# **A**dvanced **D**ata **M**anagement for Data Analysis

## *Stefan Manegold*

Data Management @ LIACS

Group leader Database Architectures
Centrum Wiskunde & Informatica (CWI)
Amsterdam

s.manegold@liacs.leidenuniv.nl
http://www.cwi.nl/~manegold/

# ADM: Agenda

- 07.09.2022: Lecture  1: **Introduction**

- 14.09.2022: Lecture  2: **SQL Recap**

    *(plus Assignment 1 [in groups; 3 weeks]: TPC-H benchmark)*

- 21.09.2022: Lecture  3: **Column-Oriented Database Systems (1/6) - Motivation & Basic Concepts**

- 28.09.2022: Lecture  4: **Column-Oriented Database Systems (2a/6) - Selected Execution Techniques (1/2)**

- 05.10.2022: Lecture  5: **Column-Oriented Database Systems (2b/6) - Selected Execution Techniques (2/2)**

    *(plus Assignment 2 [in groups; 4 weeks]: Compression techniques)*

- 12.10.2022: Lecture  6: **Column-Oriented Database Systems (3/6) - Cache Conscious Joins**

- 19.10.2022: Lecture  7: **Column-Oriented Database Systems (4/6) - "Vectorized Execution"**

- 26.10.2022: ***No lecture!***

- 02.11.2022: Lecture  8: **DuckDB: An embedded database for data science (1/2) (guest lecture & *hands-on*)**

    *(plus Assignment 3 [individual; 2 weeks]: Analysing NYC Cab dataset with DuckDB)*

- 09.11.2022: Lecture  9: **DuckDB: An embedded database for data science (2/2) (guest lecture & *hands-on*)**

- 16.11.2022: Lecture 10: **Branch Misprediction & Predication**

    *(plus Assignment 4 [individual; 2 weeks]: Predication)*

- 23.11.2022: Lecture 11: **Column-Oriented Database Systems (5/6) - Adaptive Indexing**

- 30.11.2022: Lecture 12: **Column-Oriented Database Systems (6/6) - Progressive Indexing**
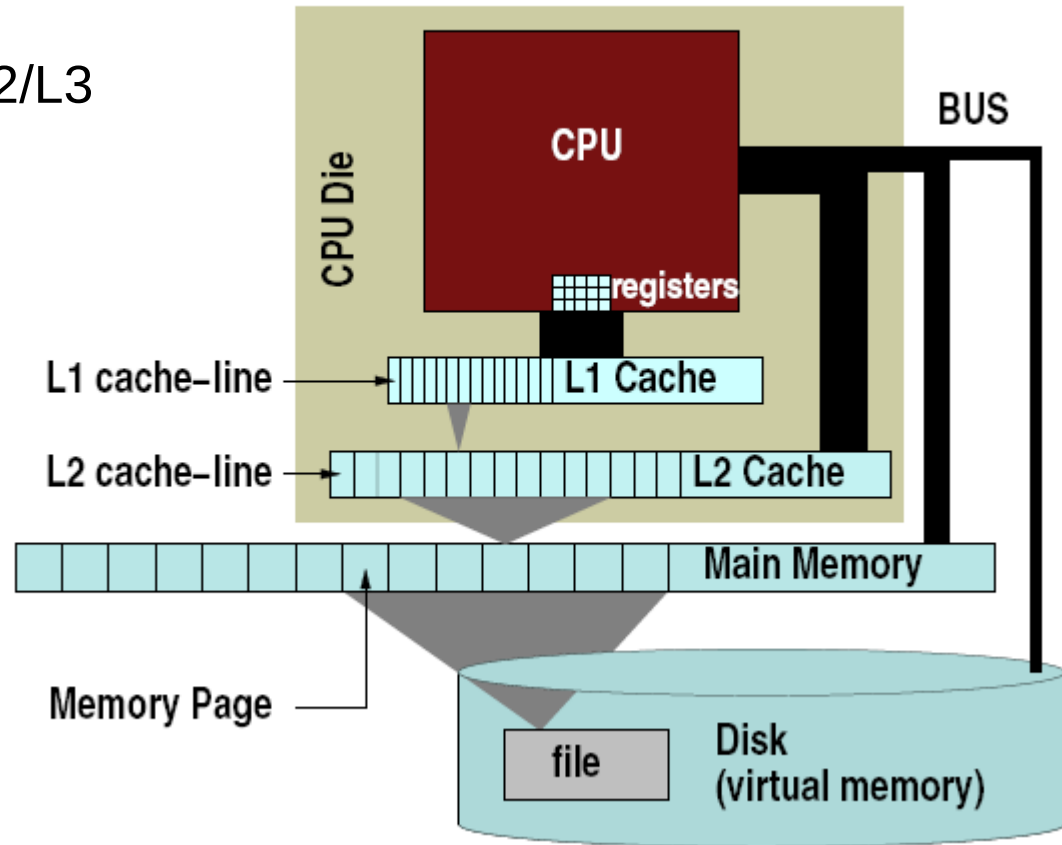
# ADM: Literature

- **Column-Oriented Database Systems (4/6) - "Vectorized Execution"**

  - "MonetDB/X100: Hyper-Pipelining Query Execution". Boncz, Zukowski, Nes. CIDR'05.

  - "Buffering Database Operations for Enhanced Instruction Cache Performance". Zhou and Ross. SIGMOD'04.

  - "Block oriented processing of relational database operations in modern computer architectures". Padmanabhan, Malkemus, Agarwal. ICDE'01.

  - "Balancing Vectorized Query Execution with Bandwidth Optimized Storage". Zukowski. PhD Thesis. CWI 2008.

# CPU Architecture

Elements:

