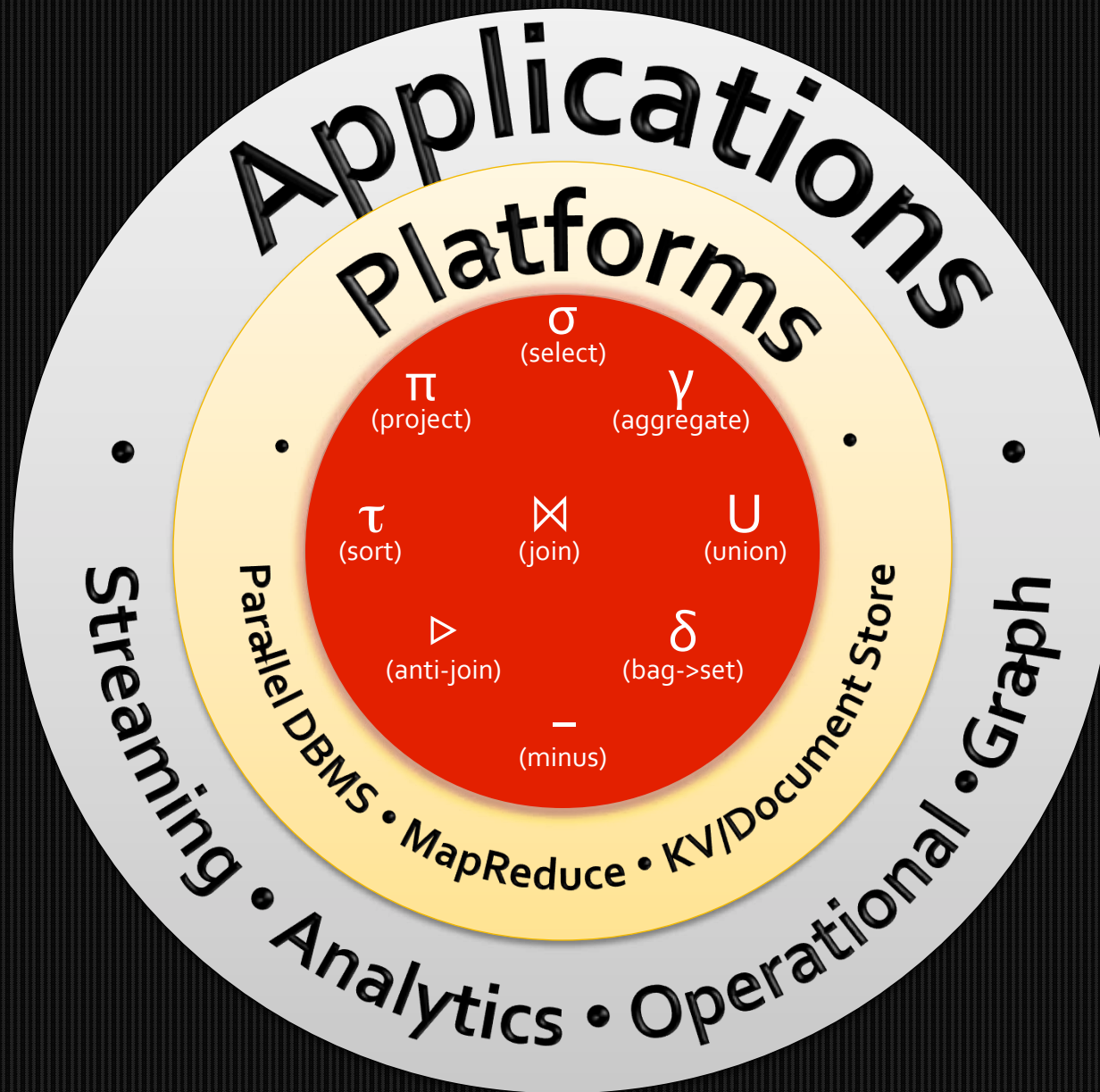# Data@ Bare Metal Speed

## Jignesh M. Patel

THE UNIVERSITY of WISCONSIN MADISON

*The Wisconsin Quickstep Project*

# Disruptive hardware trends

Want     Constraint

High Performance

Low Cost

Power

# TPC-H: Big-3 Vendors, 3TB scale



Data (40%/yr)

