

5. Column-based Storage Model

Architecture of Database Systems

Application Characteristics

OLTP (ON-LINE TRANSACTION PROCESSING)

- Mix between read-only and update queries
- Minor analysis tasks
- Used for data preservation and lookup
- Read typically only a few records at a time
- High performance by storing contiguous records in disk pages

| contract | client | date | name | price | city | product |
|----------------------|----------|------|-------|-------|------------|---------|
| 12302346 | 10042334 | | Eno | | Redmond | Car |
| 37611373 | 10067097 | | Gotz | | Berkeley | Redmond |
| update query | | | | | | |
| find client 10032112 | | | | | | |
| 95371001 | 10032112 | | Chen | | Seattle | House |
| lookup query | | | | | | |
| 51213123 | 10032423 | | Jones | | Washington | Travel |
| 54535545 | 10067823 | | Smith | | New York | House |
| 45447894 | 10013232 | | Doe | | Boston | Car |
| insert query | | | | | | |

OLTP queries:
access all
columns of
just one row.

OLAP (ON-LINE ANALYTICAL PROCESSING)

- Query-intensive DBMS applications
- Infrequent batch-oriented updates
- Complex analysis on large data volumes
- Read typically only a few attributes of large amounts of historical data in order to partition them and compute aggregates
- High performance by storing contiguous values of a single attribute

| contract | client | date | claim | city | product |
|---|--------|------|-------|------|---------|
| select those tuples sold after march 21 | | | | | |
| sum claims | | | | | |
| while grouping by city and product | | | | | |
| OLAP query: accesses only a few columns of almost all rows. | | | | | |

Introduction

RECAP – ROW-BASED RECORD MANAGEMENT

- Classic N-ary storage model (NSM), also „row-store“
- NSM = Normalized Storage Model

ADVANTAGES

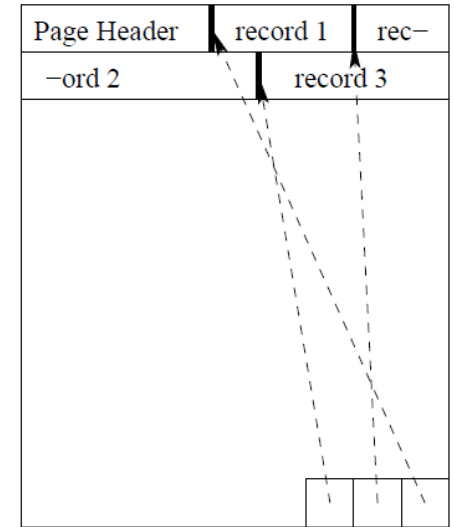
- Common database can be read with one page view
- single attributes values can simply be changed

DISADVANTAGES:

- even if only a few attributes value are needed, all attributes values have to be read
→ unnecessary IO-cost

ALTERNATIVES: COLUMN-ORIENTED STORAGE MODELS

- Decomposition of a n-ary relation in a set of projections (for example binary relation)
- Identification (and reconstruction) by a key column or position



NSM page organization

Example

ROW-BASED RECORD MANAGEMENT VERSUS COLUMN-BASED RECORD MANAGEMENT

| ProdNr | Bezeichnung | Preis |
|--------|-------------------|-------|
| 1 | Jamaica Blue | 8,55 |
| 2 | Arabica Black | 9,95 |
| 3 | New York Espresso | 10,95 |
| 4 | Guatemala Grande | 11,95 |
| 5 | Breakfast Blend | 9,90 |



| SAdr | ProdNr | SAdr | Bezeichnung | SAdr | Preis |
|------|--------|------|-------------------|------|-------|
| 0x00 | 1 | 0x00 | Jamaica Blue | 0x00 | 8,55 |
| 0x01 | 2 | 0x01 | Arabica Black | 0x01 | 9,95 |
| 0x02 | 3 | 0x02 | New York Espresso | 0x02 | 10,95 |
| 0x03 | 4 | 0x03 | Guatemala Grande | 0x03 | 11,95 |
| 0x04 | 5 | 0x04 | Breakfast Blend | 0x04 | 9,90 |

Dataset as a unit

Set of vertical projections

| | | | | | |
|-------------|------|---|------|---|---------|
| PAGE HEADER | | 1 | 0962 | | |
| 2 | 7658 | 3 | 3859 | 4 | 5523 |
| | | | | | • • • • |

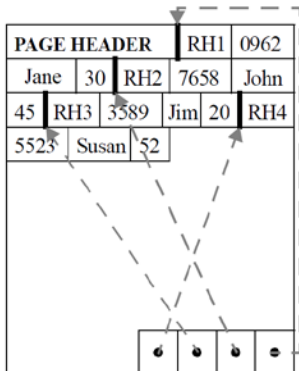
sub-relation R1

| | | | | | |
|-------------|------|---|------|---|---------|
| PAGE HEADER | | 1 | Jane | | |
| 2 | John | 3 | Jim | 4 | Susan |
| | | | | | • • • • |

sub-relation R2

| | | | | | | |
|-------------|---|----|---|----|---------|---|
| PAGE HEADER | | | | 1 | 30 | 2 |
| 45 | 3 | 20 | 4 | 52 | | |
| | | | | | • • • • | |

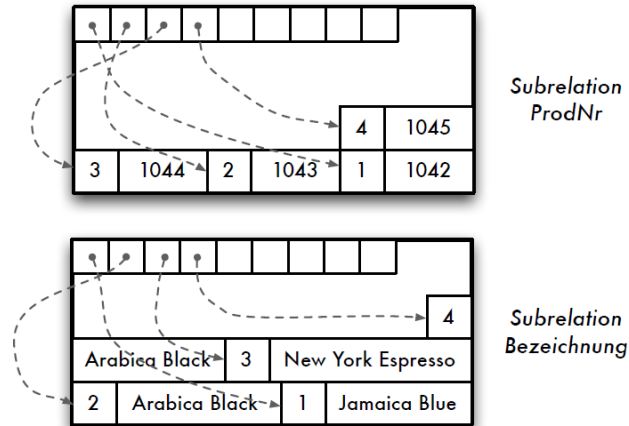
sub-relation R3



Decomposition Storage Model - DSM

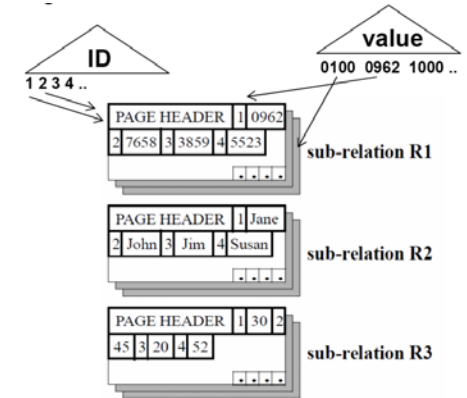
DESCRIPTION

- All values of a column (attribute) are saved consecutively
- Addressing by position / logical ID (surrogate)
- Page view (data set consisting of 2 attributes)



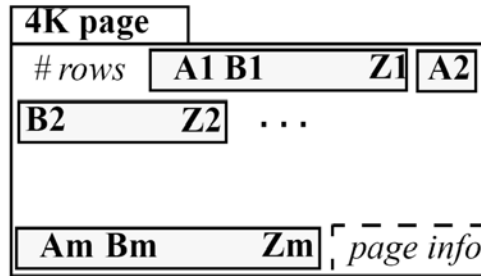
1985: DSM (DECOMPOSITION STORAGE MODEL)

- Proposed as an alternative to NSM (Normalized Storage Model)
- Decomposition storage mode, decomposes relations vertically
- 2 indexes: clustered on ID, non-clustered on value
- Speeds up queries projecting few columns
- Disadvantages: storage overhead for storing tuple IDs, expensive tuple reconstruction costs



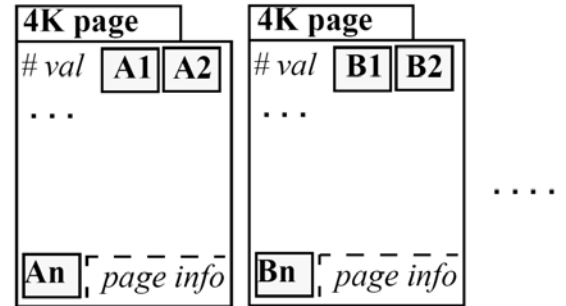
Comparison

ROW STORAGE



- + easy to add/modify a record
- might read unnecessary data

COLUMN STORAGE



- + only need to read in relevant data, more efficient scan operators
- + Compression easily possible (for example run length encoding)
- tuple writes require multiple accesses
- Reading all columns of a single row requires expensive row-reconstruction (1:1 join)

PAX – the Best of Both Worlds?