- Storage
  - CPU caches L1/L2/L3
- Registers
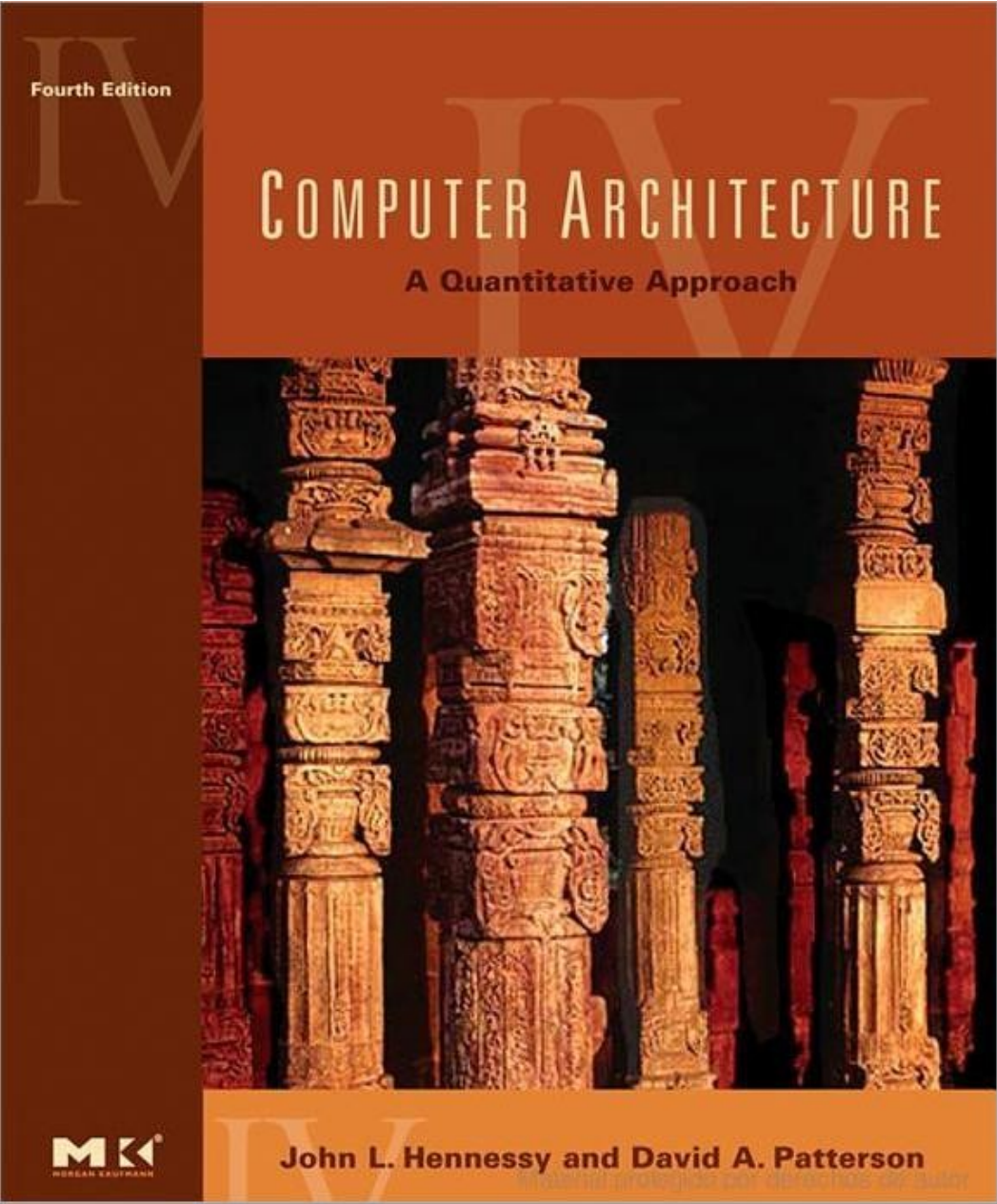- Execution Unit(s)
  - Pipelined
  - SIMD

# CPU Metrics

| Processor | 16-bit address/, bus, micro-coded | 32-bit address/ bus, micro-coded | 5-stage pipeline, on-chip I&D caches FPU | 2-way super-scalar, 64-bit bus | Out-of-order, 3-way super-scalar | Out-of-order, super-pipelined, on-chip L2 cache | Multi-core |
|---|---|---|---|---|---|---|---|
| Product | 80286 | 80386 | 80486 | Pentium | PentiumPro | Pentium4 | CoreDuo |
| Year | 1982 | 1985 | 1989 | 1993 | 1997 | 2001 | 2006 |
| Transistors (thousands) | 134 | 275 | 1,200 | 3,100 | 5,500 | 42,000 | 151,600 |
| Latency (clocks) | 6 | 5 | 5 | 5 | 10 | 22 | 12 |
| Bus width (bits) | 16 | 32 | 32 | 64 | 64 | 64 | 64 |
| Clock rate (MHz) | 12.5 | 16 | 25 | 66 | 200 | 1500 | 2333 |
| Bandwidth (MIPS) | 2 | 6 | 25 | 132 | 600 | 4500 | 21000 |
| Latency (ns) | 320 | 313 | 200 | 76 | 50 | 15 | 5 |

# CPU Metrics

| Processor | 16-bit address/, bus, micro-coded | 32... add... b... mi... co... | core |
|---|---|---|---|
| Product | 80286 | 80... | Duo |
| Year | 1982 | 1... | 6 |
| Transistors (thousands) | 134 | 2... | 600 |
| Latency (clocks) | 6 | | |
| Bus width (bits) | 16 | 3... | |
| Clock rate (MHz) | 12.5 | | 3 |
| Bandwidth (MIPS) | 2 | | 00 |
| Latency (ns) | 320 | 3... | |



Fourth Edition

COMPUTER ARCHITECTURE

A Quantitative Approach

John L. Hennessy and David A. Patterson

# DRAM Metrics

# Super-Scalar Execution (pipelining)



## Sequential execution

| | | | |
|---|---|---|---|
| Instruction fetch | IF-1 | IF-2 | IF-3 |
| Instruction decode | ID-1 | ID-2 | ID-3 |
| Execute | EX-1 | IE-2 | EX-3 |
| Write back | WB-1 | WB-2 | WB-3 |

CPU cycle

## Pipelined execution

| Instruction fetch | IF-1 | IF-2 | IF-3 | IF-4 | IF-5 | IF-6 |
| Instruction decode | | ID-1 | ID-2 | ID-3 | ID-4 | ID-5 |
| Execute | | | EX-1 | EX-2 | EX-3 | EX-4 |
| Write back | | | | WB-1 | WB-2 | WB-3 |

CPU cycle

Time

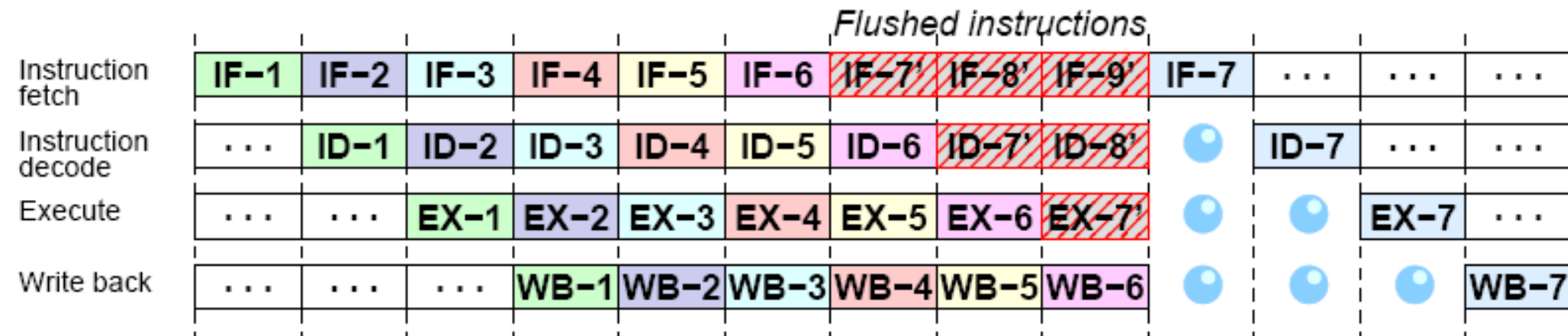(See also https://en.wikipedia.org/wiki/Instruction_pipelining )

# Hazards

- Data hazards
  - Dependencies between instructions
  - L1 data cache misses

- Control Hazards
  - Branch mispredictions
  - Computed branches (late binding)
  - L1 instruction cache misses

Result: bubbles in the pipeline



*Flushed instructions*

| | | | | | | | | Flushed instructions | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Instruction fetch | IF-1 | IF-2 | IF-3 | IF-4 | IF-5 | IF-6 | IF-7' | IF-8' | IF-9' | IF-7 | . . . | . . . | . . . |
| Instruction decode | . . . | ID-1 | ID-2 | ID-3 | ID-4 | ID-5 | ID-6 | ID-7' | ID-8' | | ID-7 | . . . | . . . |
| Execute | . . . | . . . | EX-1 | EX-2 | EX-3 | EX-4 | EX-5 | EX-6 | EX-7' | | | EX-7 | . . . |
| Write back | . . . | . . . | . . . | WB-1 | WB-2 | WB-3 | WB-4 | WB-5 | WB-6 | | | | WB-7 |

