| 00000**010** | 00000**110** | 00000**100** | 00000**001** |

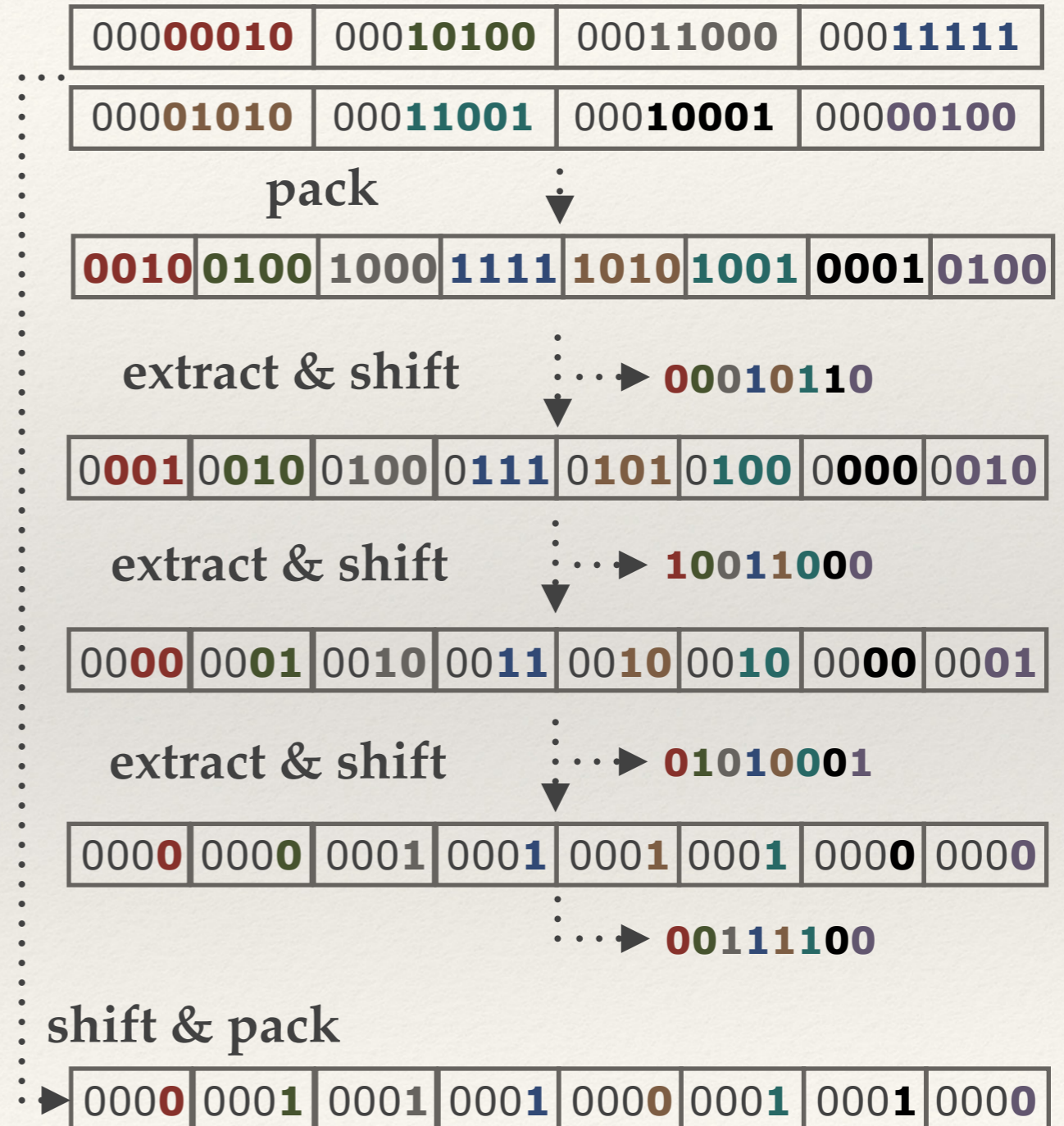# Vertical Layout

❖ *Fully* packed & word *aligned*

  ❖ Interleave bits of *k* codes

    ❖ *k* divides the processor word

  ❖ Scanning

    ❖ Evaluate *without* unpacking

    ❖ Can *skip* words early

    ❖ Increase *k* to minimize false (pre)fetches

  ❖ Packing

    ❖ Transfer *nb* bits across registers

    ❖ Can be written in *SIMD* !