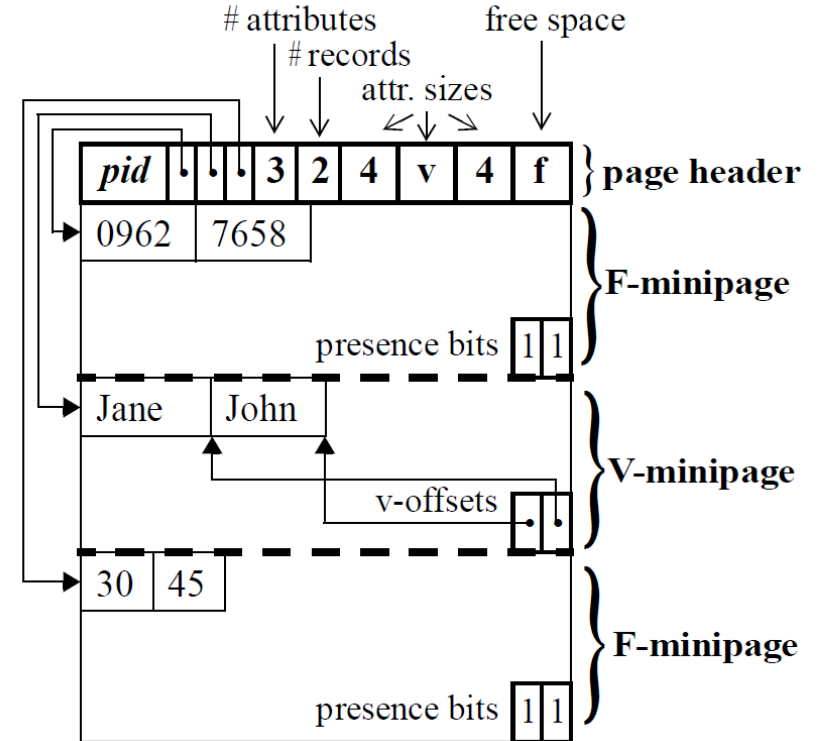
Partition Attributes Across - PAX

GOALS

- Maximizes inter-record spatial locality
- Incurs a minimal record reconstruction cost

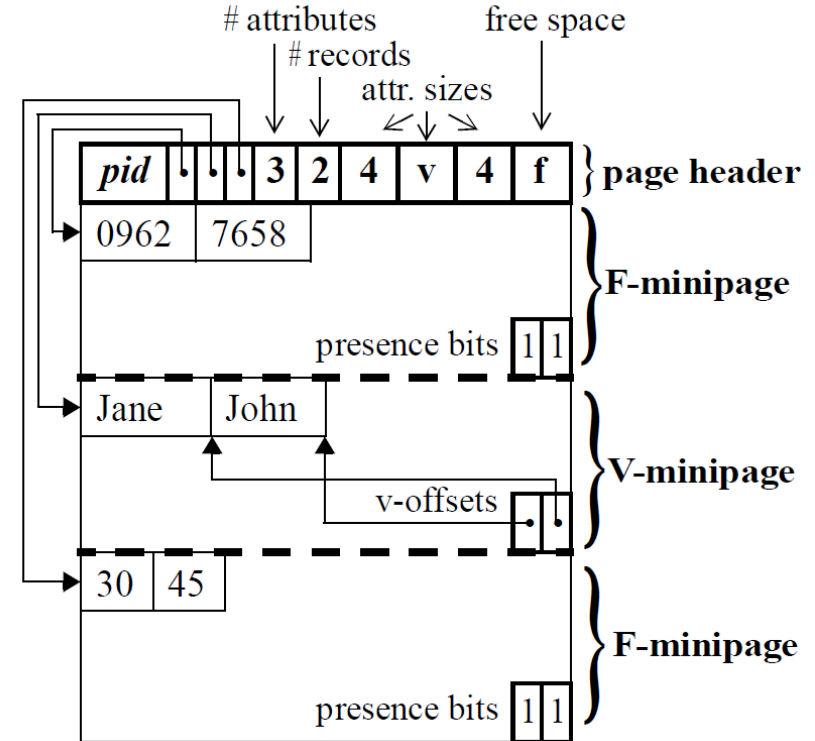
APPROACH

- compromise between NSM and DSM
- keep attributes values of each record on the same page
- using a cache-friendly algorithm for placing attributes values inside the page
 - vertically partitions the records within each page
 - storing together the values of each attribute in minipages

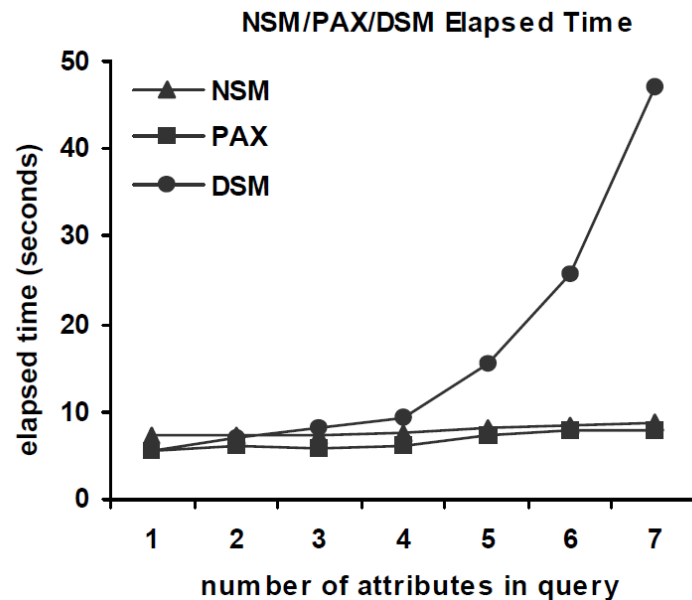


STORAGE DESIGN

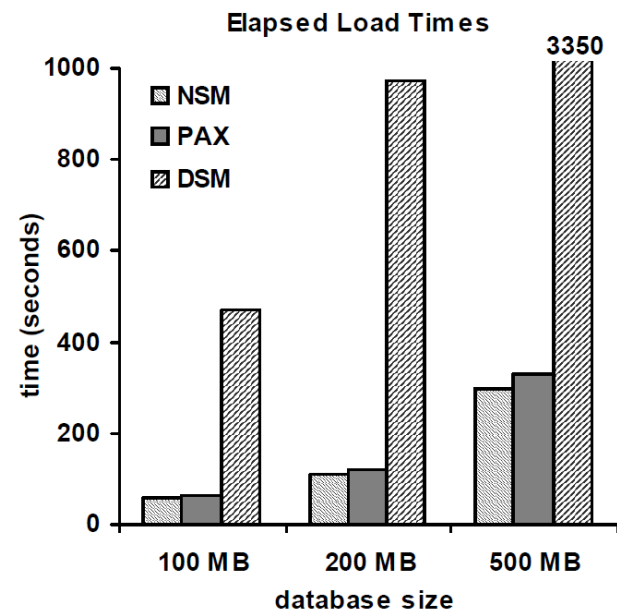
- each page is partitioned in n minipages (n attributes in a relation)
- Page Header
 - pointers to the beginning of each minipage
 - free space information
 - number of records
 - attributes sizes (fixed length or variable)
- Minipages
 - F-minipage → fixed-length attribute values, presence bits indicate the availability of attributes values for the records (if null, the attribute value is not present)
 - V-minipage → variable-length attributes values, slotted with pointers to each value, null values are denoted by null pointers



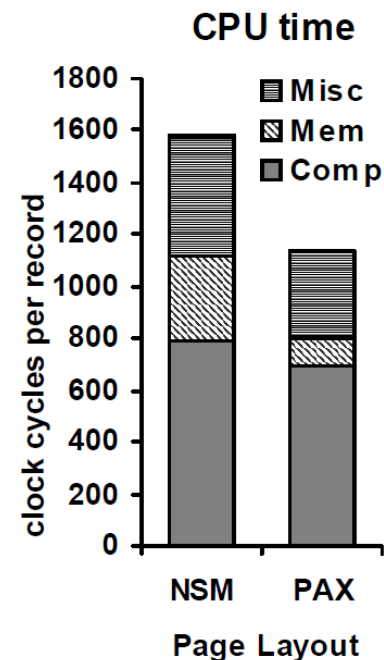
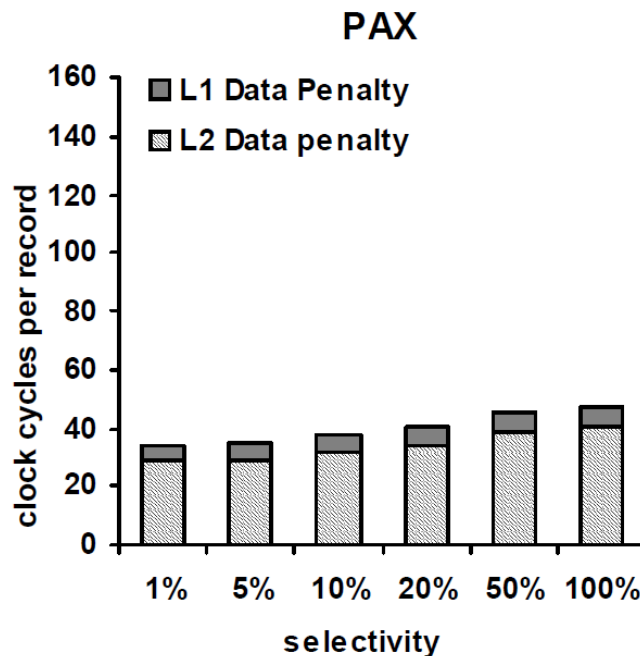
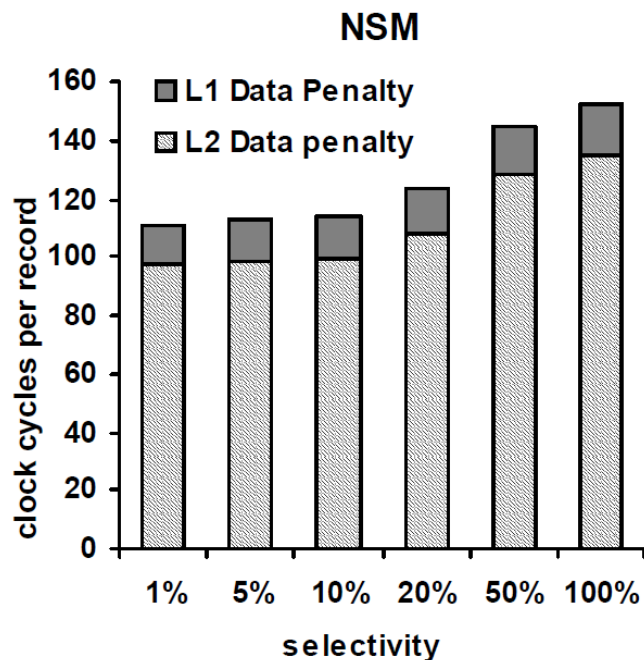
QUERY PERFORMANCE (READ)



BULK-LOADING



Cache Behavior



History & Development

From DSM to Column-Stores

1985: DECOMPOSITION STORAGE MODEL

LATE 90s – 2000s: FOCUS ON MAIN-MEMORY

PERFORMANCE

- MonetDB
- PAX: Partition Attributes Across
 - Retains NSM I/O pattern
 - Optimizes cache-to-RAM communication