Out-of-order execution addresses data hazards
- control hazards typically more expensive

(See also https://en.wikipedia.org/wiki/Instruction_pipelining )

# Database Architecture causes Hazards

DB workload execution on a modern computer
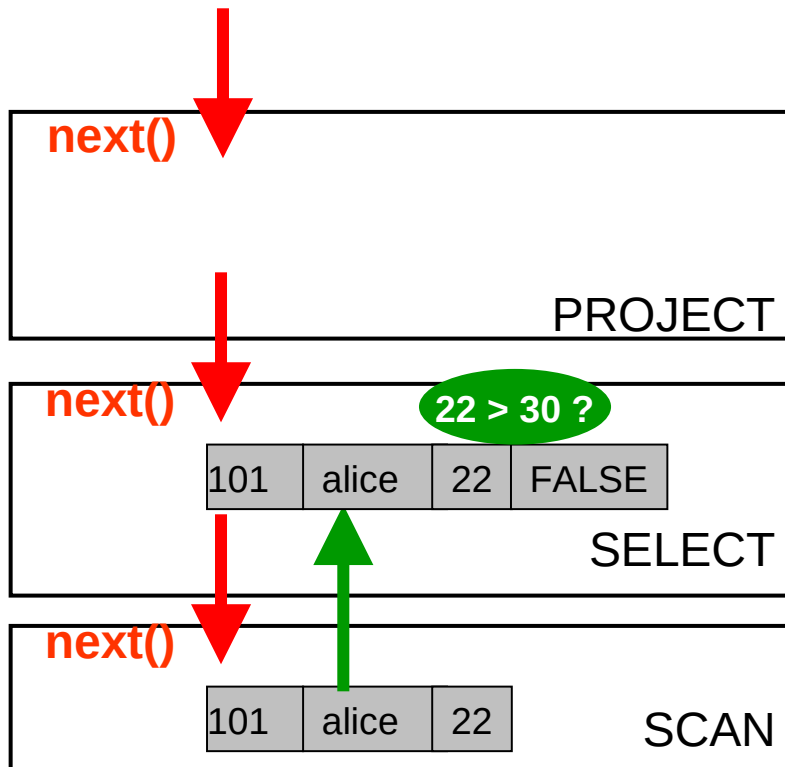
Processor — BUSY — IDLE (STALLED)



"DBMSs On A Modern Processor: Where Does Time Go? "
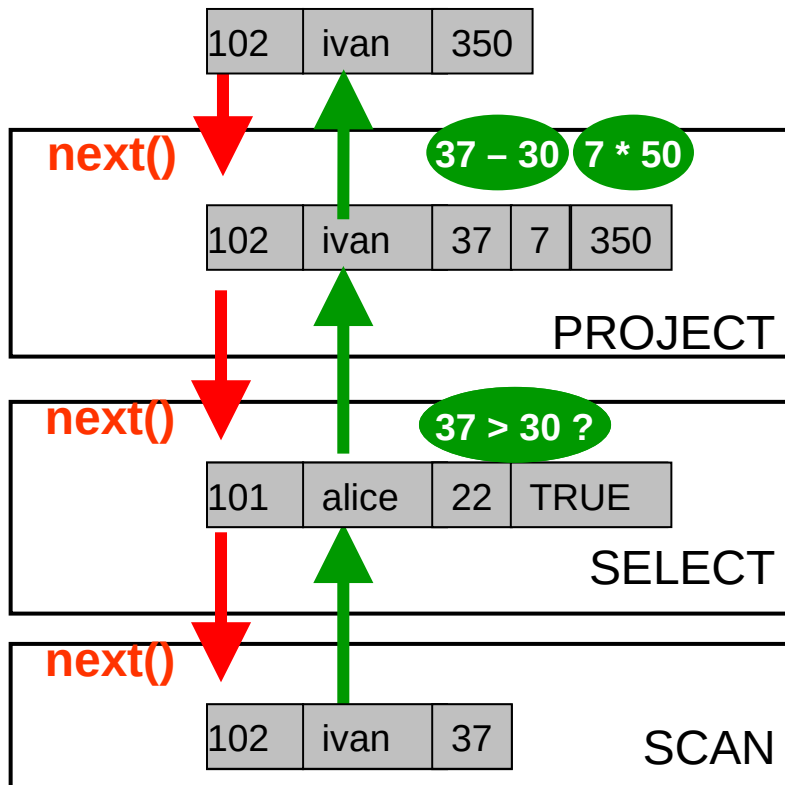Ailamaki, DeWitt, Hill, Wood, VLDB'99

# SIMD



- Single Instruction Multiple Data
  - Same operation applied on a vector of values
  - MMX: 64 bits, SSE: 128bits, AVX: 256bits
  - SSE, e.g. multiply 8 short integers
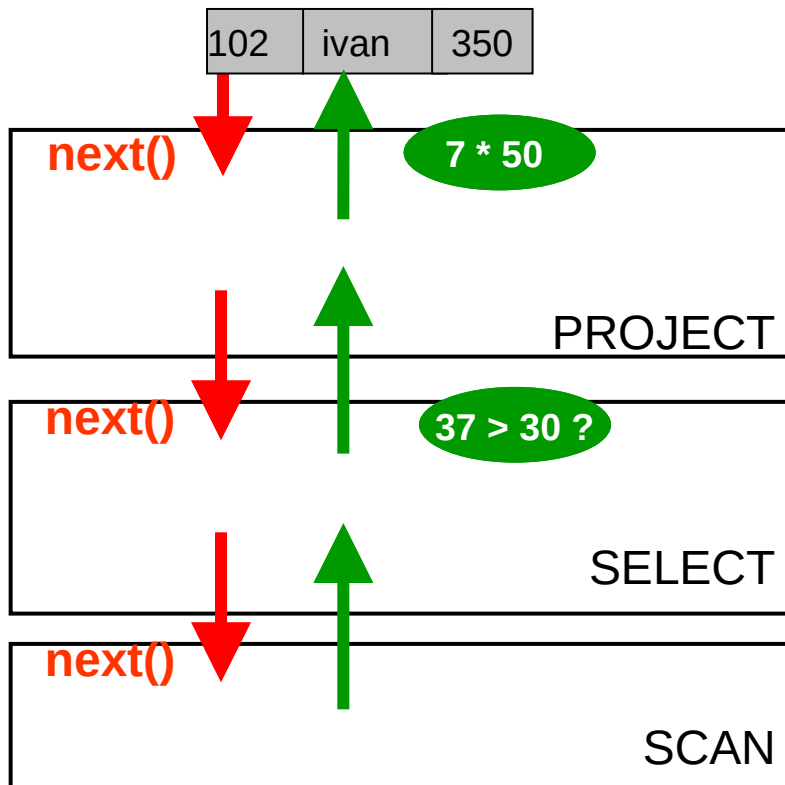
# A Look at the Query Pipeline

**next()**

PROJECT

**next()**

22 > 30 ?

| 101 | alice | 22 | FALSE |

SELECT

**next()**

| 101 | alice | 22 |

SCAN

**SELECT** id, name
(age-30)*50 AS bonus
**FROM** employee
**WHERE** age > 30

# A Look at the Query Pipeline

| 102 | ivan | 350 |
|-----|------|-----|

**next()**     ( 37 – 30 ) ( 7 * 50 )

| 102 | ivan | 37 | 7 | 350 |
|-----|------|----|----|-----|

PROJECT

**next()**     ( 37 > 30 ? )

| 101 | alice | 22 | TRUE |
|-----|-------|----|------|

SELECT

**next()**

| 102 | ivan | 37 |
|-----|------|----|

SCAN

**SELECT   id, name**
              **(age-30)*50 AS bonus**
**FROM      employee**
**WHERE   age > 30**

# A Look at the Query Pipeline

| 102 | ivan | 350 |
|-----|------|-----|

**next()**    7 * 50

PROJECT

**next()**    37 > 30 ?

SELECT

**next()**

SCAN

## Operators

Iterator interface
-open()
**-next():** tuple
-close()