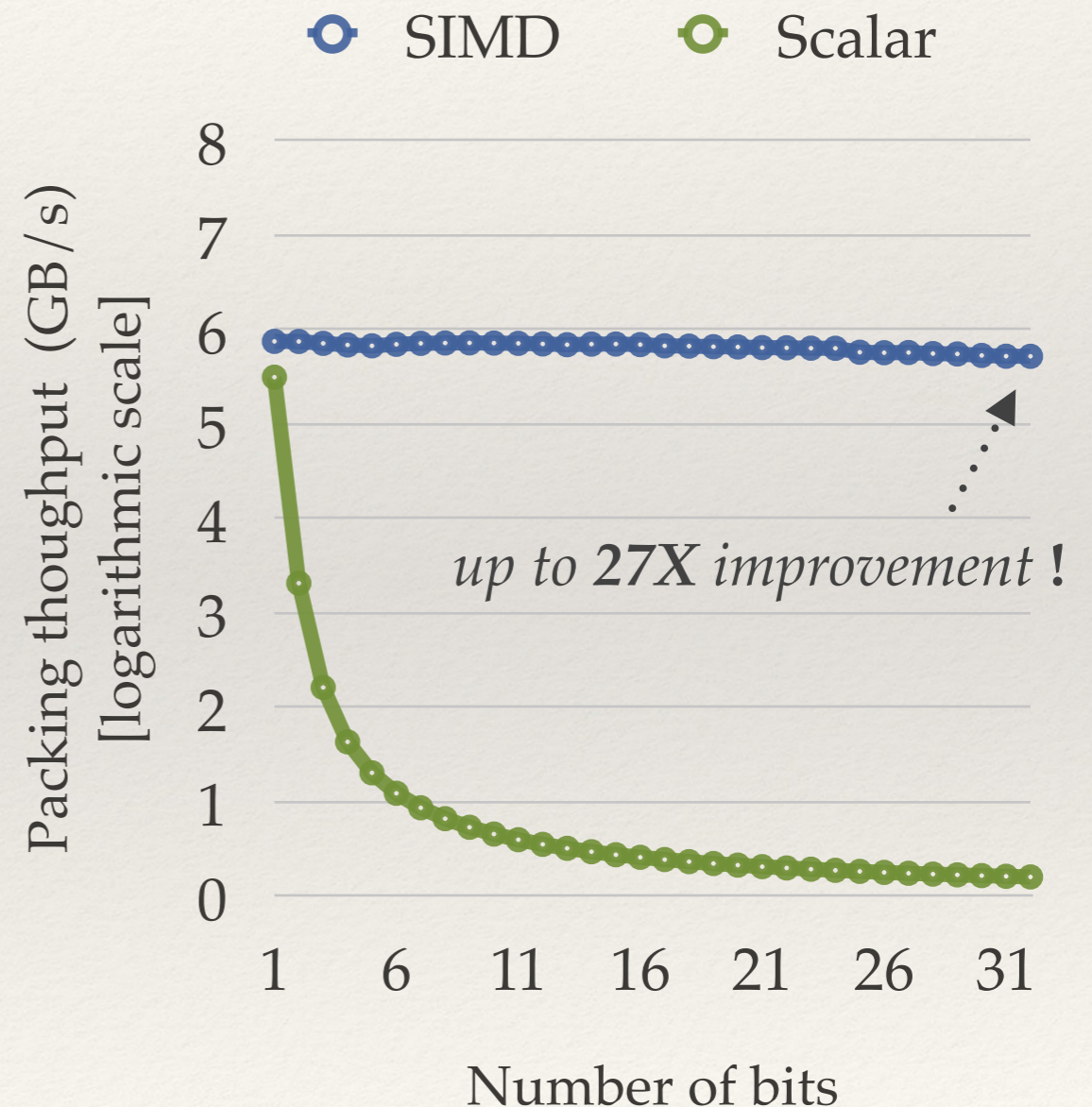    ❖ Extract bits *per byte* not per int

# Vertical Layout

❖ *Fully* packed & word *aligned*

    ❖ Interleave bits of *k* codes

        ❖ *k* divides the processor word

    ❖ Scanning

        ❖ Evaluate *without* unpacking

        ❖ Can *skip* words early

        ❖ Increase *k* to minimize false (pre)fetches

    ❖ Packing

        ❖ Transfer *nb* bits across registers

        ❖ Can be written in *SIMD* !