2005: THE (RE)BIRTH OF COLUMN-STORES

- New hardware and application realities
 - Faster CPUs, larger memories, disk bandwidth
 - Multi-terabyte Data Warehouses
- New approach: combine several techniques
 - Read-optimized, fast multi-column access, disk/CPU efficiency, light-weight compression
- Used in read oriented environments - OLAP

SOME COLUMN STORE SYSTEMS

- MonetDB, C-Store, Sybase IQ, SAP HANA, Infobright, Exasol, X100/VectorWise

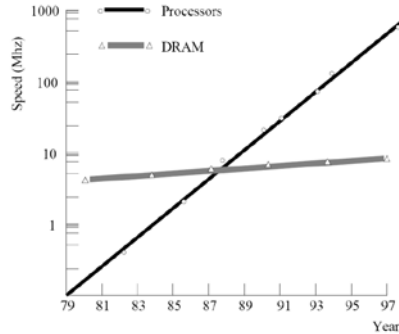
MAIN REASONS FOR COLUMNAR DATA REPRESENTATION

- Significantly higher cache hit rate
 - Data is more compact represented
- Use of HW-prefetcher
 - Sequential memory scans trigger pre-load of memory pages
- Use of SIMD-instructions
 - Multiple logical operations within one single CPU operation
- Opportunity for aggressive compression
 - Dictionary encoding as prerequisite

Hardware Development - Memory Wall

HARDWARE IMPROVEMENTS NOT EQUALLY DISTRIBUTED

- Advances in CPU speed have outpaced advances in RAM latency



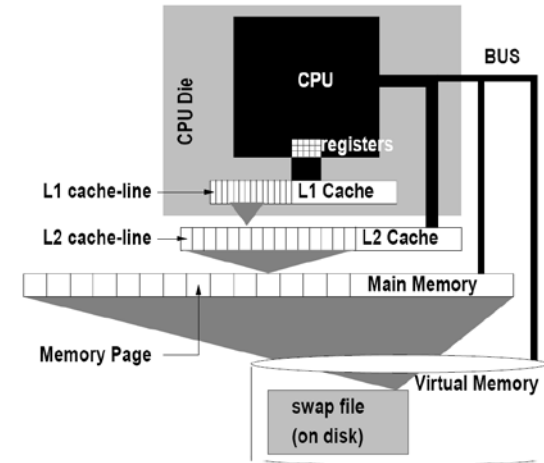
- Main-memory access has become a performance bottleneck for many computer applications
 - Bandwidth
 - Latency
 - Address translation (TLB)

→ Memory Wall

CACHE MEMORIES CAN REDUCE THE MEMORY LATENCY

WHEN THE REQUESTED DATA IS FOUND IN THE CACHE.

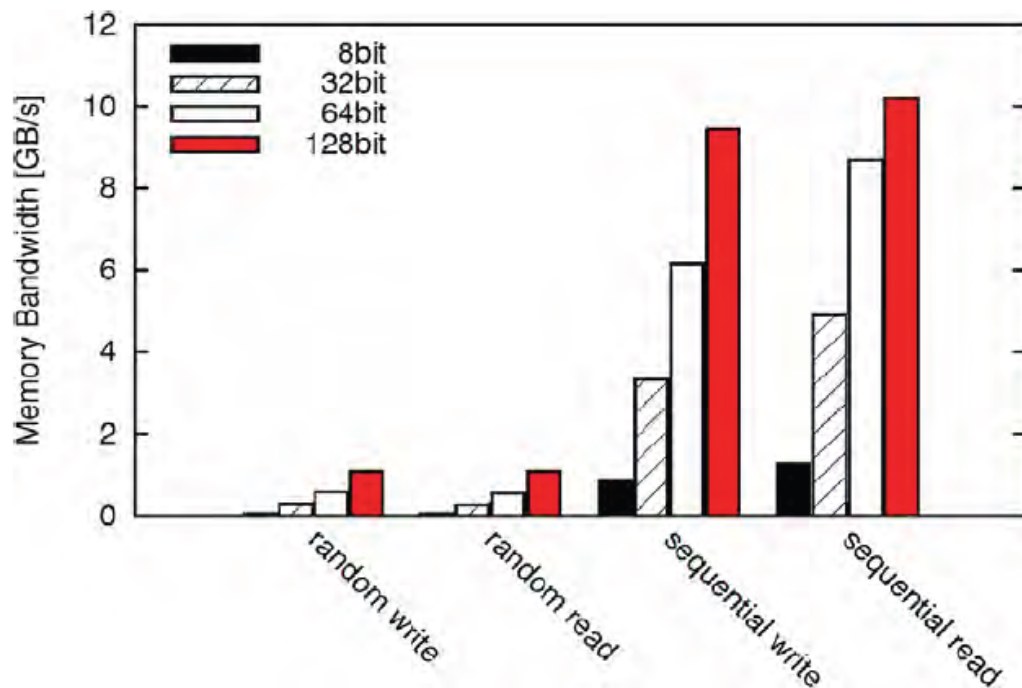
- Vertically fragmented data structures optimize memory cache usage



Memory Performance Comparison

IMPACT OF CACHES

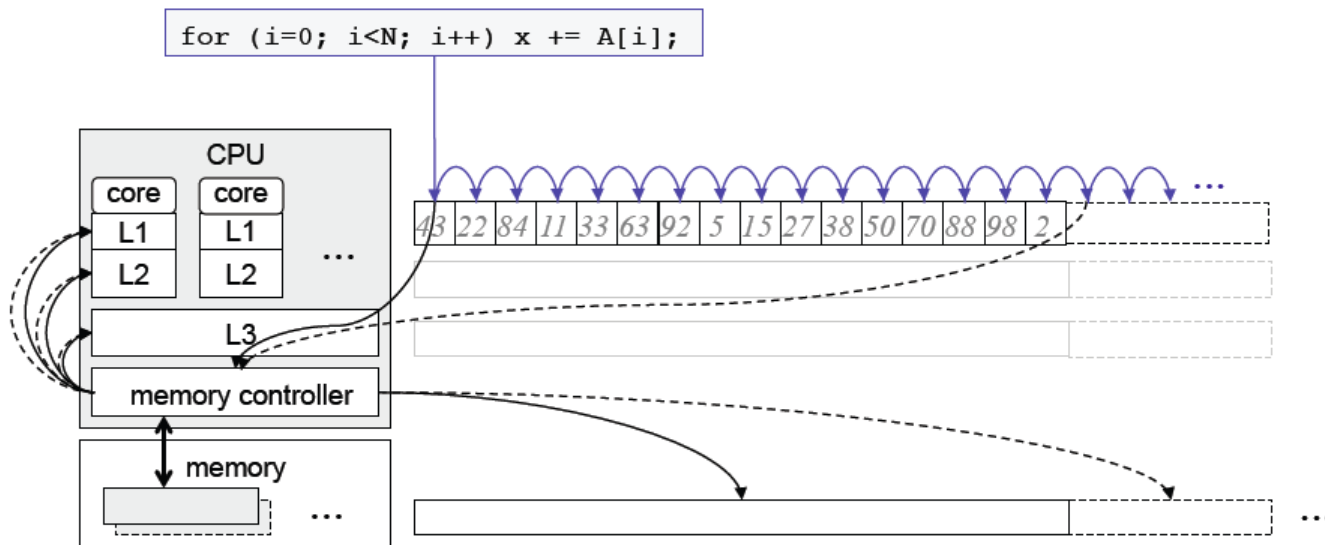
- Better cache efficiency
 - Multiple data snippets within the same cache line
- HW prefetcher
 - Sequential access trigger the pre-fetcher to pre-load subsequent memory pages



Results for a quad-core i7 2.66GHz, DDR3 1666. 32GB data accessed total.

CACHES – THE SUNNY SIDE

- Memory is physically accessed at cache line granularity, e.g. 64Byte
- Sequential