# A Look at the Query Pipeline

| 102 | ivan | 350 |

**next()**    ( 7 * 50 )

| 102 | ivan | 37 | 7 | 350 |

PROJECT

**next()**    ( 37 > 30 ? )

| 101 | alice | 22 | TRUE |

SELECT

**next()**

| 102 | ivan | 37 |    SCAN

# Primitives

Provide computational functionality

All arithmetic allowed in expressions,
e.g. Multiplication

( 7 * 50 )

`mult(int,int)` ➔ `int`

# Database Architecture causes Hazards

- Data hazards
  - Dependencies between instructions
  - L1 data cache misses

- Control Hazards
  - Branch mispredictions
  - Computed branches (late binding)
  - L1 instruction cache misses

SIMD    Out-of-order Execution

work on one tuple at a time

Large Tree/Hash Structures

Code footprint of all operators in query plan exceeds L1 cache

Data-dependent conditions

Next() late binding method calls

Tree, List, Hash traversal    Complex NSM record navigation

PROJECT

SELECT

SCAN

## Operators

Iterator interface
-open()
-next(): tuple
-close()

# TPC-H Database Schema

**PART**
- p_partkey
- p_name
- p_mfgr
- p_brand
- p_type
- p_size
- p_container
- p_retailprice
- p_comment

**PARTSUPP**
- ps_partkey
- ps_suppkey
- ps_availqty
- ps_supplycost
- ps_comment

**SUPPLIER**
- s_suppkey
- s_name
- s_address
- s_nationkey
- s_phone
- s_acctbal
- s_comment

**LINEITEM**
- l_orderkey
- l_linenumber
- l_partkey
- l_suppkey
- l_quantity
- l_extendedprice
- l_discount
- l_tax
- l_returnflag
- l_linestatus
- l_shipdate
- l_commitdate
- l_receiptdate
- l_shipinstructions
- l_shipmode
- l_comment

**ORDERS**
- o_orderkey
- o_custkey
- o_orderstatus
- o_totalprice
- o_orderdate
- o_orderpriority
- o_clerk
- o_shippriority
- o_comment

**CUSTOMER**
- c_custkey
- c_name
- c_address
- c_nationkey
- c_phone
- c_acctbal
- c_clerk
- c_mktsegment
- c_comment

**NATION**
- n_nationkey
- n_name
- n_regionkey
- n_comment

**REGION**
- r_regionkey
- r_name
- r_comment

# TPC-H Query 1

```
select
        l_returnflag,
        l_linestatus,
        sum(l_quantity) as sum_qty,
        sum(l_extendedprice) as sum_base_price,
        sum(l_extendedprice * (1 - l_discount)) as sum_disc_price,
        sum(l_extendedprice * (1 - l_discount) * (1 + l_tax)) as sum_charge,
        avg(l_quantity) as avg_qty,
        avg(l_extendedprice) as avg_price,
        avg(l_discount) as avg_disc,
        count(*) as count_order
from
        lineitem
where
        l_shipdate <= date '1998-12-01' - interval '90' day (3)
group by
        l_returnflag,
        l_linestatus
order by
        l_returnflag,
        l_linestatus;
```

# DBMS Computational Efficiency

TPC-H 1GB, query 1

- selects 98% of fact table, computes net prices and aggregates all

- Results:
  - C program:      ?
  - MySQL:          26.2s
  - DBMS "X":       28.1s

"MonetDB/X100: Hyper-Pipelining Query Execution " Boncz, Zukowski, Nes, CIDR'05

# DBMS Computational Efficiency

TPC-H 1GB, query 1

- selects 98% of fact table, computes net prices and aggregates all
- Results:
  - C program:     **0.2s**
  - MySQL:        26.2s
  - DBMS "X":     28.1s

"MonetDB/X100: Hyper-Pipelining Query Execution " Boncz, Zukowski, Nes, CIDR'05
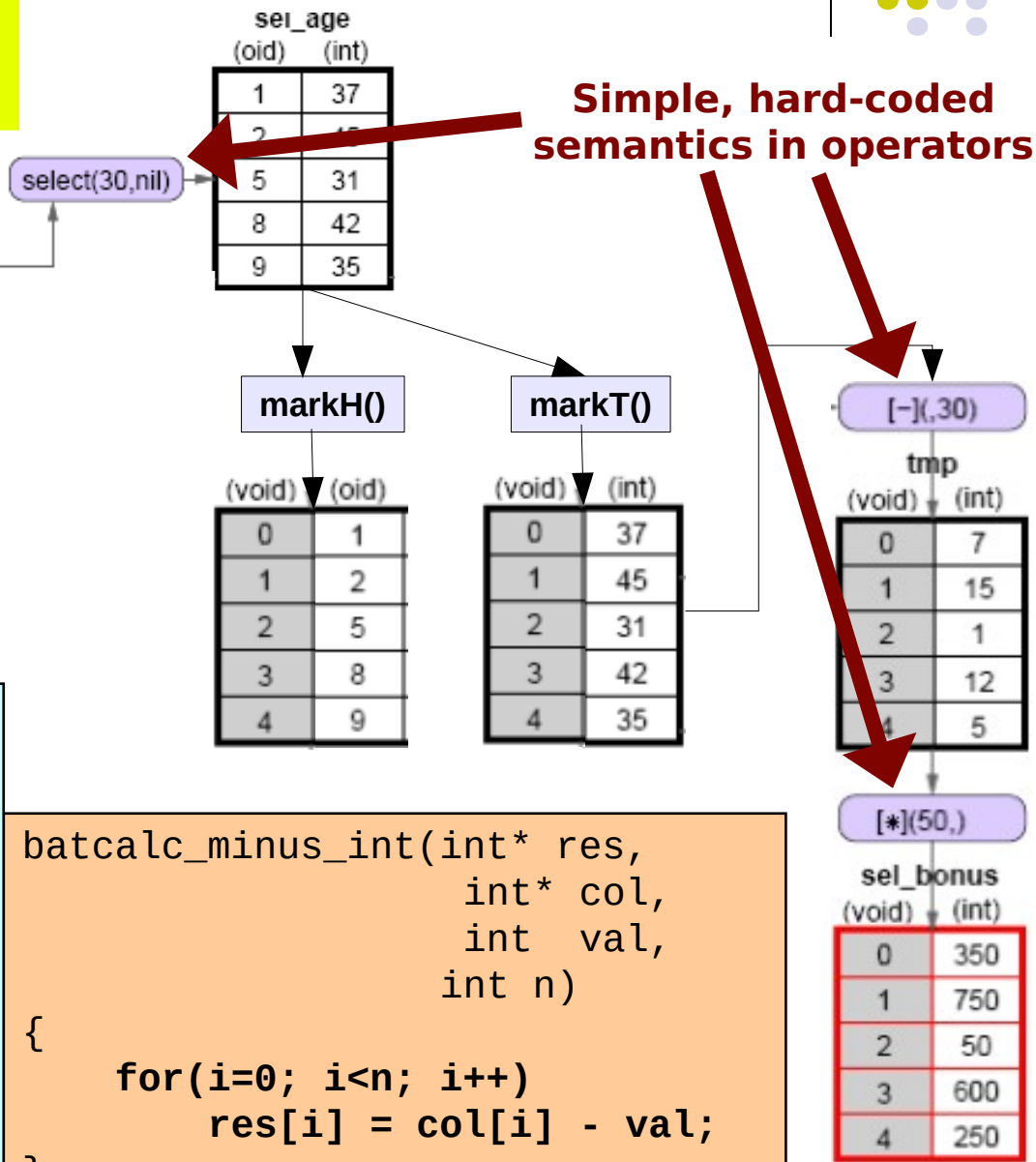
# a column-store