Performance/Core

Normalized Growth

20

10

0

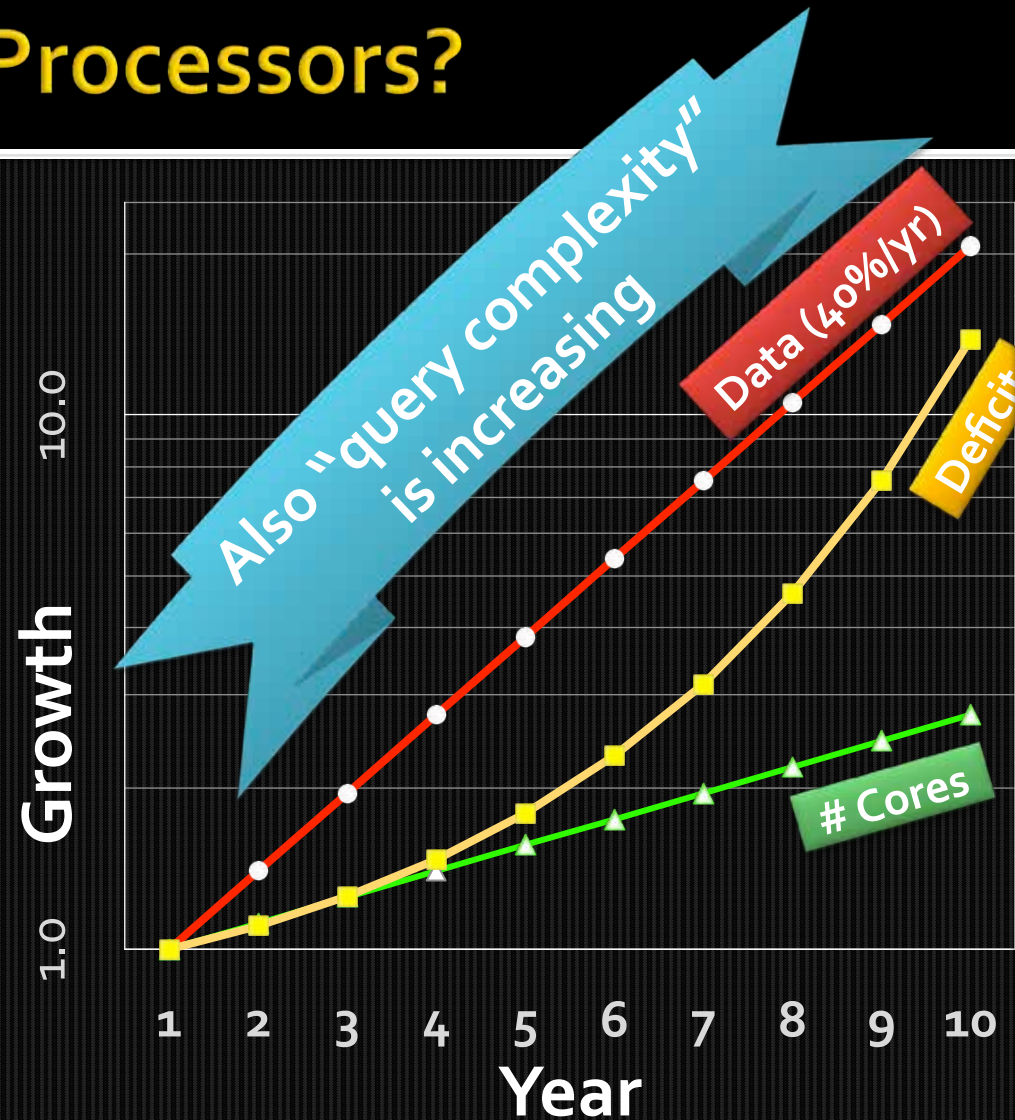Aug-03  Dec-04  May-06  Sep-07  Feb-09  Jun-10  Nov-11  Mar-13

# TPC-H: Growth driven by moving data to main memory

# What's Next for Processors?

- Future processor design?
  1. Keep adding cores (~40% per generation)
  2. Heterogeneous cores
  3. Programmable functional units
- But, systems must stay within a power budget
- Data growth continues unabated

*Also "query complexity" is increasing*

Growth

10.0

1.0

Data (40%/yr)

Deficit

# Cores

1  2  3  4  5  6  7  8  9  10

**Year**

## Need to do more with less.

# Quickstep

| Goal | • Run data analytics @ hardware speeds |
| --- | --- |
| Short-term | • Run @ the speed of hardware today |
| Long-term | • Hardware-software co-design for data kernels |

# Scan: A Key Data Processing Kernel

**What?**
- Scan a column of a table applying some predicate

**Why?**
- A key primitive in database
- "The" critical kernel in main memory analytic systems

**How?**
- Conserve memory bandwidth: **BitWeaving** the data
- Use every bit of data that is brought to the procecer efficiently using **intra-cycle parallelism**

# Focus on Column Scan (can be generalized)

## Traditional Row Store

| shipdate | ... | discount | quantity |
|---|---|---|---|
| Mar-12-2013 | | 5% | 5 |
| Jan-08-2013 | | 2% | 4 |
| Apr-29-2013 | | 10% | 3 |
| May-14-2013 | | 0% | 6 |
| ... | ... | ... | ... |
| Feb-28-2013 | | 5% | 0 |

One big file

## Column Store

| shipdate |
|---|
| Mar-12-2013 |
| Jan-08-2013 |
| Apr-29-2013 |
| May-14-2013 |
| ... |
| Feb-28-2013 |

File: 1

• • •

| discount |
|---|
| 5% |
| 2% |
| 10% |
| 0% |
| ... |
| 5% |

File: n-1

| quantity |
|---|
| 5 |
| 4 |
| 3 |
| 6 |
| ... |
| 0 |

File: n

16 bits

**Order-preserving compression**

Column Codes: | 5 | 4 | 3 | 6 | 2 | 7 | 1 | 0 | . . .

3 bits

BitWeaving/V

Code

First batch of Processor Words
*(batch size = code size in bits)*

Next batch of processor words

BitWeaving/H

| Code 1 | Code 5 | Code 9 | Code 13 | Code 17 | Code 21 | Code 25 | Code 29 |

| Code 2 | Code 6 | Code 10 | Code 14 | Code 18 | Code 22 | Code 26 | Code 30 |

| Code 3 | Code 7 | Code 11 | Code 15 | Code 19 | Code 23 | Code 27 | Code 31 |

| Code 4 | Code 8 | Code 12 | Code 16 | Code 20 | Code 24 | Code 28 | Code 32 |

First batch of Processor Words
*(batch size = code size in bits)*

Next batch of processor words

# Framework – Example

```
SELECT SUM(l_discount * l_price) FROM lineitem
WHERE l_shipdate BETWEEN Date AND Date + 1 year
    AND l_discount BETWEEN Discount – 0.01 AND Discount + 0.01
    AND l_quantity < Quantity
```