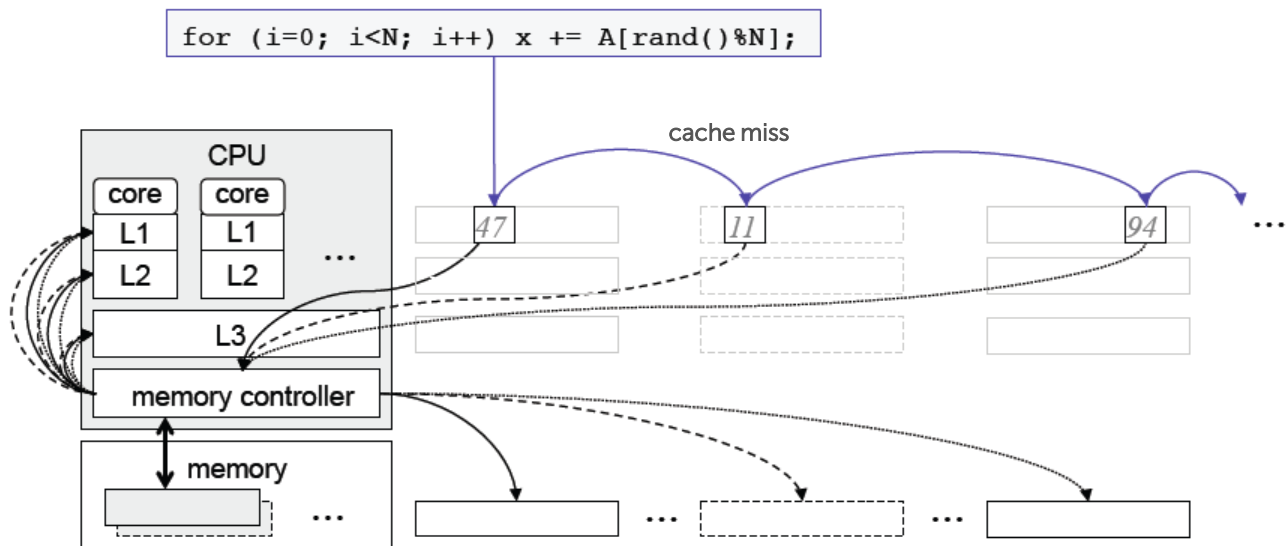


→ linear memory access maximizes cache & bandwidth utilization

The Role of Caches

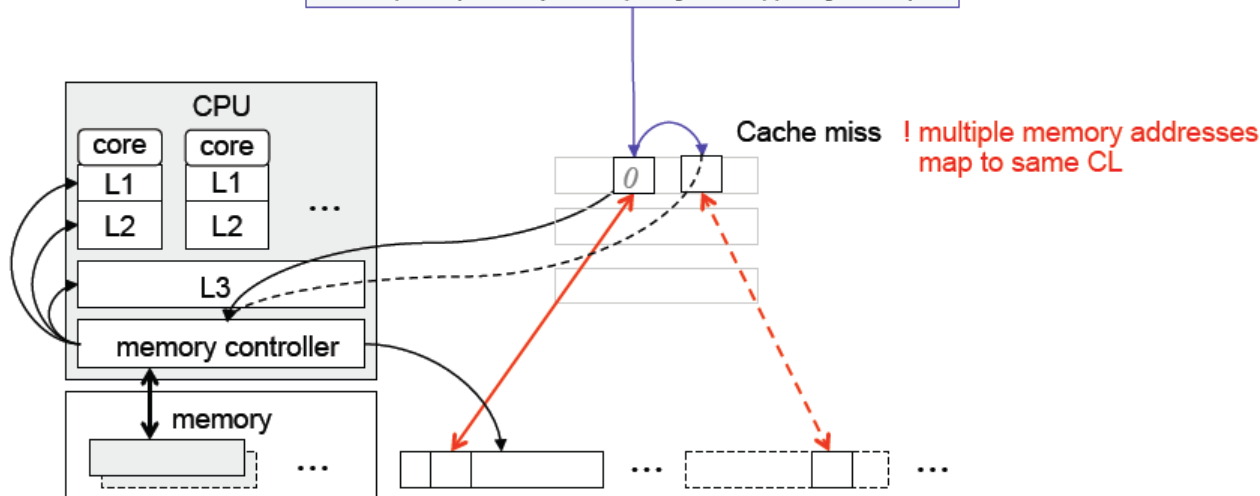
CACHES – THE BAD SIDE

- Memory is physically accessed at cache line granularity, e.g. 64Byte
- Random memory access:



→ Random memory access wastes up to 98.5%^{*} of bandwidth

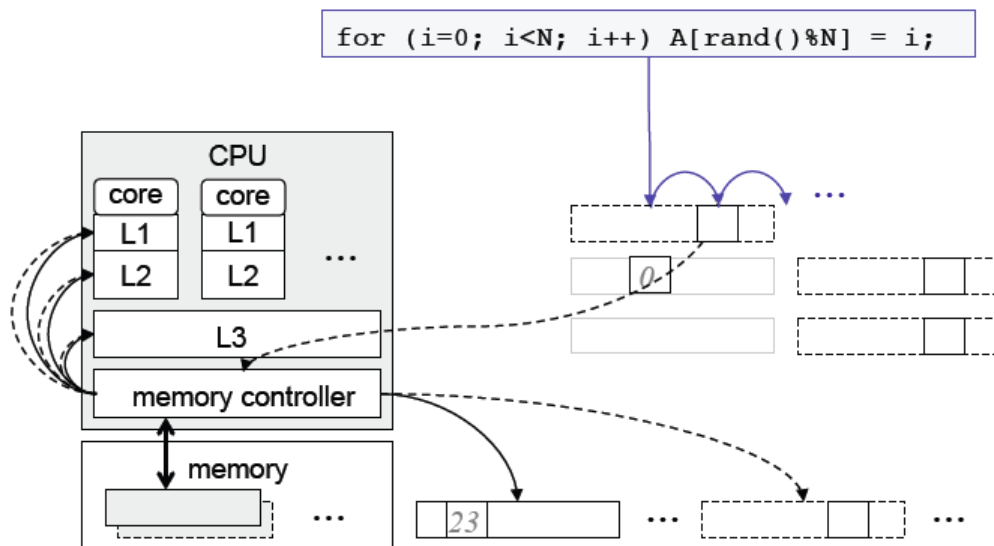
- Memory is physically accessed at cache line granularity, e.g. 64Byte
- Writes effectively turn into read-modify-write
 - Many memory addresses map into the same cache line(s)
 - "Dirty" cache line needs to be evicted before new one loads



The Role of Caches

CACHES – THE UGLY

- Memory is physically accessed at cache line granularity, e.g. 64Byte
- Writes effectively turn into read-modify-write
 - Many memory addresses map into the same cache line(s)
 - “Dirty” cache line needs to be evicted before new one loads



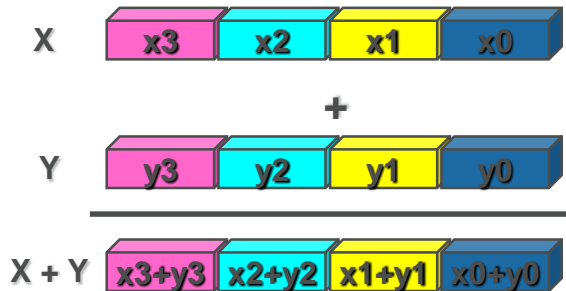
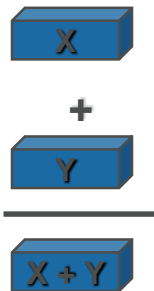
SIMD = Single Instruction Multiple Data

WHAT IS IT?

- Extension of the ISA. Data Types and instructions for the parallel computation on short vectors of integers or floats

Scalar processing

- traditional mode
- one operation produces one result



SIMD processing

- with SSE
- one operation produces multiple results

WHY DO THEY EXIST?

- Useful: Many applications have the necessary fine-grain parallelism
Then: speedup by a factor close to vector length
- Doable: Chip designers have enough transistors to play with

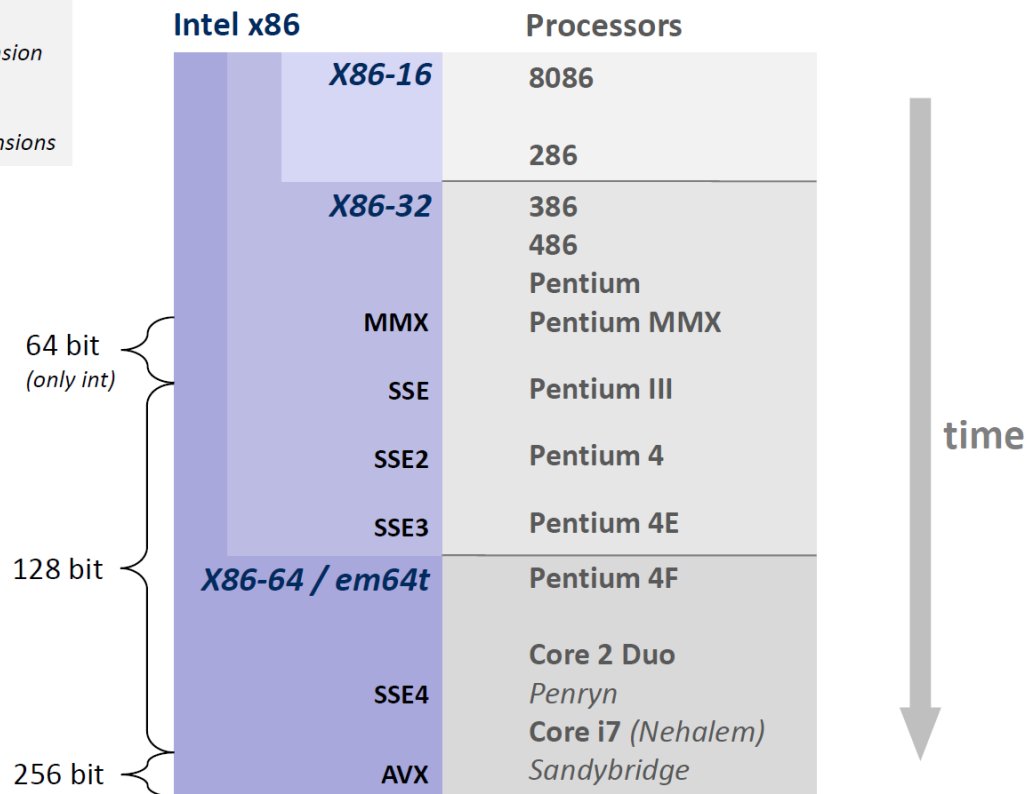
Development

SSE:

Streaming SIMD extension

AVX:

Advanced vector extensions



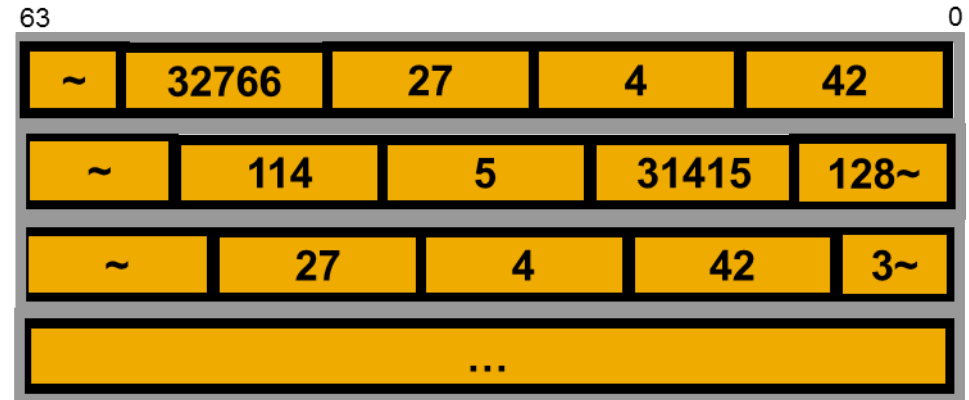
Example: Value Extraction

EXAMPLE

- Dictionary encoding provide token with length of: 15 bit
- Input: vector of 64 bit values (concatenated 15 bit values)
- Work: extracted values are compared (rangeFrom \leq value $<$ rangeTo)
- Output: 32 bit integer holding the index of a match