- ~~"save disk I/O when scan-intensive queries need a few columns"~~

- "avoid an expression interpreter to improve computational efficiency"

# RISC Relational Algebra

```
SELECT    id, name, (age-30)*50 as bonus
FROM      people
WHERE     age > 30
```

**Simple, hard-coded semantics in operators**

people_id

| (void) | (int) |
|--------|-------|
| 0 | 101 |
| 1 | 102 |
| 2 | 104 |
| 3 | 105 |
| 4 | 108 |
| 5 | 109 |
| 6 | 112 |
| 7 | 113 |
| 8 | 114 |
| 9 | 115 |

people_name

| (void) | (str) |
|--------|-------|
| 0 | Alice |
| 1 | Ivan |
| 2 | Peggy |
| 3 | Victor |
| 4 | Eve |
| 5 | Walter |
| 6 | Trudy |
| 7 | Bob |
| 8 | Zoe |
| 9 | Charlie |

people_age

| (void) | (int) |
|--------|-------|
| 0 | 22 |
| 1 | 37 |
| 2 | 45 |
| 3 | 25 |
| 4 | 19 |
| 5 | 31 |
| 6 | 27 |
| 7 | 29 |
| 8 | 42 |
| 9 | 35 |

select(30,nil)

sel_age

| (oid) | (int) |
|-------|-------|
| 1 | 37 |
| 2 | |
| 5 | 31 |
| 8 | 42 |
| 9 | 35 |

**markH()**

| (void) | (oid) |
|--------|-------|
| 0 | 1 |
| 1 | 2 |
| 2 | 5 |
| 3 | 8 |
| 4 | 9 |

**markT()**

| (void) | (int) |
|--------|-------|
| 0 | 37 |
| 1 | 45 |
| 2 | 31 |
| 3 | 42 |
| 4 | 35 |

[−](,30)

tmp

| (void) | (int) |
|--------|-------|
| 0 | 7 |
| 1 | 15 |
| 2 | 1 |
| 3 | 12 |
| 4 | 5 |

[*](50,)

sel_bonus

| (void) | (int) |
|--------|-------|
| 0 | 350 |
| 1 | 750 |
| 2 | 50 |
| 3 | 600 |
| 4 | 250 |

## CPU ☺?  Give it "nice" code !

- few dependencies (control,data)
- CPU gets out-of-order execution
- compiler can e.g. generate SIMD

**One loop for an entire column**
- no per-tuple interpretation
- arrays: no record navigation
- better instruction cache locality

```
batcalc_minus_int(int* res,
                  int* col,
                  int  val,
                  int n)
{
    for(i=0; i<n; i++)
        res[i] = col[i] - val;
}
```

# monetdb a column-store

- ~~"save disk I/O when select-intensive queries – need a few columns"~~
- "avoid an expression interpreter to improve computational efficiency"

How?

- RISC query algebra: hard-coded semantics
  - Decompose complex expressions in multiple operations
- Operators only handle **simple arrays**
  - No code that handles slotted buffered record layout
- Relational algebra becomes **array manipulation language**
  - Often SIMD for free

  - Plus: use of *cache-conscious* algorithms for Sort/Aggr/Join

# a Faustian pact

- You want efficiency
    - Simple hard-coded operators
- I take scalability
    - Result materialization



| | |
|---|---|
| C program: | 0.2s |
| MonetDB: | 3.7s |
| MySQL: | 26.2s |
| DBMS "X": | 28.1s |

# Pipelining vs Materialization

102 | ivan | 350

**next()**   7 * 50

102 | ivan | 37 | 7 | 350

PROJECT

**next()**   37 > 30 ?

101 | alice | 22 | TRUE

SELECT

**next()**

102 | ivan | 37   SCAN



**MATERIALIZED**
**intermediate**
**results**

# Pipelining vs Materialization



| 102 | ivan | 350 | | |

**next()**  **7 * 50**

| 102 | ivan | 37 | 7 | 350 |

PROJECT

**next()**  **37 > 30 ?**

| 101 | alice | 22 | TRUE |

SELECT

**next()**

| 102 | ivan | 37 | SCAN |

sel_age
(oid)   (int)

| 1 | 37 |
| 2 | 45 |
| 5 | 31 |
| 8 | 42 |
| 9 | 35 |

select(30,NII)    [−](,30)

tmp
(void)   (int)

| 0 | 7 |
| 1 | 15 |
| 2 | 1 |
| 3 | 12 |
| 4 | 5 |

**MATERIALIZED intermediate results**

=>    **vectorized query processing**

**MonetDB spin-off: MonetDB/X100**  vectorwise

**Observations:**

next() called much less often ➔ more time spent in primitives less in overhead

primitive calls process an array of values in a **loop**:

vectorwise

## "Vectorized In Cache Processing"

**vector = array of ~1000**

**processed in a tight loop**

**CPU cache Resident**

| 102 | ivan | 350 |
| 104 | peggy | 750 |

next()

**- 30**  **\* 50**

| 102 | ivan | 37 | 7 | 350 |
| 104 | peggy | 45 | 15 | 750 |

PROJECT

**> 30 ?**

next()

| 101 | alice | 22 | FALSE |
| 102 | ivan | 37 | TRUE |
| 104 | peggy | 45 | TRUE |
| 105 | victor | 25 | FALSE |

SELECT

next()

| 101 | alice | 22 |
| 102 | ivan | 37 |
| 104 | peggy | 45 |
| 105 | victor | 25 |

SCAN

vectorwise

**Observations:**

next() called much less often ➔ more time spent in primitives less in overhead

primitive calls process an array of values in a **loop**:

**CPU Efficiency depends on "nice" code**
- out-of-order execution
- few dependencies (control,data)
- compiler support

**Compilers like simple loops over arrays**
- loop-pipelining
- automatic SIMD

**Observations:**

next() called much less often ➜ more time spent in primitives less in overhead

primitive calls process an array of values in a **loop**:

**CPU Efficiency depends on "nice" code**
- out-of-order execution
- few dependencies (control,data)
- compiler support

**Compilers like simple loops over arrays**
- loop-pipelining
- automatic SIMD

> 30 ?

| |
|---|
| FALSE |
| TRUE |
| TRUE |
| FALSE |

```
for(i=0; i<n; i++)

    res[i] = (col[i] > x)
```