BitWeaving/H

l_price → Aggregation ← RID List: 9, 15 ← Result bit vector

l_discount

BitWeaving/H

Result bit vector

AND

Result bit vector

AND

Result bit vector

l_quantity

Result bit vector

l_shipdate

l_discount

BitWeaving/V

BitWeaving/H

BitWeaving/V

# BitWeaving/V

Column Codes:

| 10 | 12 | 3 | 6 | 9 | 7 | 1 | 0 |
|----|----|---|---|---|---|---|---|

| | 10 | 12 | 3 | 6 | 9 | 7 | 1 | 0 |
|---|----|----|---|---|---|---|---|---|
| Word 1 | 1 | 1 | 0 | 0 | 1 | 0 | 0 | 0 |
| Word 2 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 |
| Word 3 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 0 |
| Word 4 | 0 | 0 | 1 | 0 | 1 | 1 | 1 | 0 |

The first (most significant) bits of 8 consecutive codes

The second bits of 8 consecutive codes

The third bits of 8 consecutive codes

The last (least significant) bits of 8 consecutive codes

# BitWeaving/V - early pruning

**Column Codes:**

| 10 | 12 | 3 | 6 | 9 | 7 | 1 | 0 |
|----|----|---|---|---|---|---|---|

**Constant**

| 5 |
|---|

**Predicate**

| a < 5 |
|-------|

| 1 | 1 | 0 | 0 | 1 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 |
| 1 | 0 | 1 | 1 | 0 | 1 | 0 | 0 |
| 0 | 0 | 1 | 0 | 1 | 1 | 1 | 0 |

Constant bits: 0, 1, 0, 1

✘ ✘ ? ? ✘ ? ? ?

✘ ✘ ✔ ? ✘ ? ✔ ✔

**Result Bit Vector**

0 0 1 1 0 0 1 1

✘ ✘ ✔ ✔ ✘ ✘ ✔ ✔

**Early Pruning: terminate the predicate evaluation on a segment, when all results have been determined.**

# BitWeaving/V - Early Pruning Model



Early pruning probability: 96% at bit position 4

t size: 64 codes, code size: 32 bits

Early pruning probability: 98% at bit position 8

Early Pruning Probability P(b)

- Fill factor:100%
- Fill factor: 10%
- Fill factor:  1%

This cut-off mechanism allows for efficient evaluation of conjunction/disjunction predicates

# BitWeaving/H - Example

Code size (3 + 1 bits)

Segment 1

| c1 | c2 | c3 | c4 | c5 | c6 | c7 | c8 | c9 | c10 | c11 | c12 | c13 | c14 | c15 | c16 | c17 | c18 | c19 |
|----|----|----|----|----|----|----|----|----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|

Predicate evaluation is done on the 4 codes in parallel

Segment 1

Memory space

word 1    | c1 | c5 | c9 | c13 |    < 5   →   **?**

word 2    | c2 | c6 | c10 | c14 |

word 3    | c3 | c7 | c11 | c15 |

word 4    | c4 | c8 | c12 | c16 |

Word size
(16 bits)

# BitWeaving/H: Less Than Predicate

Uses only 3 instructions! Without the delimiter, we would need ~12 instructions...

$$X = \left( c_1 c_5 c_9 c_{13} \right)$$

$$Y = \left( 5555 \right)$$

$$\left( Y + \left( X \oplus M1 \right) \right) \wedge M2$$

$M1 = 0111\ 0111\ 0111\ 0111$
$M2 = 1000\ 1000\ 1000\ 1000$

|         | c5=7  | c9=6  | c13=2 |
|---------|-------|-------|-------|
| 0001    | 0111  | 0110  | 0010  |
| 0101    | 0101  | 0101  | 0101  |
| 1000    | 0000  | 0000  | 1000  |

**Works for arbitrary code sizes & word sizes!**

# BitWeaving/H - Example

Code size (3 + 1 bits)

Segment 1

| c1 | c2 | c3 | c4 | c5 | c6 | c7 | c8 | c9 | c10 | c11 | c12 | c13 | c14 | c15 | c16 | c17 | c18 | c19 |
|----|----|----|----|----|----|----|----|----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|

Segment 1

Memory space

word 1

| c1 | c5 | c9 | c13 |
|----|----|----|-----|

< 5 ?  →  1000100000001000  >>0

c1  1000100000001000

word 2

| c2 | c6 | c10 | c14 |
|----|----|-----|-----|

< 5 ?  →  1000000000000000  >>1

c2  0100000000000000

word 3

| c3 | c7 | c11 | c15 |
|----|----|-----|-----|

< 5 ?  →  0000100010001000  >>2

c3  0000001000100010

word 4

| c4 | c8 | c12 | c16 |
|----|----|-----|-----|

< 5 ?  →  1000000000001000  >>3

c4  0001000000000001

$\sum v$

Word size (16 bits)

**Result bit vector computed efficiently with this layout!**

1101101000101011

# Evaluation

## SYSTEM

- Intel Xeon X5650
  - 64 bits ALU
  - 128 bits SIMD
  - 12MB L3 Cache
- 24GB memory
- Single threaded execution