        ❖ Extract bits *per byte* not per int
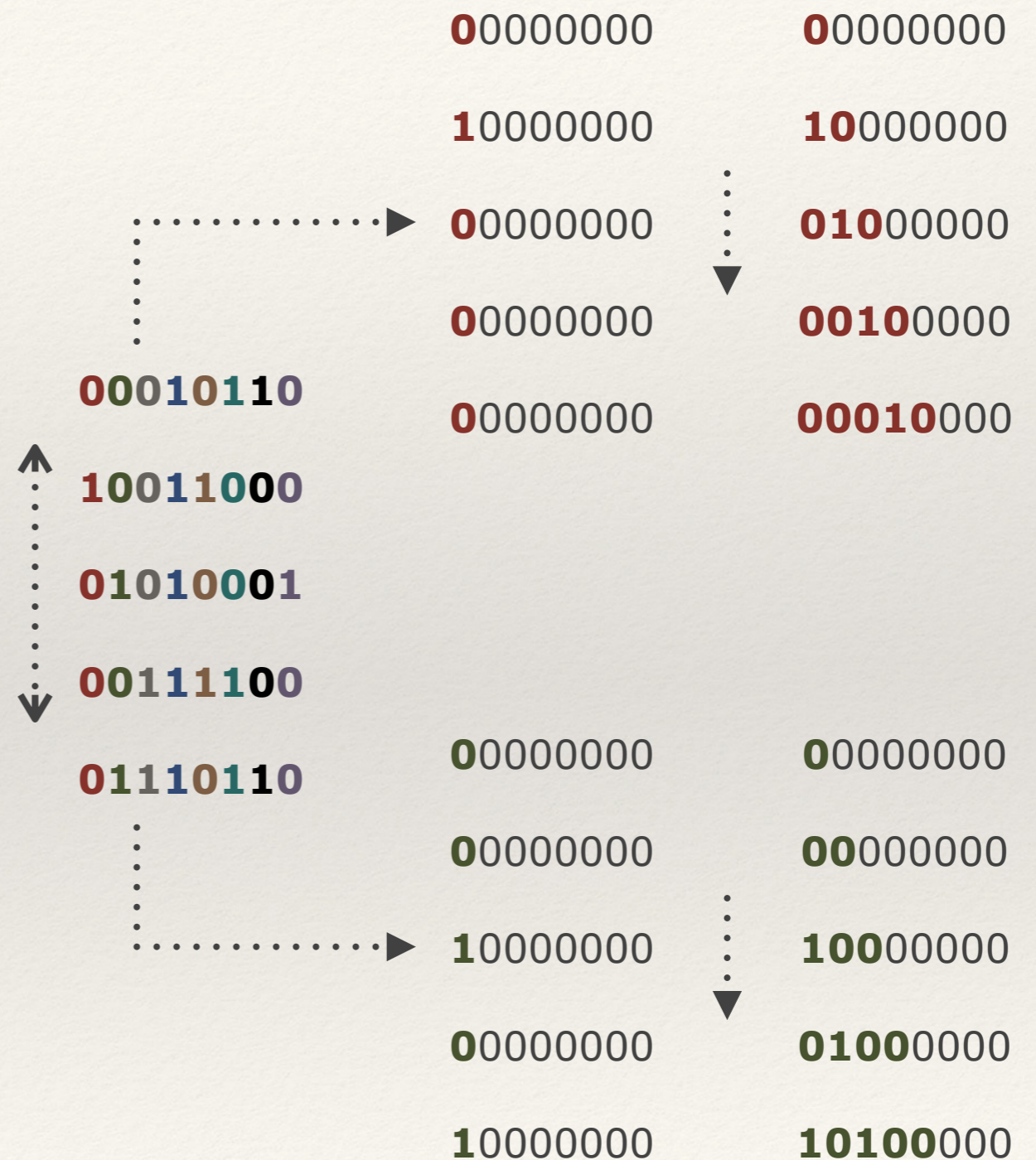


*up to 27X improvement* !

# Vertical Layout

- *Fully* packed & word *aligned*
  - Interleave bits of *k* codes
    - *k* divides the processor word
  - Scanning
    - Evaluate *without* unpacking
    - Can *skip* words early
    - Increase *k* to minimize false (pre)fetches
  - Packing
    - Transfer *nb* bits across registers
    - Can be written in *SIMD* !
    - Extract bits *per byte* not per int
  - Unpacking
    - Transfer *nb* bits across registers

**0**0000000   **0**0000000
**1**0000000   **1**0000000
**0**0000000   **01**000000
**0**0000000   **001**00000
**0**0000000   **0001**0000

**00010110**
**10011000**
**01010001**
**00111100**

**01110110**

**0**0000000   **0**0000000
**0**0000000   **00**000000
**1**0000000   **100**00000
**0**0000000   **0100**0000
**1**0000000   **10100**000