- 30

| |
|---|
| 7 |
| 15 |

```
for(i=0; i<n; i++)

    res[i] = (col[i] - x)
```

* 50

| |
|---|
| 350 |
| 750 |

```
for(i=0; i<n; i++)

    res[i] = (col[i] * x)
```

vectorwise

**Tricks being played:**

**- Late materialization**

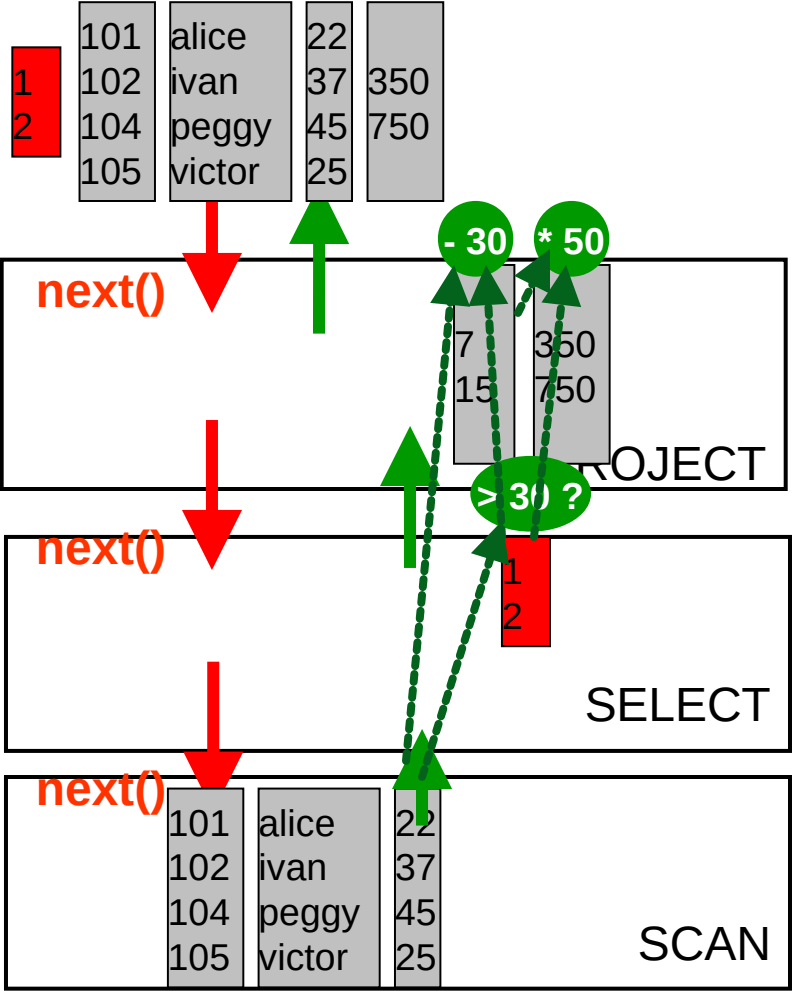**- Materialization avoidance using selection vectors**

| 1 2 | 101 102 104 105 | alice ivan peggy victor | 22 37 45 25 | 350 750 |

**next()**

- 30   * 50

| | 7 15 | 350 750 |

PROJECT

> 30 ?

**next()**

| 1 2 |

SELECT

**next()**

| 101 102 104 105 | alice ivan peggy victor | 22 37 45 25 |

SCAN

vectorwise

```
map_mul_flt_val_flt_col(
    float *res,
    int*   sel,
    float  val,
    float *col, int n)

{

    for(int i=0; i<n; i++)
            res[i] = val * col[sel[i]];
}
```

selection vectors used  to reduce vector copying

contain selected positions



next()

next()

- 30    * 50

> 30 ?

PROJECT

SELECT

next()

SCAN

101 alice 22
102 ivan 37 350
104 peggy 45 750
105 victor 25

1
2

7 350
15 750

1
2

101 alice 22
102 ivan 37
104 peggy 45
105 victor 25

# MonetDB/X100

- Both efficiency
  - Vectorized primitives
- and scalability..
  - Pipelined query evaluation

| | |
|---|---|
| C program: | 0.2s |
| MonetDB/X100: | 0.6s |
| MonetDB: | 3.7s |
| MySQL: | 26.2s |
| DBMS "X": | 28.1s |

# Memory Hierarchy

vectorwise

X100 query engine

tax

map_mul_flt_val_flt_col

PROJECT          0.19

selection vector

select_lt_int_col_int_val

SELECT          25

SCAN    age    name    salary

CPU cache

RAM    ColumnBM (buffer manager)

(raid) Disk(s)

Small Fast Expensive

CPU

registers

~10 GB/s 2–20 cycles

CPU Cache

2–3 GB/s 150–250 cycles

Main Memory

40–400 MB/s millions of cycles

Harddrive

Large Slow Cheap

# Memory Hierarchy

Vectors are only the in-cache representation

RAM & disk representation might actually be different

(vectorwise uses both PAX & DSM)

"MonetDB/X100: Hyper-Pipelining Query
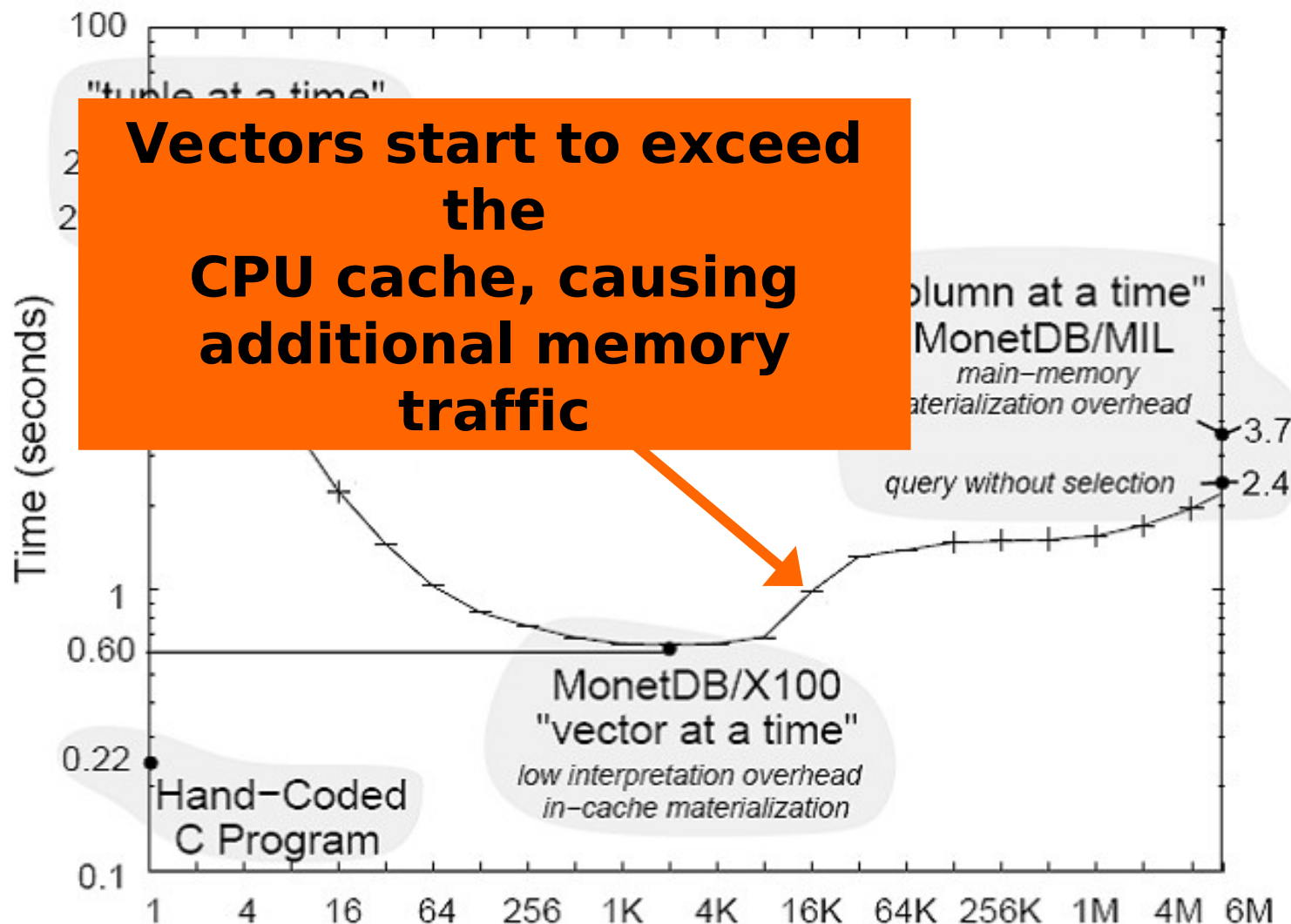Execution " Boncz, Zukowski, Nes, CIDR'05