## WORKLOAD

1. Synthetic
   - ```
     SELECT COUNT(*)
     FROM R
     WHERE R.a < C
     ```
   - 1 billion tuples
   - Uniform distribution
   - Selectivity: 10%
2. TPC-H @ SF=10
   - scan only with materialized join results

# Evaluation: Micro-benchmark

# Evaluation: Micro-benchmark



2X: SIMD parallelism

Cycles / code

Size of code (# bits)

Naive

SIMD

# Evaluation: Micro-benchmark

# Evaluation: Micro-benchmark



3X-4X speedup over BL:
1) Use the extra (delimiter) bit
2) Easy to produce the result bit vector with the HBP layout

Legend: Naive, SIMD, BL, BitWeaving/H

X-axis: Size of code (# bits)
Y-axis: Cycles / code

# Evaluation: Micro-benchmark



Chart axes: Cycles / code (y-axis, 0 to 10) vs Size of code (# bits) (x-axis, 0 to 32)

Legend:
- Naive
- SIMD
- BL
- BitWeaving/H
- BitWeaving/V

2X speedup: Early pruning

Many more experiments in the paper

# WideTable

## Customer

| cid | cname | gender | address |
|-----|-------|--------|---------|
| 1 | Andy | M | 100 Main st. |
| 2 | Kate | F | 20 10$^{th}$ blvd. |
| 3 | Bob | M | 300 5$^{th}$ ave. |

## Product

| pid | pname |
|-----|-------|
| 1 | Milk |
| 2 | Coffee |
| 3 | Tea |

## Buy

| cid | pid | status |
|-----|-----|--------|
| 1 | 2 | S |
| 2 | 2 | F |
| 3 | 3 | S |
| 1 | 2 | S |

| cid | cname | gender | address | pid | pname | status |
|-----|-------|--------|---------|-----|-------|--------|
| 1 | Andy | M | 100 Main st. | 2 | Coffee | S |
| 2 | Kate | F | 20 10$^{th}$ blvd. | 2 | Coffee | F |
| 3 | Bob | M | 300 5$^{th}$ ave. | 3 | Tea | S |
| 1 | Andy | M | 100 Main st. | 2 | Coffee | S |
| NULL | NULL | NULL | NULL | 1 | Milk | NULL |

**WideTable**

# WideTable

WideTable

Denormalization

Column-store

Packed Scan

Dictionary encoding

**Now we can run analytical workloads (e.g. TPC-H) using simple BitWeaved scans**

# WideTable Techniques

**Quickstep/WT**

Data