# Vertical Layout

- *Fully* packed & word *aligned*
  - Interleave bits of $k$ codes
    - $k$ divides the processor word
  - Scanning
    - Evaluate *without* unpacking
    - Can *skip* words early
    - Increase $k$ to minimize false (pre)fetches
  - Packing
    - Transfer *nb* bits across registers
    - Can be written in *SIMD* !
    - Extract bits *per byte* not per int
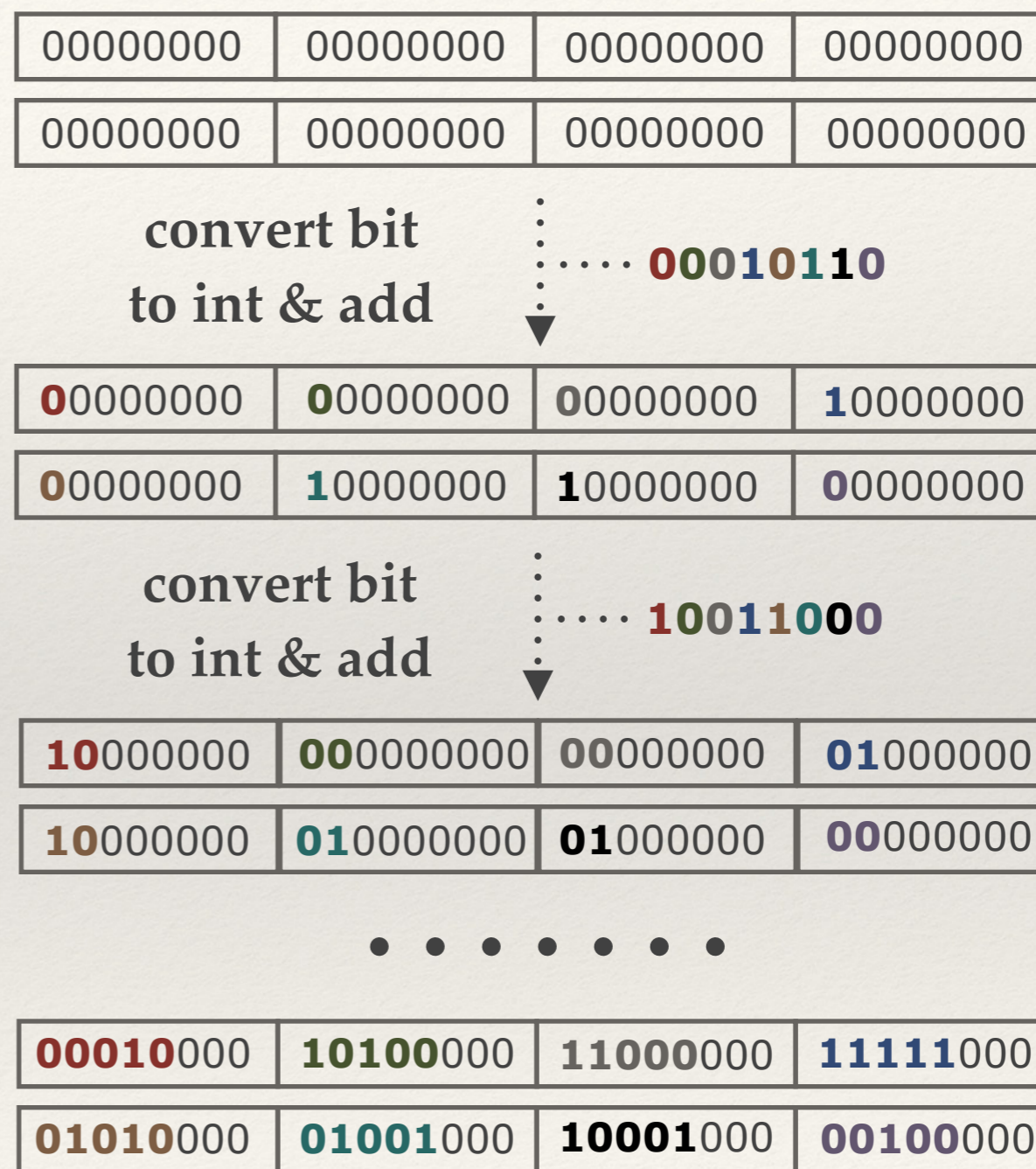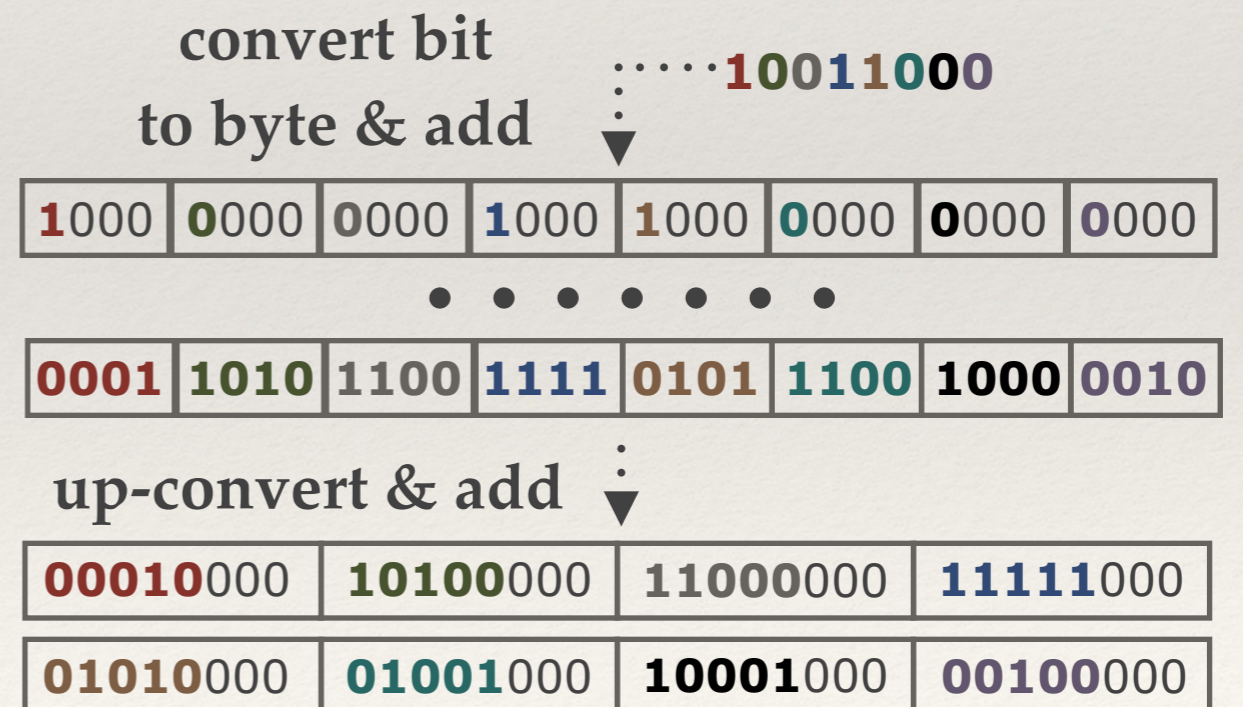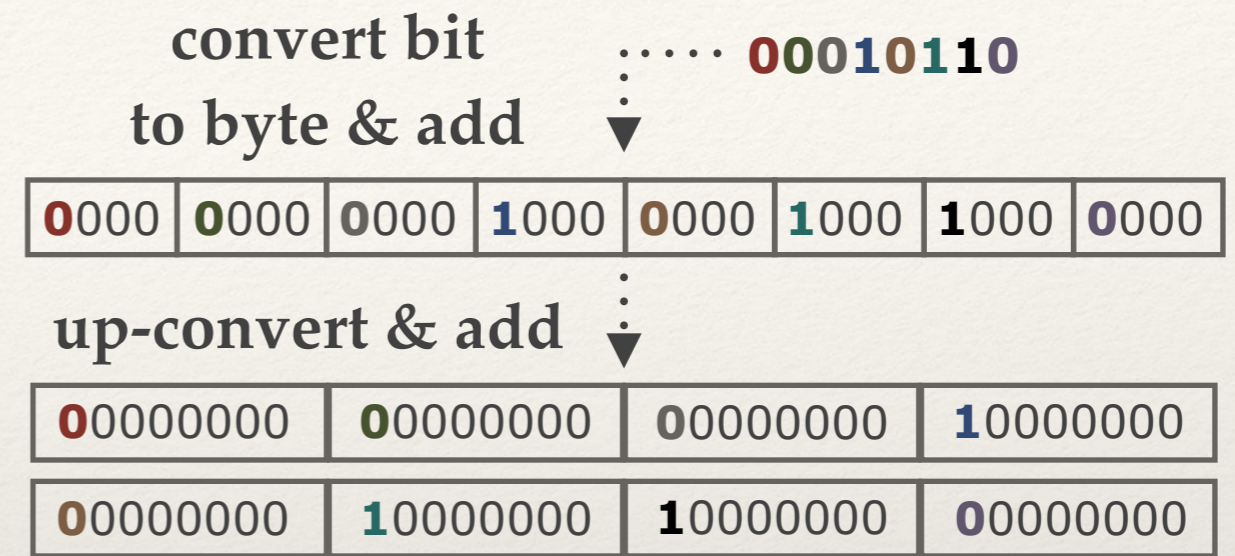  - Unpacking
    - Transfer *nb* bits across registers
    - Can be written in *SIMD* !