# Varying the Vector size

# Varying the Vector size



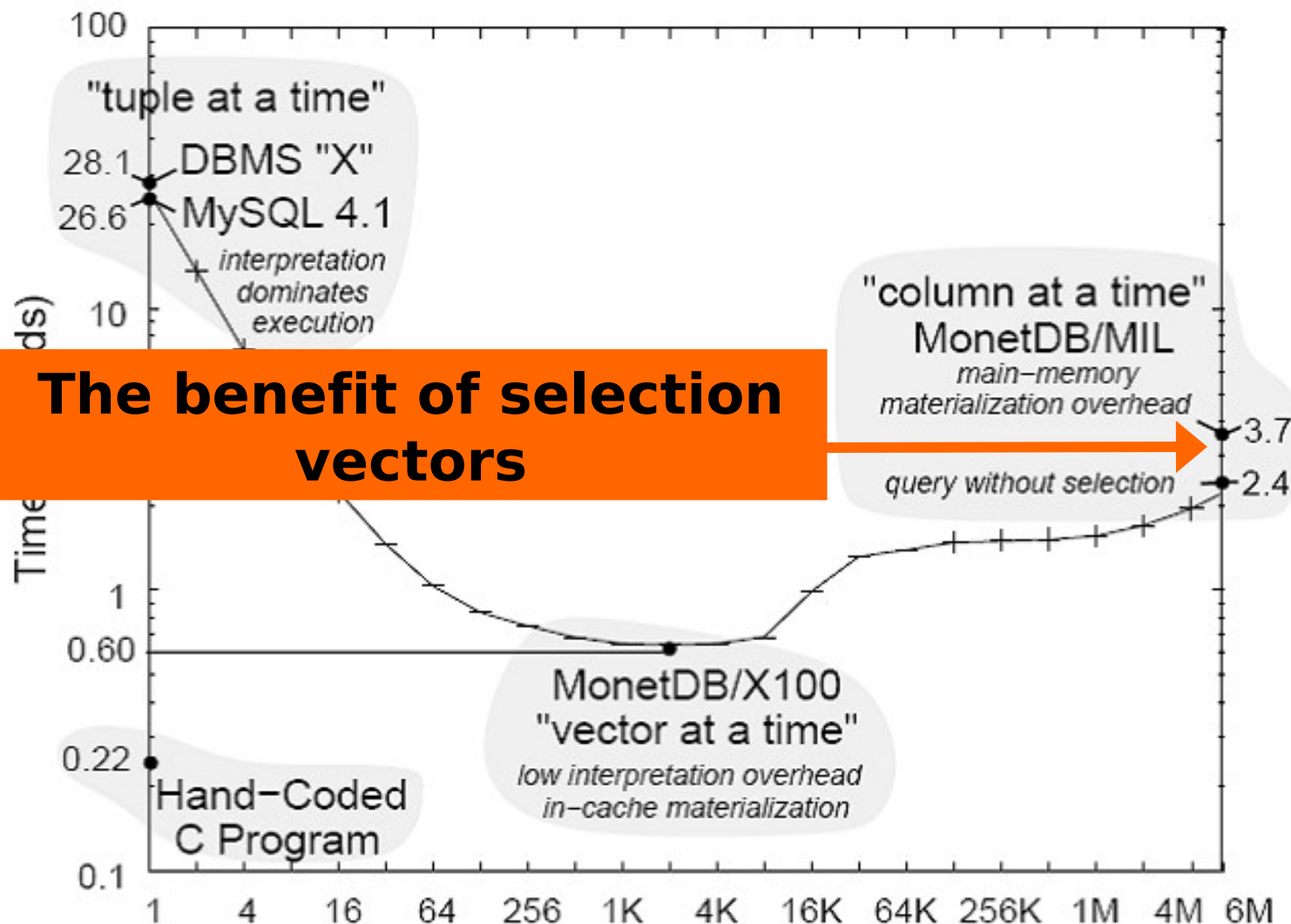Less and less iterator.next() and primitive function calls ("interpretation overhead")
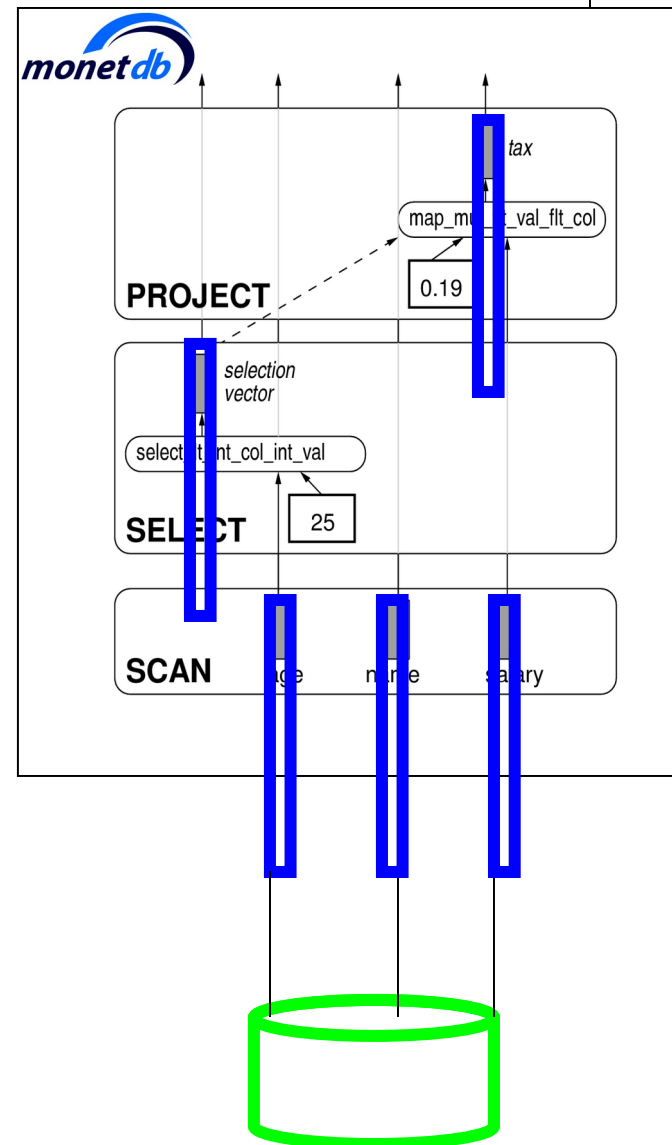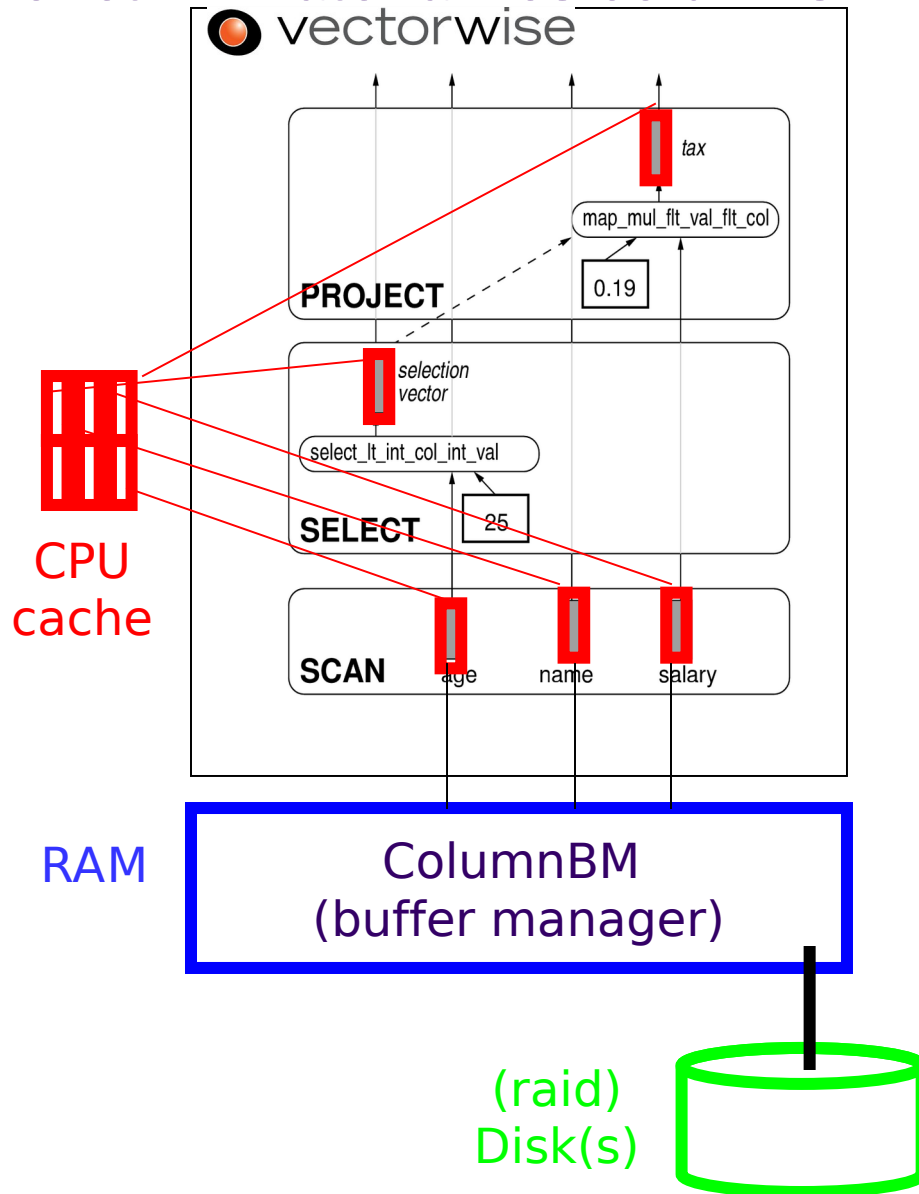
# Varying the Vector size