Schema Graph

**Denormalizer**
- Schema transformer
- Flatten using $\bowtie$

**Query Transformer**

$$R \ltimes S \equiv \pi_{R.*}(R \bowtie S)$$

$$R \bowtie \ldots \bowtie S \equiv \sigma_{R.p \neq NULL \wedge \ldots \wedge S.p \neq NULL}(R \bowtie \ldots \bowtie S)$$

$$\gamma_{\ldots}(R) \bowtie S \equiv \gamma_{\ldots, S.*}(R \bowtie S)$$

$$R \equiv \pi_{R.*}(R \bowtie S)$$

$$\pi_{a_1}(R) \equiv \pi_{a_1}(\ldots(\pi_{a_n}(R))), \text{ where } a_1 \subseteq \ldots \subseteq a_n$$

**Query**

Bit-Weaved Scans

**{WideTables}**

Results

# Schema Graph



**Schema Graph** → **Schema Tree**

**WideTable = (Region ⋈ Nation ⋈ Customer) ⋈ (Region ⋈ Nation ⋈ Product ⋈ Buy)**

SMW = {WideTables}
e.g. for TPC-H, SMW={lineItemWT, ordersWT, partsuppWT, customerW

# TPC-H Queries

| TPC-H Queries | Joins | Nested Queries | Non-FK joins | WideTable |
|---|---|---|---|---|
| Q1, Q6 | | | | LineitemWT |