| | | | |
|---|---|---|---|
| 00000000 | 00000000 | 00000000 | 00000000 |
| 00000000 | 00000000 | 00000000 | 00000000 |

**convert bit to int & add** ······ 00010110

| | | | |
|---|---|---|---|
| 00000000 | 00000000 | 00000000 | 10000000 |
| 00000000 | 10000000 | 10000000 | 00000000 |

**convert bit to int & add** ····· 10011000

| | | | |
|---|---|---|---|
| 10000000 | 000000000 | 00000000 | 01000000 |
| 10000000 | 010000000 | 01000000 | 00000000 |

• • • • • •

| | | | |
|---|---|---|---|
| 00010000 | 10100000 | 11000000 | 11111000 |
| 01010000 | 01001000 | 10001000 | 00100000 |

# Vertical Layout

❖ *Fully* packed & word *aligned*

   ❖ Interleave bits of *k* codes

      ❖ *k* divides the processor word

   ❖ Scanning

      ❖ Evaluate *without* unpacking

      ❖ Can *skip* words early

      ❖ Increase *k* to minimize false (pre)fetches

   ❖ Packing

      ❖ Transfer *nb* bits across registers

      ❖ Can be written in *SIMD* !

      ❖ Extract bits *per byte* not per int

   ❖ Unpacking

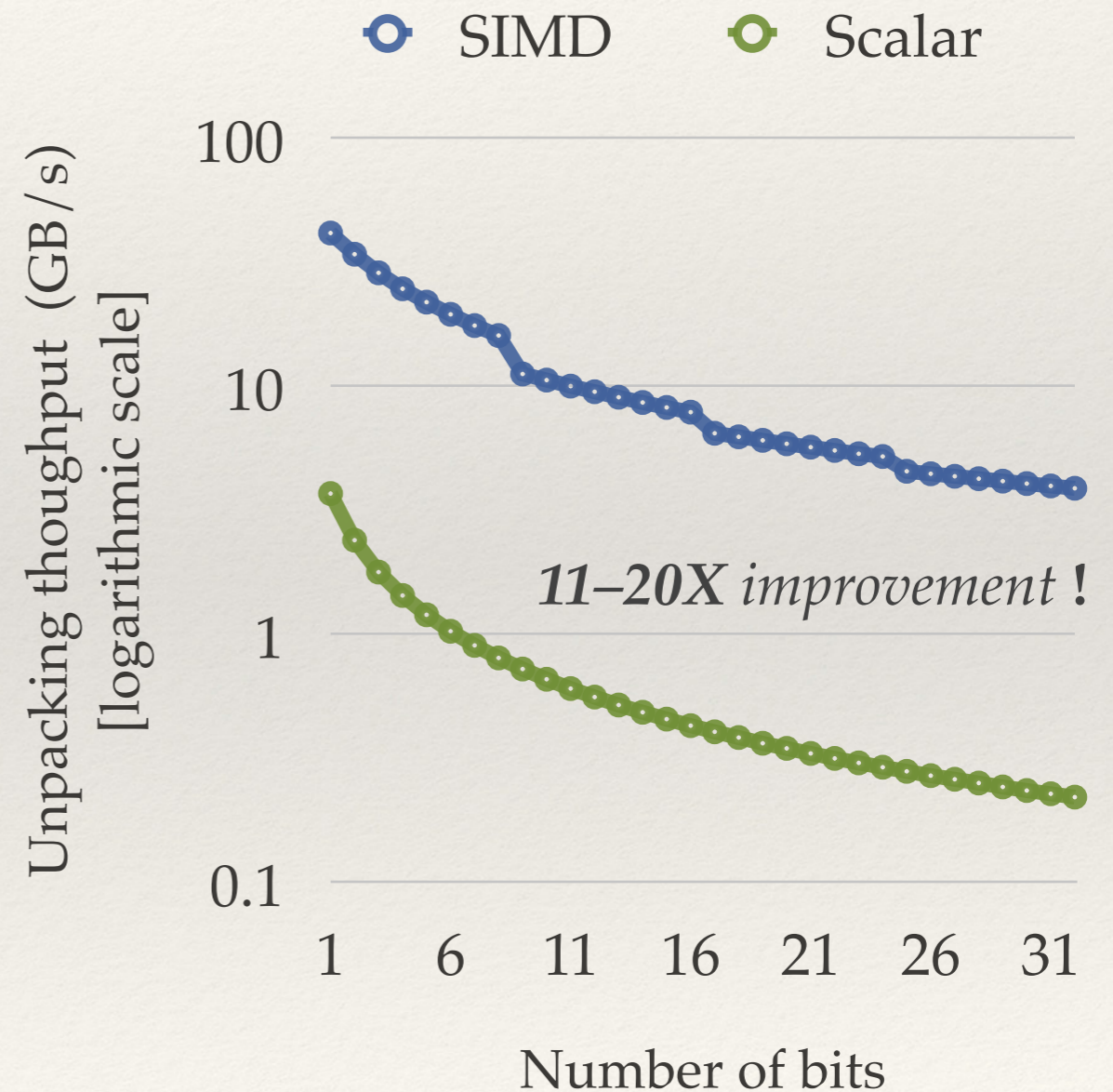      ❖ Transfer *nb* bits across registers