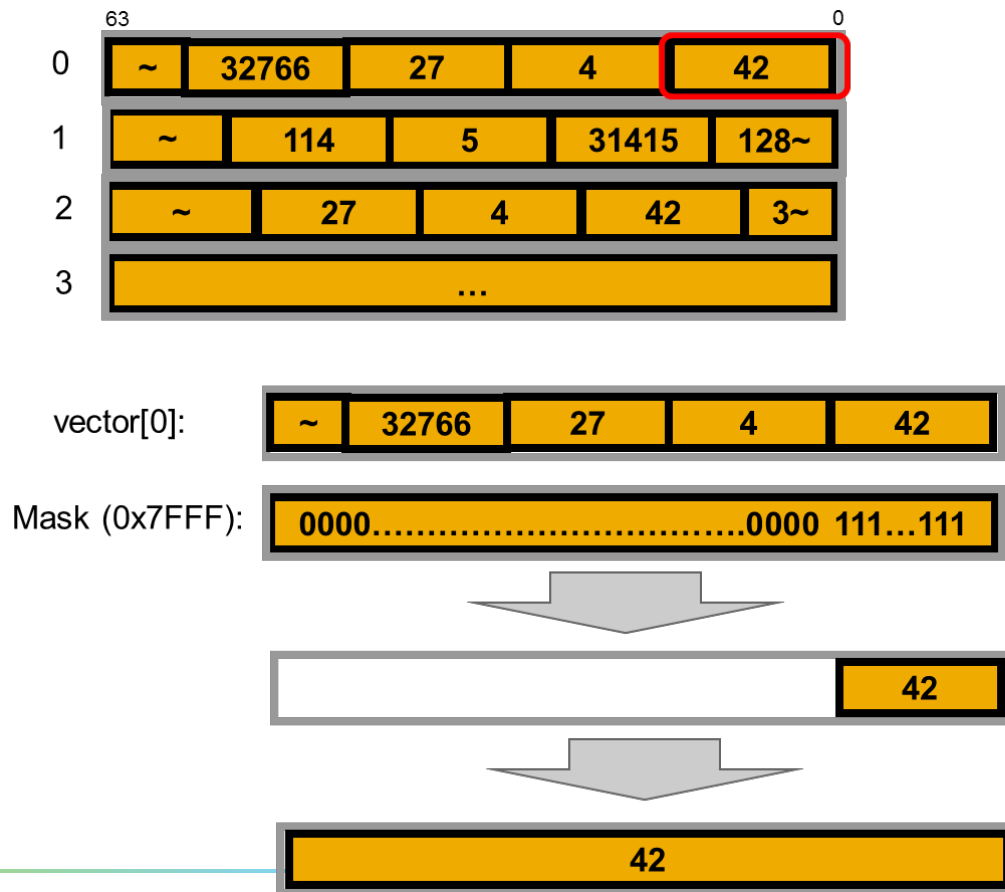
CODE

- Unrolled extraction loop to extract 64 values at once
- 15 x 64 bit input integer \rightarrow 64 x 15 bit values



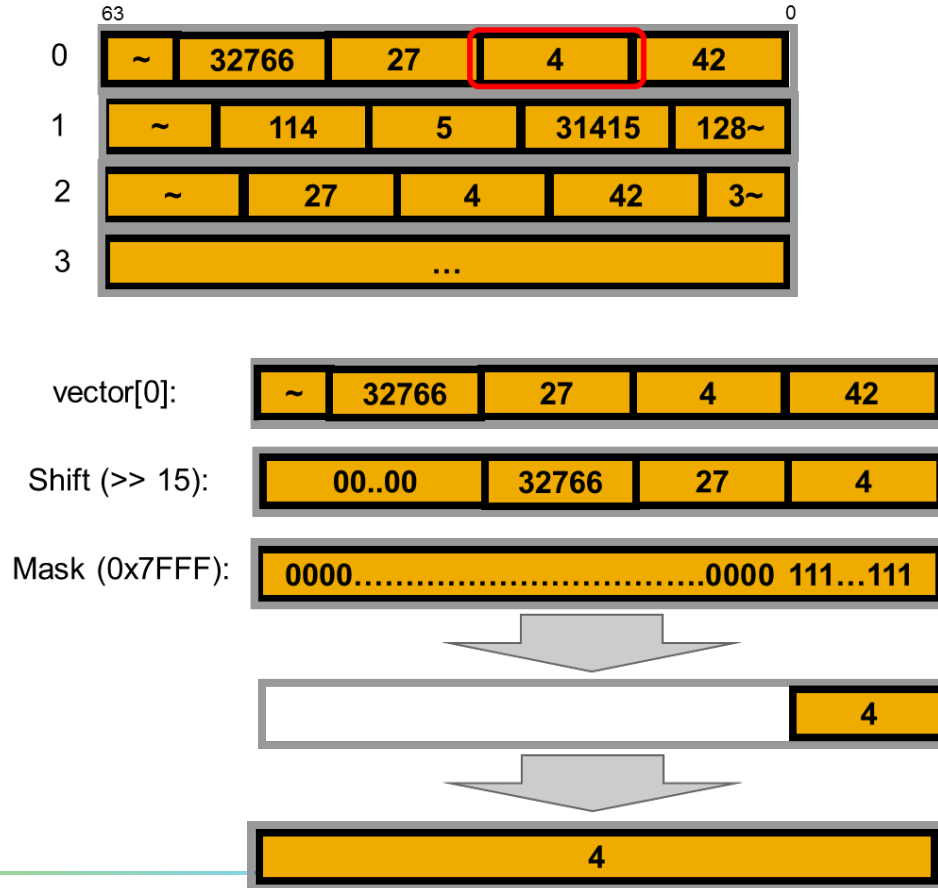
Example: Value Extraction

STEP 1:



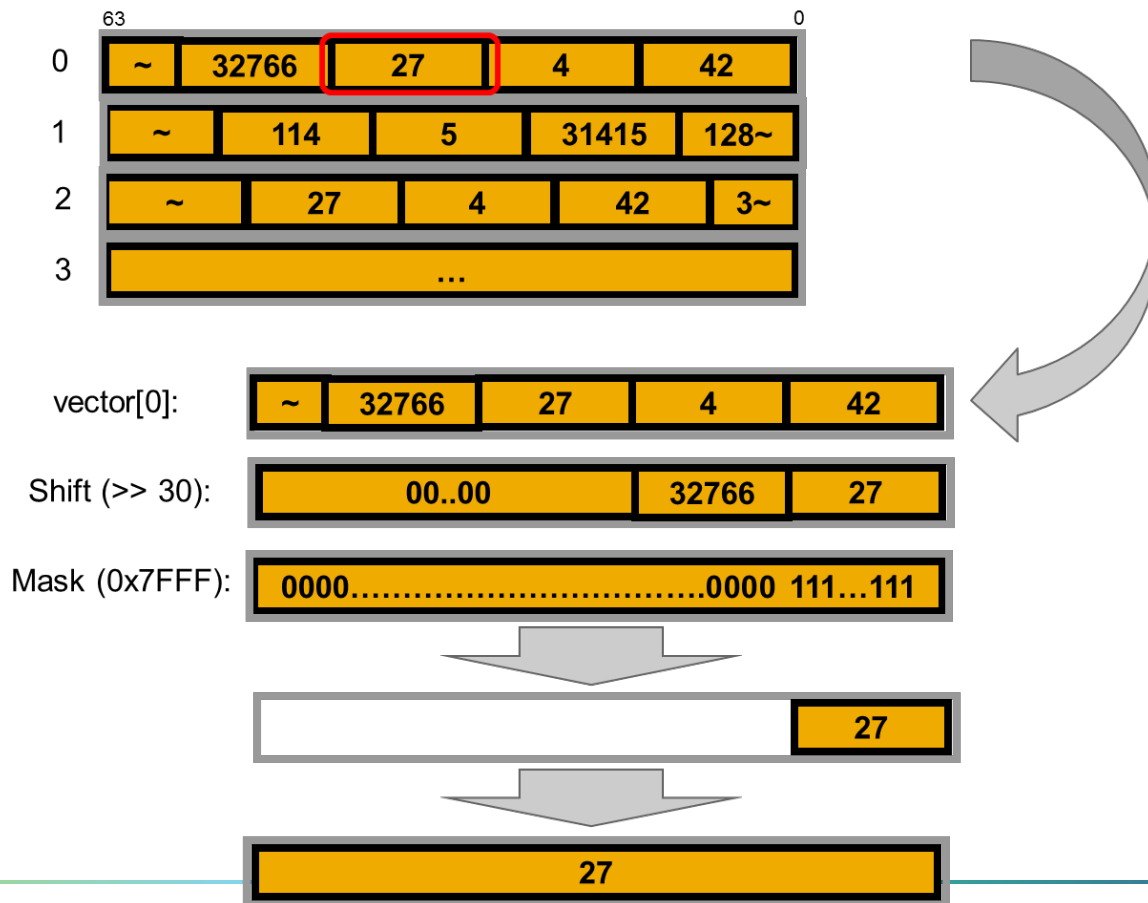
Example: Value Extraction

STEP 2:



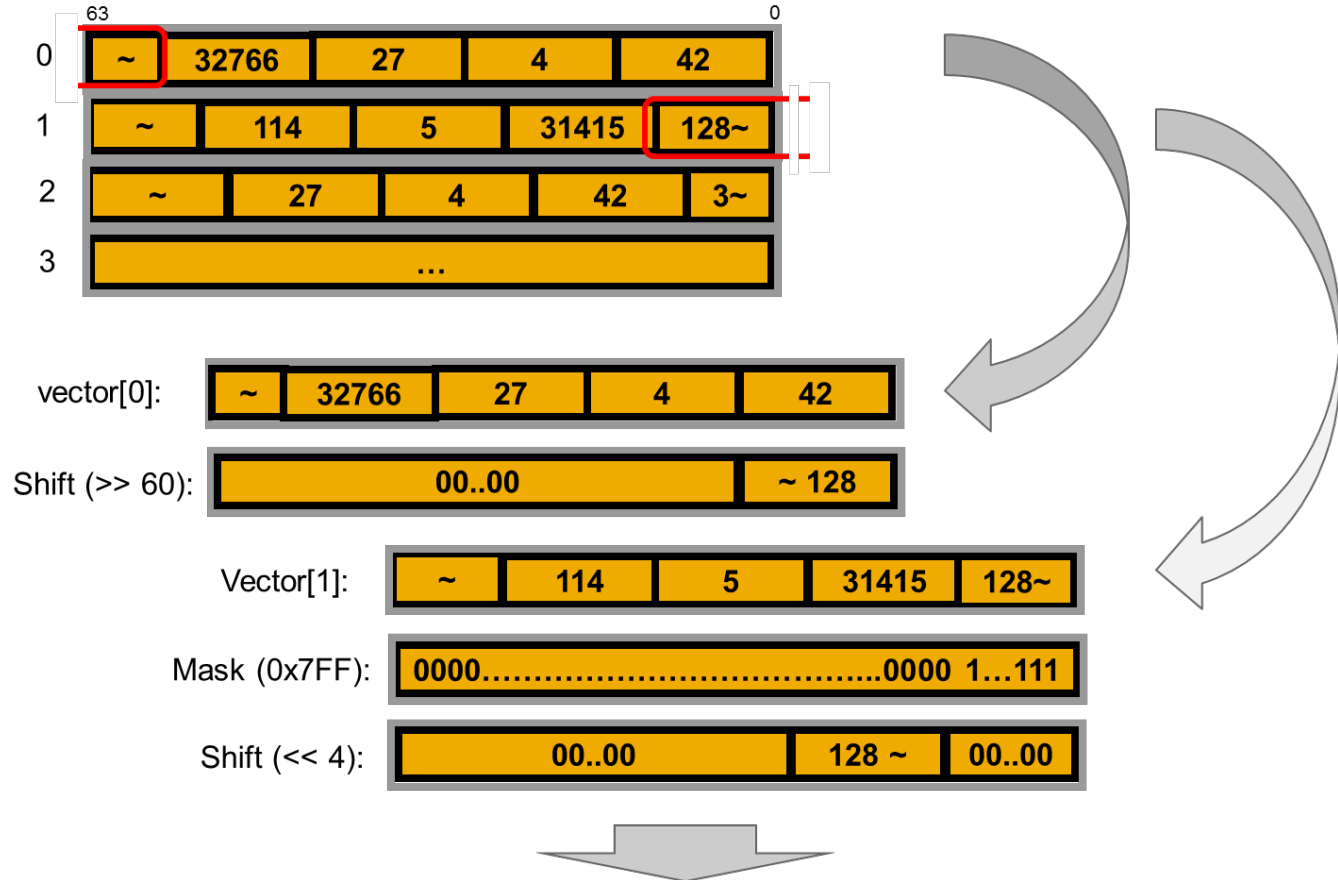
Example: Value Extraction

STEP 3:



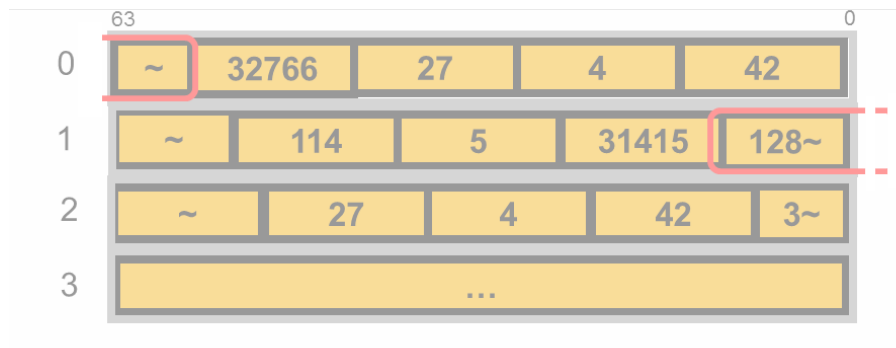
Example: Value Extraction

STEP 4:



Example: Value Extraction

STEP 4 CONT.:



vector[0]:



(shifted)

Vector[1]:



(shifted & masked)

OR:



BitWeaving

IDEA

- Fast scan method for column stores
- Two types of BitWeaving
 - BitWeaving/H (Horizontal bit organization)
 - BitWeaving/V (Vertical bit organization)

STORAGE LAYOUT

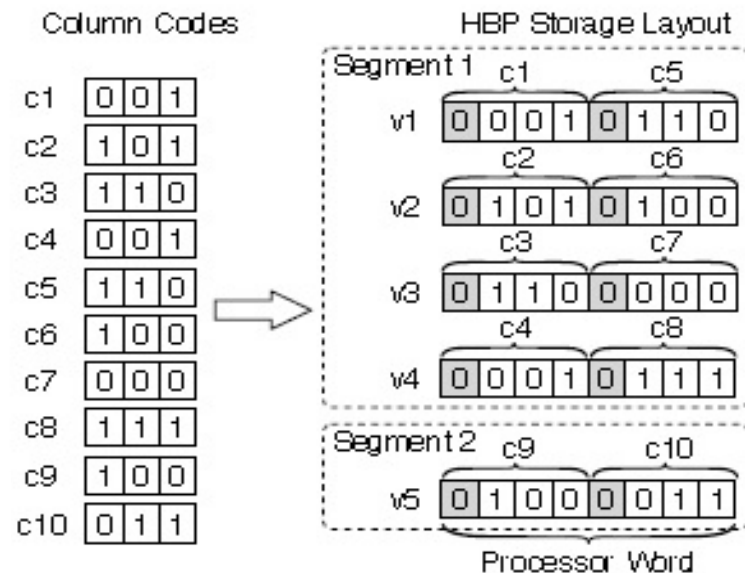
- Packs codes “horizontally” into processor words
- Uses an extra bit (delimiter bit) in each code
- Column-scalar scan: parallel predicate evaluation on packed codes

BitWeaving/H – Storage Layout

HORIZONTAL LAYOUT

- k denote the number of bits needed to encode a code
- each code is stored in a $(k + 1)$ -bit section whose leftmost bit is used as a delimiter between adjacent codes.
- w denotes the width of a processor word
→ inside the processor word, $\lfloor w/(k+1) \rfloor$ sections are concatenated together
- The column is divided into fixed-length segments, each of which contains $(k+1) * \lfloor w/(k+1) \rfloor$ codes

$$\begin{array}{l}
 \text{Query: } c < 5 \\
 \begin{array}{l}
 X = \begin{matrix} v_1(c_1, c_5) & v_2(c_2, c_6) & v_3(c_3, c_7) & v_4(c_4, c_8) & v_5(c_9, c_{10}) \\
 (0001 \ 0110)_2 & (0101 \ 0100)_2 & (0110 \ 0000)_2 & (0001 \ 0111)_2 & (0100 \ 0011)_2 \\
 Y = (0101 \ 0101)_2 & (0101 \ 0101)_2 & (0101 \ 0101)_2 & (0101 \ 0101)_2 & (0101 \ 0101)_2 \\
 mask = (0111 \ 0111)_2 & (0111 \ 0111)_2 & (0111 \ 0111)_2 & (0111 \ 0111)_2 & (0111 \ 0111)_2 \\
 X \oplus mask = (0110 \ 0001)_2 & (0010 \ 0011)_2 & (0001 \ 0111)_2 & (0110 \ 0000)_2 & (0011 \ 0100)_2 \\
 Y + (X \oplus mask) = (1011 \ 0110)_2 & (0111 \ 1000)_2 & (0110 \ 1100)_2 & (1011 \ 0101)_2 & (1001 \ 1001)_2 \\
 Z = (Y + (X \oplus mask)) \wedge \neg mask = (1000 \ 0000)_2 & (0000 \ 1000)_2 & (0000 \ 1000)_2 & (1000 \ 0000)_2 & (1000 \ 1000)_2
 \end{matrix}
 \end{array}
 \end{array}$$



BitWeaving/H storage layout ($k = 3$; $w = 8$)
Delimiter bits are marked in gray.

STORAGE LAYOUT

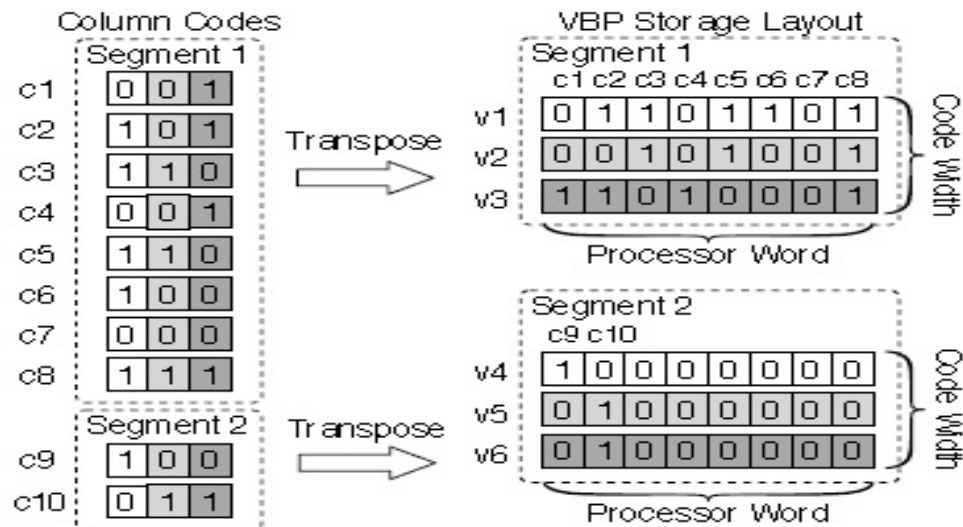
- Bit-level columnar data organization, i.e. its like a bit-level columnar store.