| Q3, Q5, Q7-Q10, Q12, Q14, Q19 | × | | | LineitemWT |
| Q4, Q15, Q17, Q18, Q20 | × | × | | LineitemWT |
| Q21 | × | × | × | --- |
| Q2, Q11, Q16 | × | × | | PartsuppWT |
| Q13 | × | | | OrdersWT |
| Q22 | × | × | | OrdersWT |

# Evaluation

## SYSTEM

- Intel Xeon E5-2620 × 2

- 2.0 GHz

- 12 cores / 24 threads

- 15MB L3 Cache

- 32G. 1600MHz DDR3

## BENCHMARK

- SF: 10 (~10GB)

- SMW =

| | |
|---|---|
| lineItemWT | 5.4 GB |
| ordersWT | 0.7 GB |
| partsuppWT | 0.2 GB |
| customerWT | 0.05 GB |
| dictionaries | 0.8 GB |
| filter columns | 1.3 GB |
| TOTAL | 8.5GB |

# Speedup over MonetDB: Single Thread

**WideTable over 10X faster than MonetDB for about half of the 21 queries**

# Speedup over MonetDB: 12 Threads

WideTable scales better



Speedup over MonetDB

| Query | Value |
|-------|-------|
| Q1 | 2.1 |
| Q2 | 6.3 |
| Q3 | 29.6 |
| Q4 | 8.7 |
| Q5 | 37.9 |
| Q6 | 2.2 |
| Q7 | 20.3 |
| Q8 | 35.1 |
| Q9 | 10.9 |
| Q10 | 7.3 |