      ❖ Can be written in *SIMD* !

      ❖ Insert bits *per byte* not per int

**convert bit to byte & add**  ····· **00010110**

| 0000 | 0000 | 0000 | 1000 | 0000 | 1000 | 1000 | 0000 |
|------|------|------|------|------|------|------|------|

**up-convert & add**

| 00000000 | 00000000 | 00000000 | 10000000 |
|----------|----------|----------|----------|

| 00000000 | 10000000 | 10000000 | 00000000 |
|----------|----------|----------|----------|

**convert bit to byte & add**  ····· **10011000**

| 1000 | 0000 | 0000 | 1000 | 1000 | 0000 | 0000 | 0000 |
|------|------|------|------|------|------|------|------|

• • • • • • •

| 0001 | 1010 | 1100 | 1111 | 0101 | 1100 | 1000 | 0010 |
|------|------|------|------|------|------|------|------|

**up-convert & add**

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

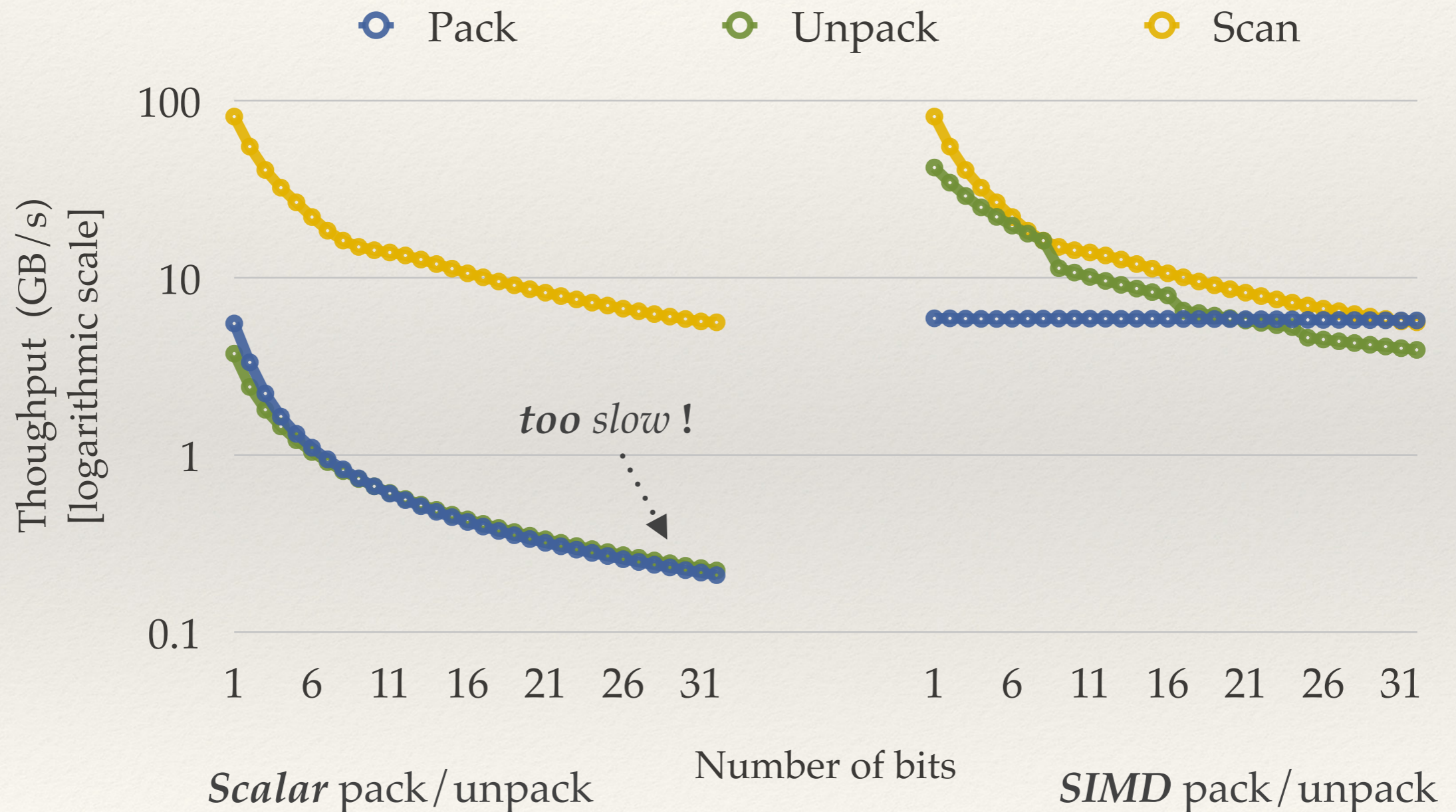| 01010000 | 01001000 | 10001000 | 00100000 |
|----------|----------|----------|----------|