COLUMN-SCALAR SCAN

- Predicate evaluation is converted to logical computation on these "words of bits"

EXAMPLE

- middle bits of codes are marked in light gray, whereas the least significant bits are marked in dark gray



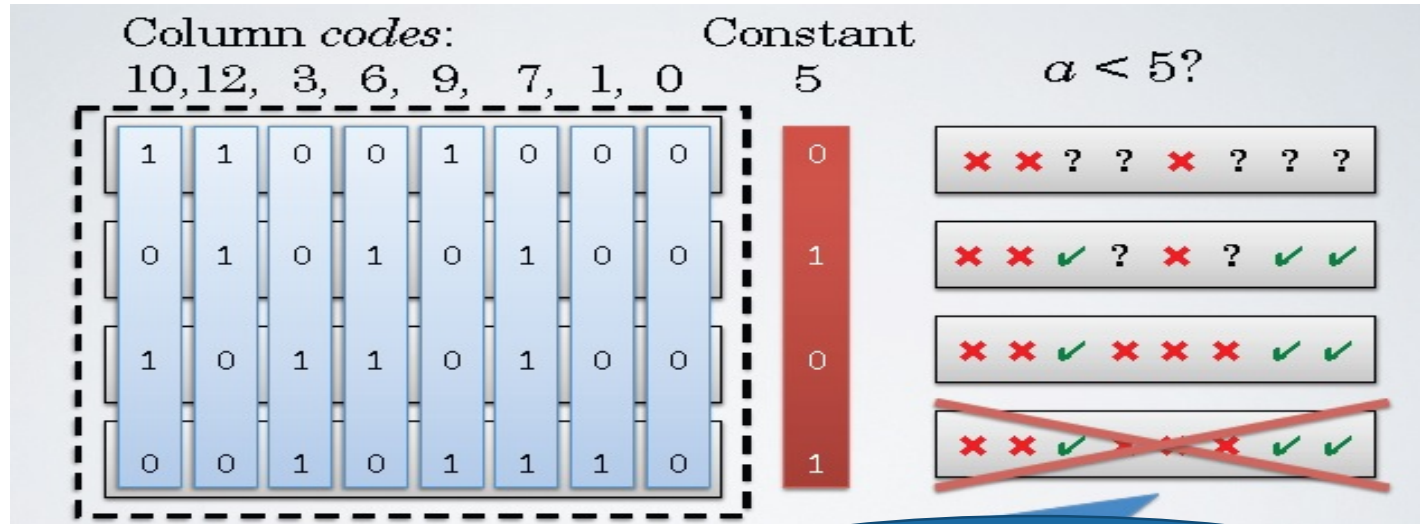
Example BitWeaving/V

QUERY

- $A < 5$

EARLY PRUNING

- Predicate eval may stop as soon as all results are identified



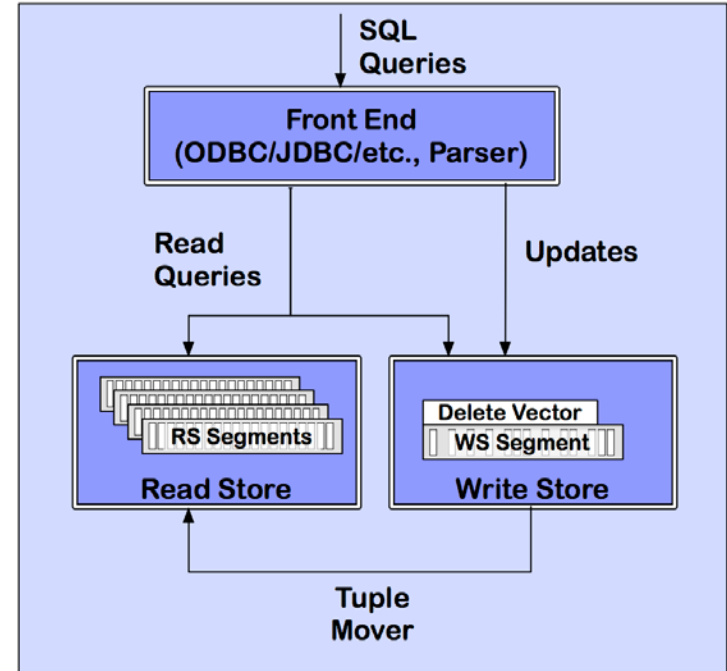
Early Pruning: terminate the predicate evaluation on a segment, when all results have been determined.



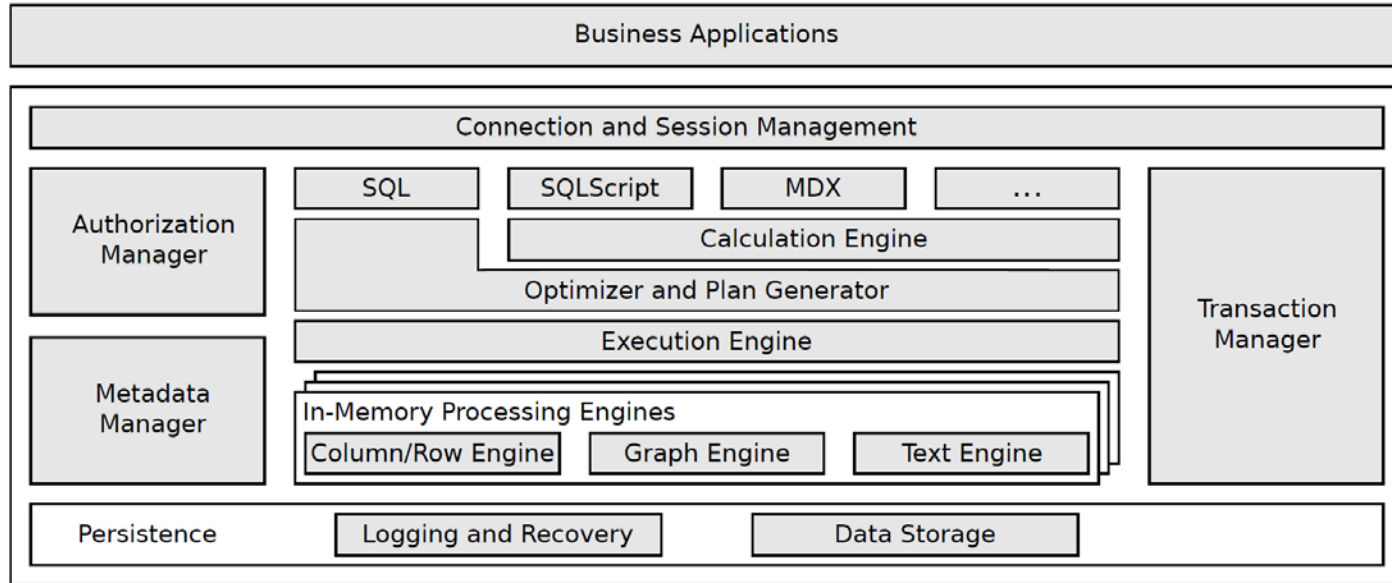
Architecture of commercial products

VERTICA ANALYTIC DATABASE

- DBMS Optimized for Next-Generation Data Warehousing (OLAP)
- Hybrid Store consisting of two distinct storage structures
 - WOS (Write-Optimized Store): fits into main memory and is designed to efficiently support insert and update operations; WOS is unsorted and uncompressed
 - ROS (Read-Optimized Store): bulk of the data; sorted and compressed; making it efficient to read and query
- Tuple Mover
 - Moves data out of the WOS and into ROS
- Structure
 - WOS and ROS are organized as DMS



ARCHITECTURE



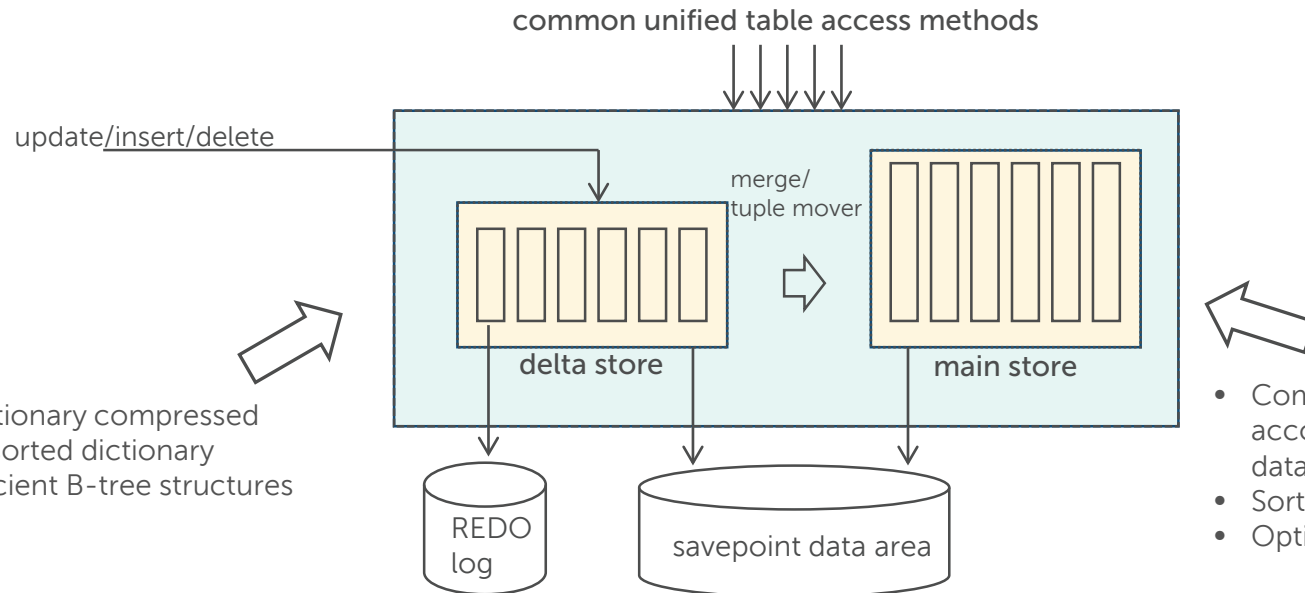
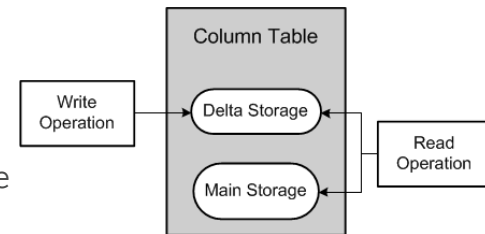
Hybrid Storage Architecture for Column Stores