| Q11 | 5.3 |
| Q12 | 21.2 |
| Q13 | 3.2 |
| Q14 | 3.0 |
| Q15 | 11.6 |
| Q16 | 18.3 |
| Q17 | … to 182 |
| Q18 | 1.0 |
| Q19 | 15.0 |
| Q20 | 2.9 |
| Q22 | |

# Conclusions and Future Work

Transformative architectural changes at all levels (CPU, memory subsystem, I/O subsystem) is underway

Need to rethink data processing kernels
- Run @ current bare metal speed

Need to think of hardware software co-design

**Big Data Hardware**

**Big Data Software**

# Thanks!



| Craig Chasseur | Anusha Dasarakothapalli | Harshad Deshmukh | Jing Fan | Yinan Li | James Paton | Navneet Potti | Sangmin Shin | Qiang Zeng |

**The Quickstep Team**



**Funding**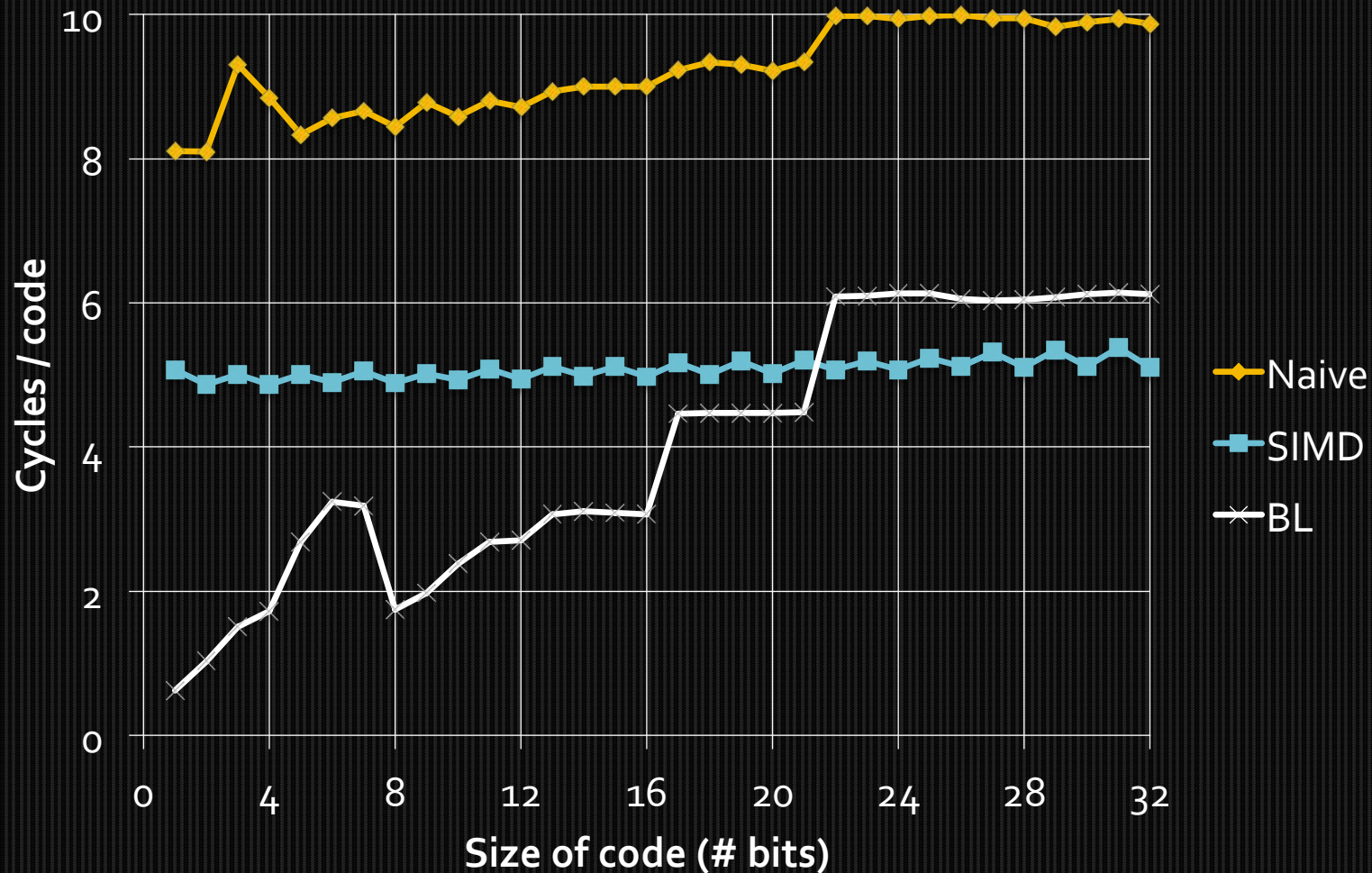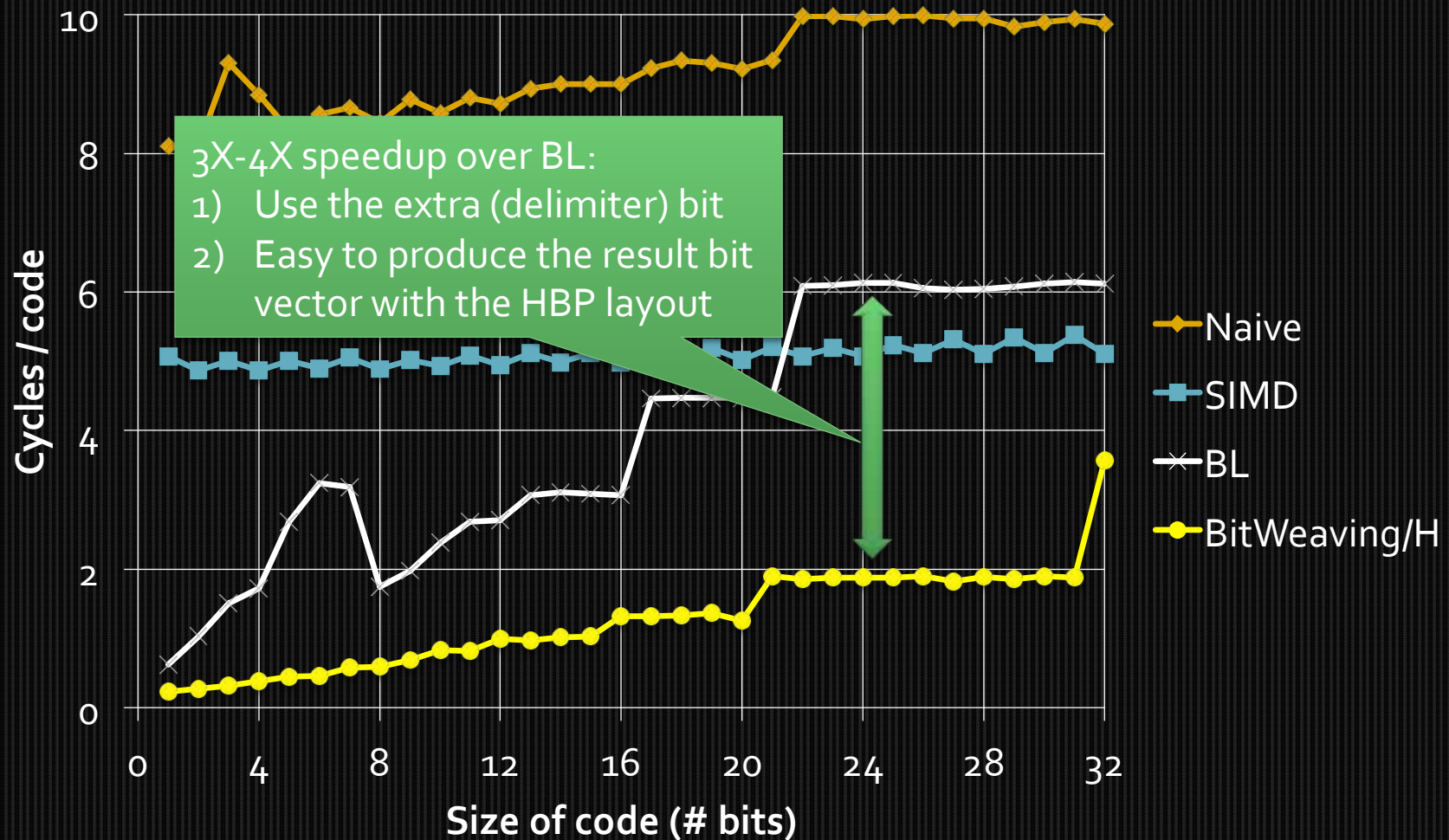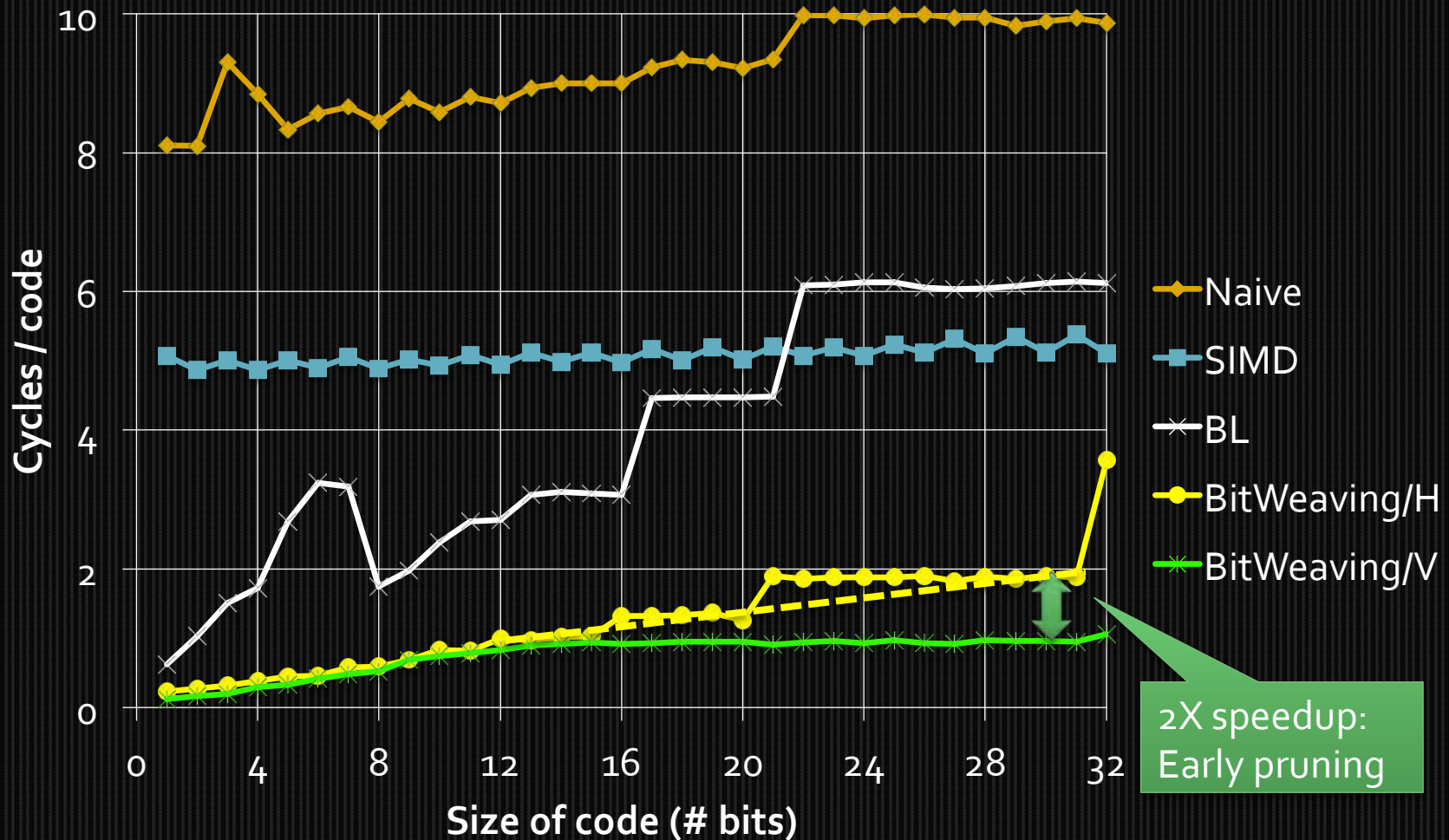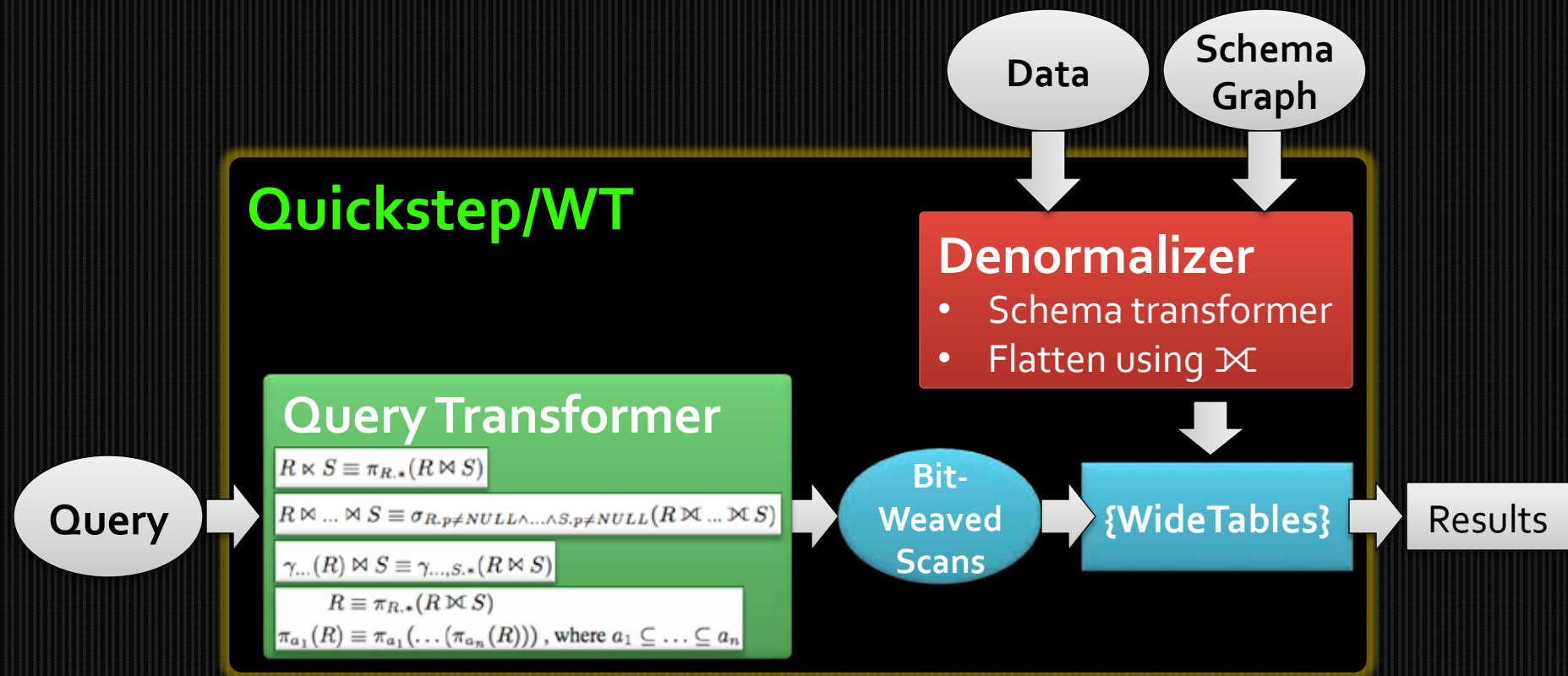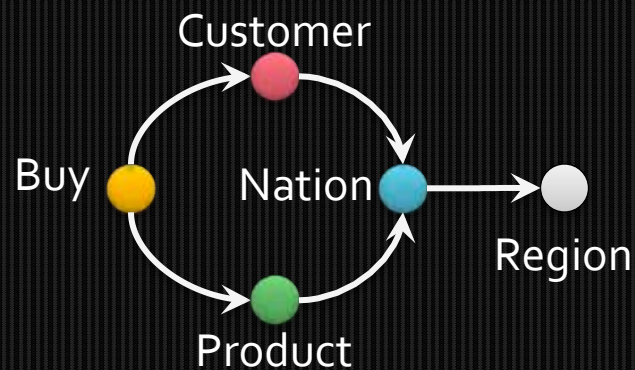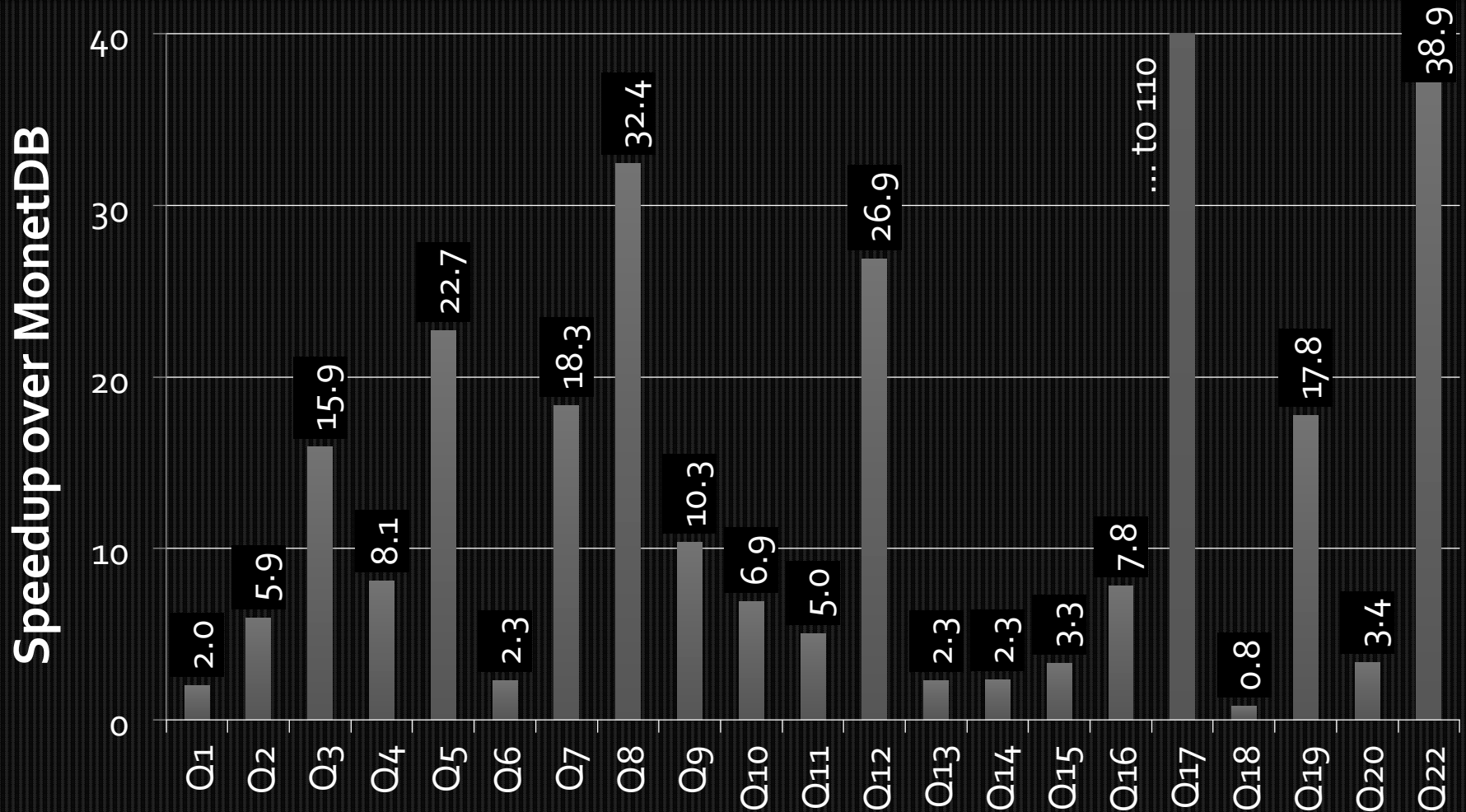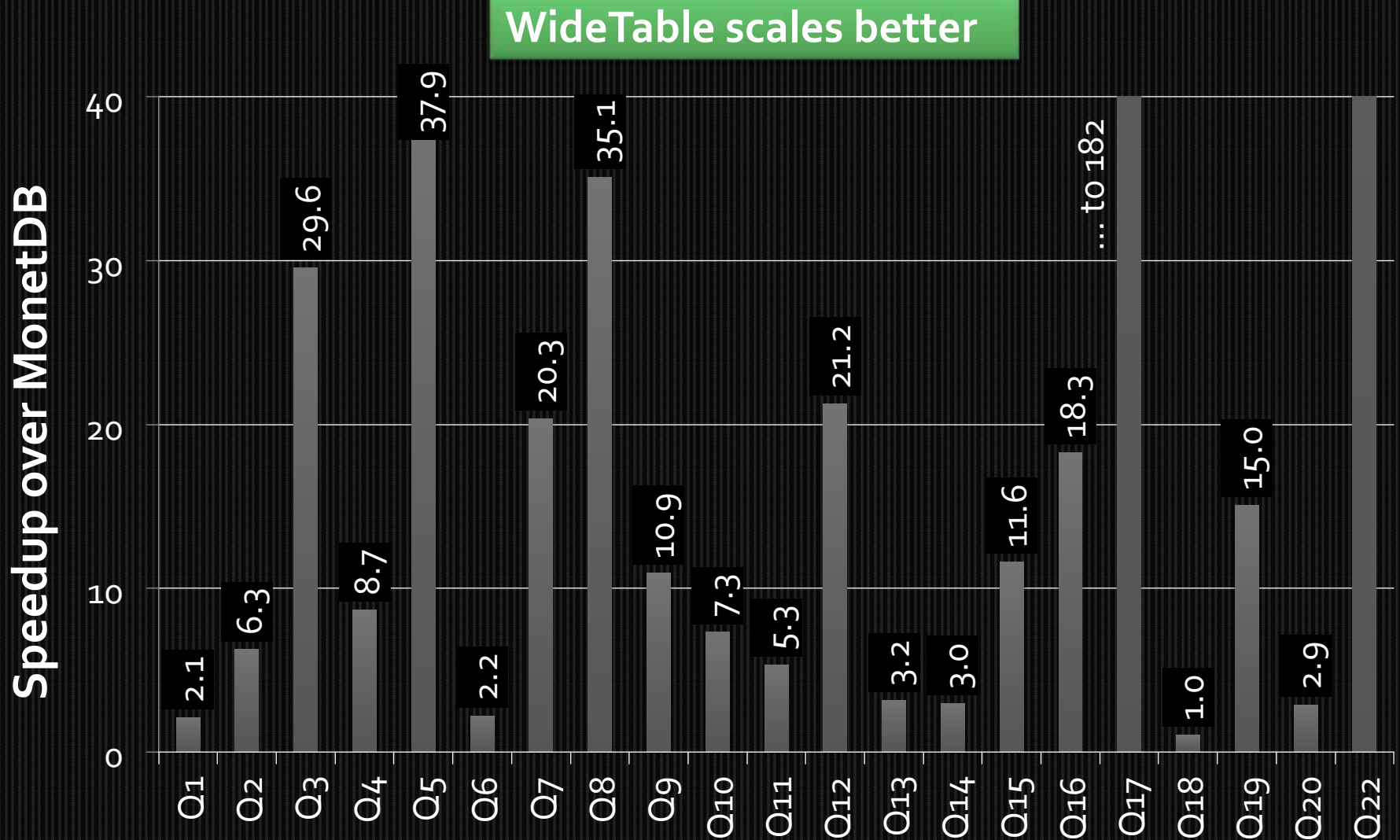