# Vertical Layout

- *Fully* packed & word *aligned*
  - Interleave bits of *k* codes
    - *k* divides the processor word
  - Scanning
    - Evaluate *without* unpacking
    - Can *skip* words early
    - Increase *k* to minimize false (pre)fetches
  - Packing
    - Transfer *nb* bits across registers
    - Can be written in *SIMD* !
    - Extract bits *per byte* not per int
  - Unpacking
    - Transfer *nb* bits across registers
    - Can be written in *SIMD* !
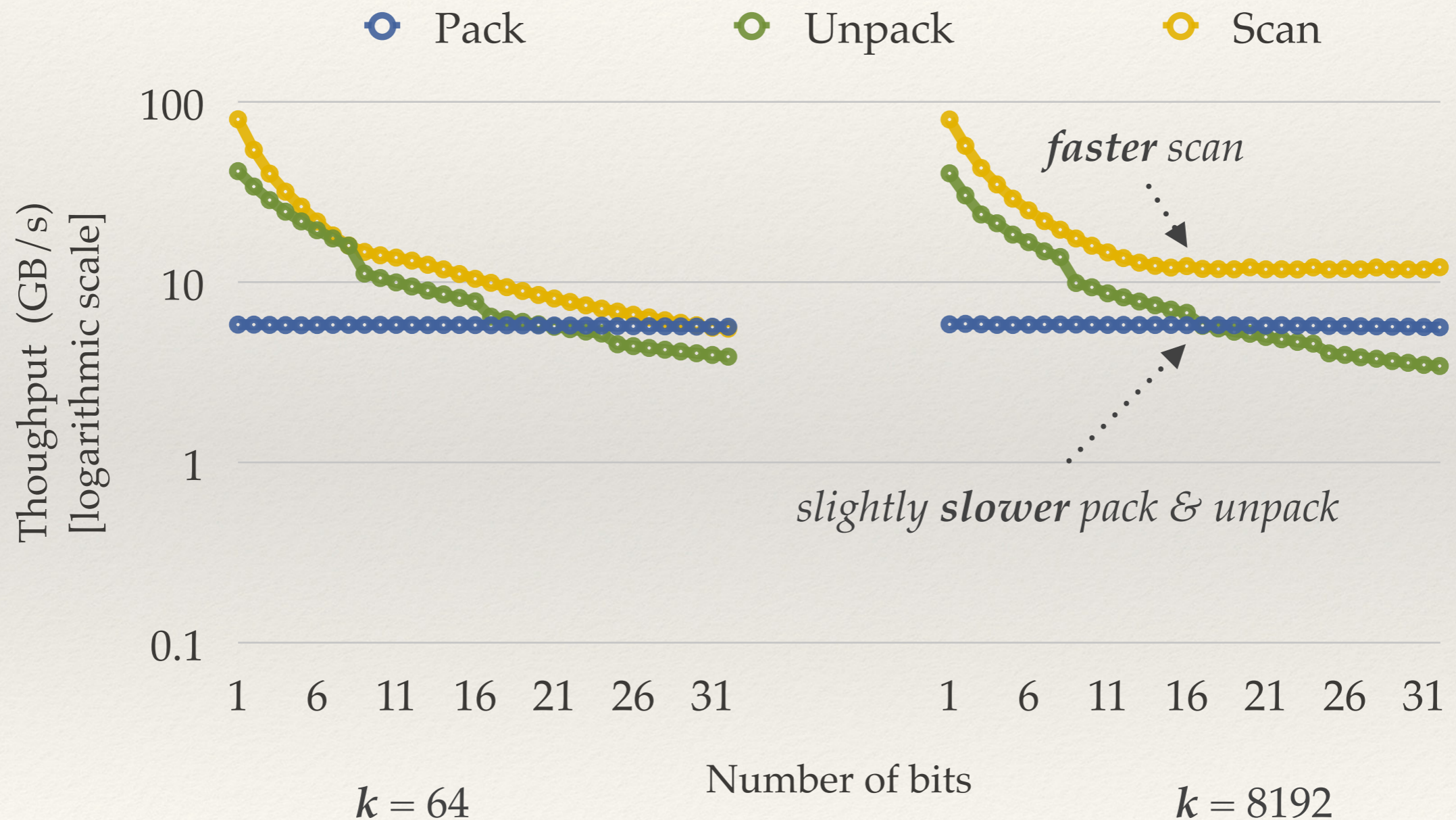    - Insert bits *per byte* not per int

# Vertical Layout

- Scalar to *SIMD* for packing & unpacking  ($k = 64$)
  - Scalar scan



Pack    Unpack    Scan

*too* slow !