USE OF COMPRESSION IMPLIES TWO STORES

- Write optimized store (WOS)
- Read optimized (compressed) store (ROS)

MERGE PROCESS

- Moves data from delta to main store

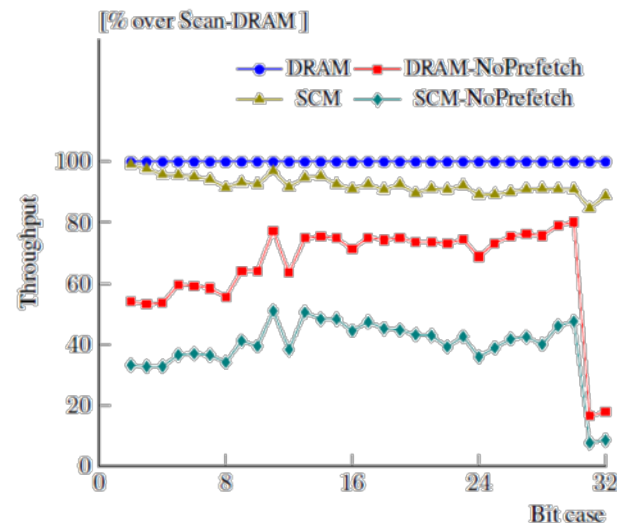
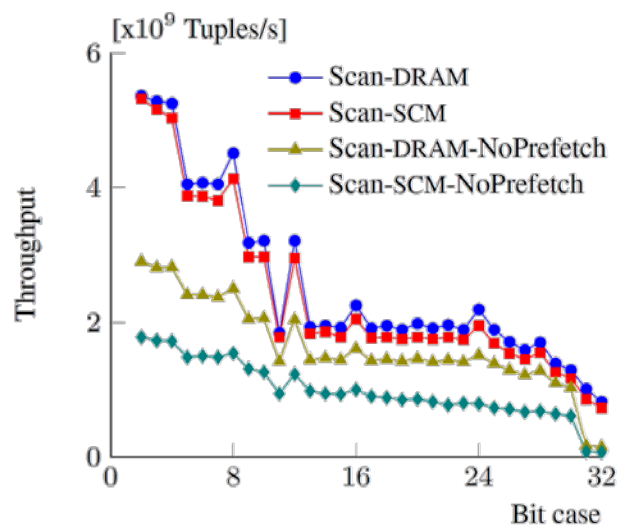


NVRAM for ROS-Structure

MICRO EXPERIMENTS: SIMD SCAN PERFORMANCE ON DRAM AND SCM

- With prefetching: average penalty for using SCM instead of DRAM is only 8%.
- Without prefetching: average penalty for using SCM instead of DRAM is 41%.

NOTE FOR OPERATORS WITH SEQUENTIAL MEMORY ACCESS PATTERNS, SCM PERFORMS ALMOST AS GOOD AS DRAM



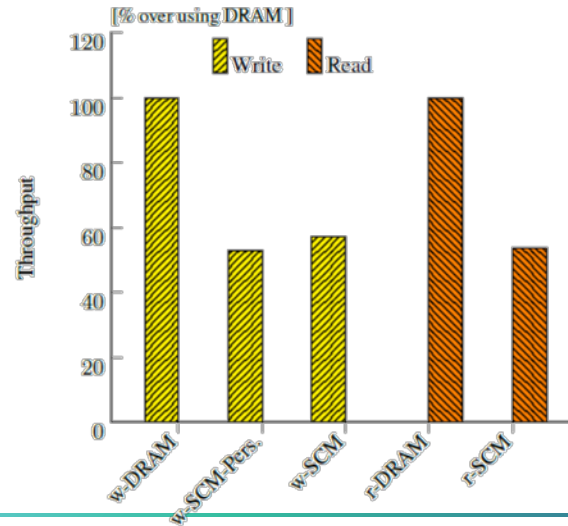
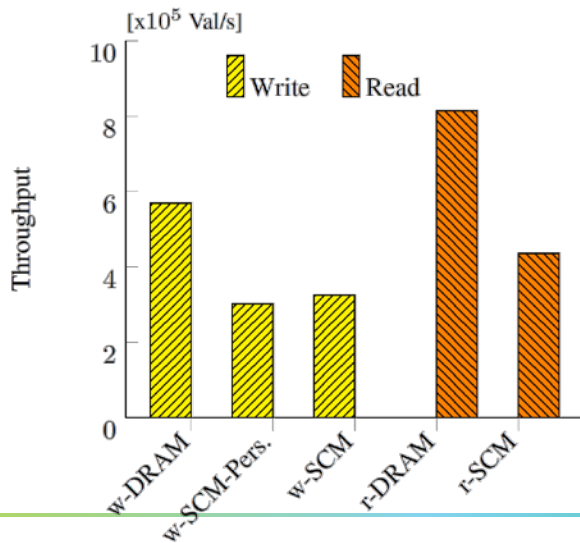
NVRAM for WOS-Structure

MICRO EXPERIMENT

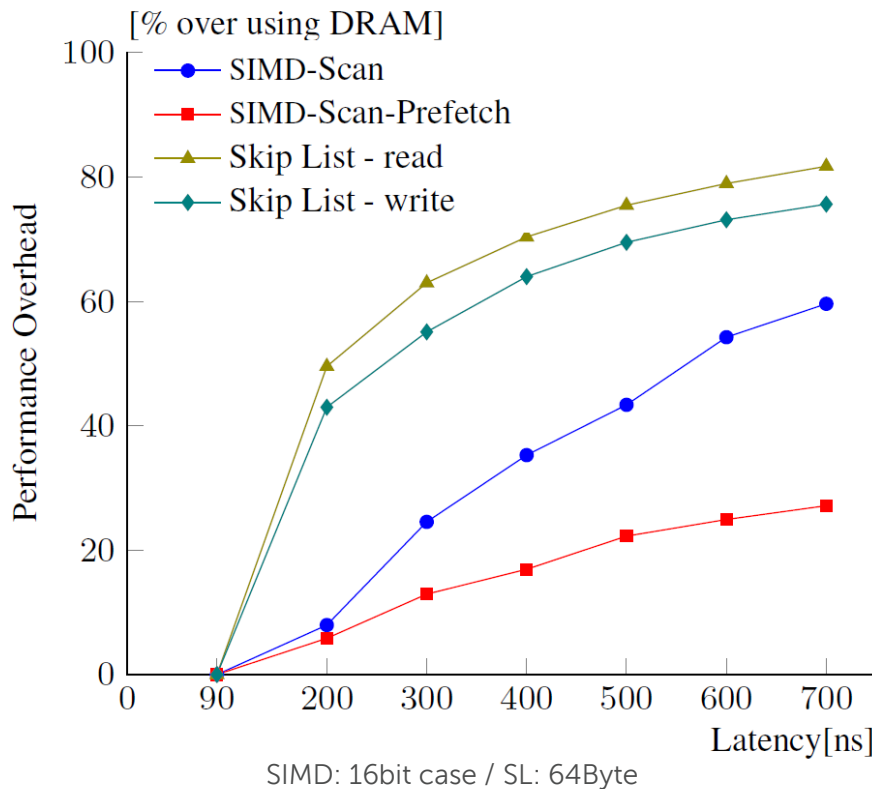
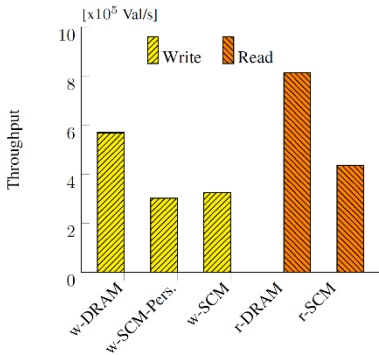
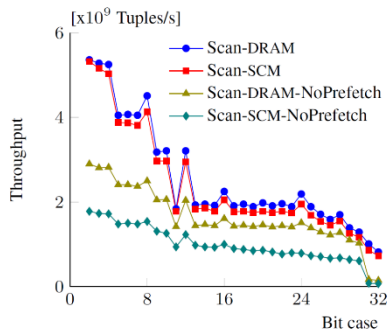
- Skip List read/write performance on DRAM and SCM
- 47% penalty for reads, and 43-47% penalty for writes for using SCM instead of DRAM.

OPERATIONS WITH RANDOM MEMORY ACCESS PATTERNS ARE EXPENSIVE IN SCM

- Writing persistent and concurrent data structures is NOT trivial



Impact of NVRAM latency



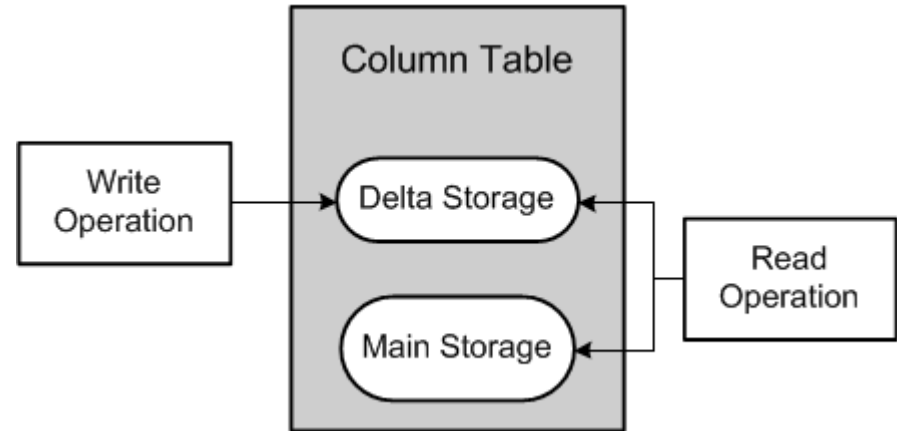
SAP HANA – Column Store

MAIN AND DELTA STORE

- Main Store: main part of the data; compressed data
- Delta Store: all data changes are written; basic compression and optimized for write access

MERGE PROCESS

- Moves data from delta to main store



THE DELTA MERGE OPERATION