**Vectors start to exceed the CPU cache, causing additional memory traffic**

# Varying the Vector size



The benefit of selection vectors

# MonetDB materializes columns

CPU cache

RAM

ColumnBM
(buffer manager)

(raid)
Disk(s)

**vectorwise**

"MonetDB/X100: Hyper-Pipelining Query Execution" Boncz, Zukowski, Nes, CIDR'05

# Benefits of Vectorized Processing

"Buffering Database Operations for Enhanced Instruction Cache Performance" Zhou, Ross, SIGMOD'04

- **100x less Function Calls**
  - iterator.next(), primitives
- **No Instruction Cache Misses**
  - High locality in the primitives

"Block oriented processing of relational database operations in modern computer architectures" Padmanabhan, Malkemus, Agarwal, ICDE'01

- **Less Data Cache Misses**
  - Cache-conscious data placement
- **No Tuple Navigation**
  - Primitives are record-oblivious, only see arrays
- **Vectorization allows algorithmic optimization**
  - Move activities out of the loop ("strength reduction")
- **Compiler-friendly function bodies**
  - Loop-pipelining, automatic SIMD

# Vectorizing Relational Operators

- Project

- Select

  - Exploit selectivities, test buffer overflow

- Aggregation

  - Ordered, Hashed

- Sort

  - Radix-sort nicely vectorizes

- Join

  - Merge-join + Hash-join

# ADM: Literature

- **Column-Oriented Database Systems (4/6) - "Vectorized Execution"**

  - "MonetDB/X100: Hyper-Pipelining Query Execution". Boncz, Zukowski, Nes. CIDR'05.

  - "Buffering Database Operations for Enhanced Instruction Cache Performance". Zhou and Ross. SIGMOD'04.

  - "Block oriented processing of relational database operations in modern computer architectures". Padmanabhan, Malkemus, Agarwal. ICDE'01.

  - "Balancing Vectorized Query Execution with Bandwidth Optimized Storage". Zukowski. PhD Thesis. CWI 2008.

# ADM: Agenda

- <u>07.09.2022:</u> Lecture  1: **Introduction**

- <u>14.09.2022:</u> Lecture  2: **SQL Recap**

  *(plus Assignment 1 [in groups; 3 weeks]: TPC-H benchmark)*

- <u>21.09.2022:</u> Lecture  3: **Column-Oriented Database Systems (1/6) - Motivation & Basic Concepts**

- <u>28.09.2022:</u> Lecture  4: **Column-Oriented Database Systems (2a/6) - Selected Execution Techniques (1/2)**

- <u>05.10.2022:</u> Lecture  5: **Column-Oriented Database Systems (2b/6) - Selected Execution Techniques (2/2)**

  *(plus Assignment 2 [in groups; 4 weeks]: Compression techniques)*

- <u>12.10.2022:</u> Lecture  6: **Column-Oriented Database Systems (3/6) - Cache Conscious Joins**

- <u>19.10.2022:</u> Lecture  7: **Column-Oriented Database Systems (4/6) - "Vectorized Execution"**

- ~~<u>26.10.2022:</u>~~ ***No lecture!***

- <u>02.11.2022:</u> Lecture  8: **DuckDB: An embedded database for data science (1/2) (guest lecture & _hands-on_)**

  *(plus Assignment 3 [individual; 2 weeks]: Analysing NYC Cab dataset with DuckDB)*

- <u>09.11.2022:</u> Lecture  9: **DuckDB: An embedded database for data science (2/2) (guest lecture & _hands-on_)**

- <u>16.11.2022:</u> Lecture 10: **Branch Misprediction & Predication**

  *(plus Assignment 4 [individual; 2 weeks]: Predication)*

- <u>23.11.2022:</u> Lecture 11: **Column-Oriented Database Systems (5/6) - Adaptive Indexing**

- <u>30.11.2022:</u> Lecture 12: **Column-Oriented Database Systems (6/6) - Progressive Indexing**

***Bring your own laptop!***

# ADM: Literature (5/6)

- **DuckDB: An embedded database for data science**

  - "DuckDB: an Embeddable Analytical Database". Mark Raasveldt & Hannes Mühleisen. SIGMOD'19. Demo.

  - "Data Management for Data Science - Towards Embedded Analytics". Mark Raasveldt & Hannes Mühleisen. CIDR'20.

  - "Integrating Analytics with Relational Databases". Mark Raasveldt. PhD Thesis, Leiden University & CWI, 2020.

  - https://duckdb.org