Thoughput (GB/s) [logarithmic scale]

*Scalar* pack/unpack     Number of bits     *SIMD* pack/unpack

# Vertical Layout

- ❖ Increasing **k** to the L1 cache size ($k = 8192$)
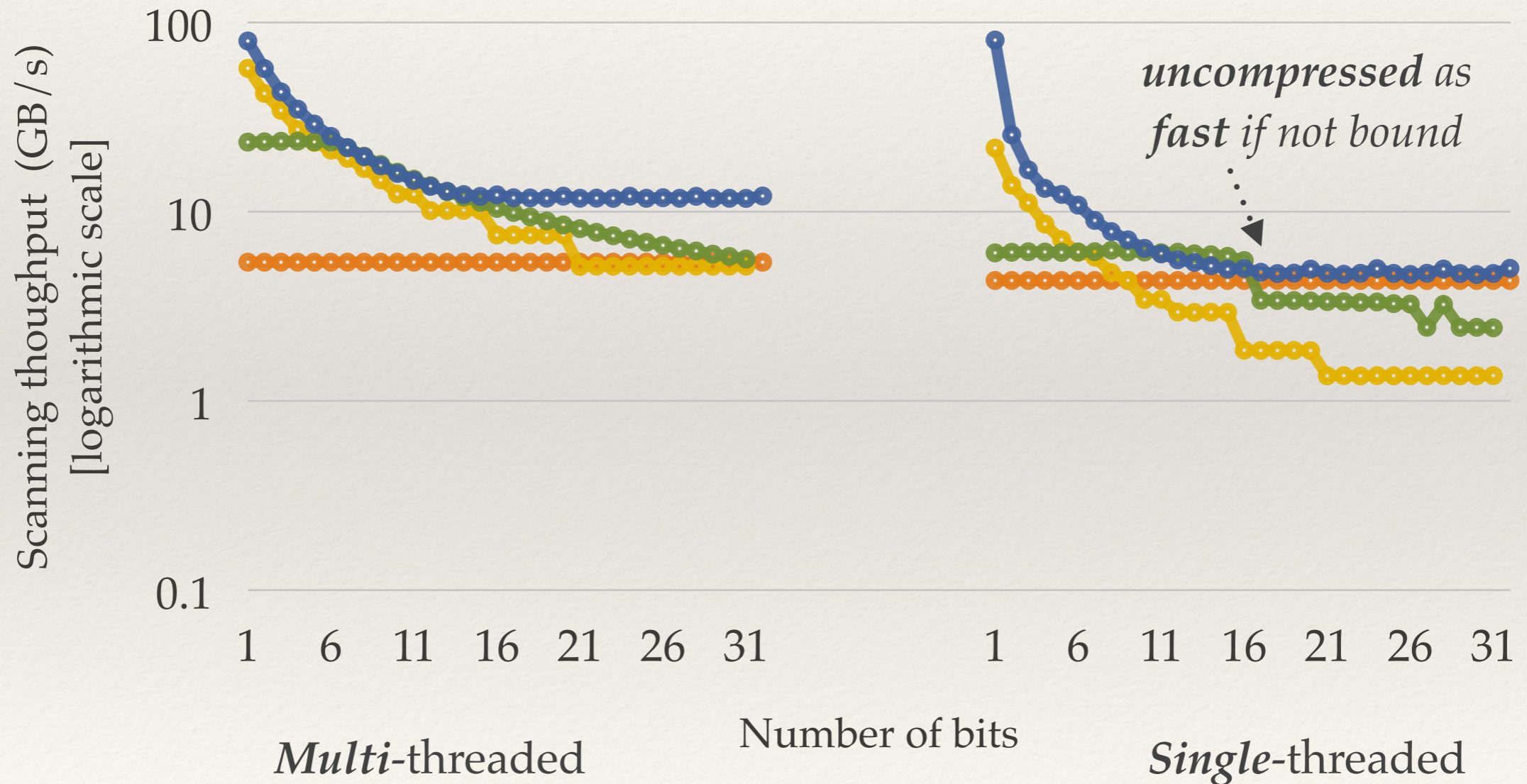  - ❖ SIMD scanning

# Vertical Layout

- If not memory **bound**
  - Using 1 thread



**Vertical (k = 8192)**    **Horizontal full**
**Horizontal word**    **Uncompressed**

Scanning thoughput (GB/s) [logarithmic scale]

*uncompressed as fast if not bound*

Number of bits

*Multi*-threaded    *Single*-threaded

# Conclusions

❖ Horizontal layouts

  ❖ Fully packed

    ❖ No wasted space but somewhat *slow*

    ❖ Can optimize unpacking & scanning with *SIMD*

# Conclusions

❖ Horizontal layouts

  ❖ Fully packed

    ❖ No wasted space but somewhat *slow*

    ❖ Can optimize unpacking & scanning with *SIMD*

  ❖ Word aligned

    ❖ Fast *scalar* scans but not optimal due to wasted space

# Conclusions

❖ Horizontal layouts

    ❖ Fully packed

        ❖ No wasted space but somewhat *slow*

        ❖ Can optimize unpacking & scanning with *SIMD*

    ❖ Word aligned

        ❖ Fast *scalar* scans but not optimal due to wasted space

❖ Vertical layout

    ❖ Known techniques

        ❖ Fast scalar scans *without* wasting space

        ❖ Very slow *scalar* packing & unpacking

# Conclusions

❖ Horizontal layouts

  ❖ Fully packed

    ❖ No wasted space but somewhat *slow*

    ❖ Can optimize unpacking & scanning with *SIMD*

  ❖ Word aligned

    ❖ Fast *scalar* scans but not optimal due to wasted space


❖ Vertical layout

  ❖ Known techniques

    ❖ Fast scalar scans *without* wasting space

    ❖ Very slow *scalar* packing & unpacking

  ❖ New techniques

    ❖ Fast packing & unpacking using *SIMD*

    ❖ Maximize bit transfers by using the *smallest* SIMD lanes

    ❖ Increase $k$ to skip cache lines effectively