



**Hasso
Plattner
Institut**

IT Systems Engineering | Universität Potsdam

In-Memory Data Structures and Databases

Jens Krueger

Enterprise Platform and Integration Concepts
Hasso Plattner Institute

What to take home from this talk?

2

Answer to the following questions:

- What makes an in-memory database fast?
- What are differences of an in-memory database to disk-based systems?
- How does the physical data representation affect the performance of a in-memory database?
- How to leverage sequential data access?
- How can compression improve read access?

Recap

Jens Krueger

Enterprise Platform and Integration Concepts
Hasso Plattner Intitute

Recap: Workload Characteristics

4

OLTP	OLAP/DSS
Full row operations	Retrieve small number of columns
Simple Queries	Complex Queries
Detail Row Retrieval	Aggregation and Group By
Inserts/Updates/Selects	Mainly Selects
Short Transactions	Long Transactions
Small Found Sets	Large Found Sets
Pre-determined Queries	Adhoc Queries
Real Time Updates	Batch Updates
„Source of Truth“	Alternative representation

Recap: Trends in Enterprise Apps

5

Today's Enterprise Applications

- Complex processes
- Increased data set (but real-world events driven)
- Separated into OLTP and OLAP

Enterprise data management

- Wide schemas
- Sparse data with limited domain
- Workload consists of complex, analytic-style queries
- Workload is mostly:
 - Set processing
 - Read access
 - Insert instead of updates

 **Mixed Workload**

Question

6

Why is an in-memory database faster than a fully cached disk-based database?

Excursus: Disk-based Databases

Jens Krueger

Enterprise Platform and Integration Concepts
Hasso Plattner Intitute

Excursus: Magnetic Disks

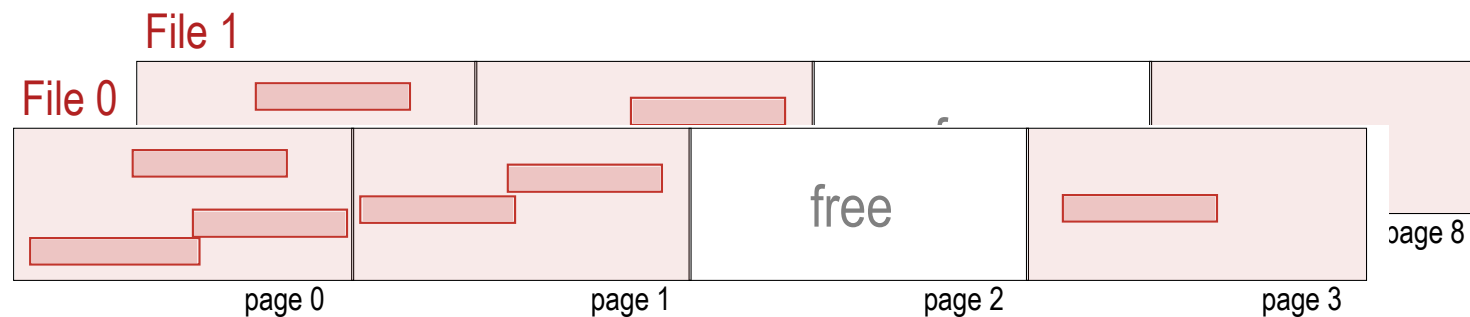
8

- Random Access (even though slow)
- Inexpensive
- Non-volatile
- Parts of an magnetic disk
 - Platter: covered with magnetic recording material
(turning)
 - Track: logical division of platter surface
 - Sector: hardware division of tracks
 - Block: OS division of tracks
Typical block sizes: 512B, 2KB, 4KB
 - Read/write head
(moving)

Files on Disk

9

- Metadata defines
 - Tables
 - Attributes
 - Data Types
- Stored are (data)
 - Logs
 - Records (== tuple)
 - Indices
- Data is stored in files
 - A file has one or more pages
 - A page contains of one or more records.



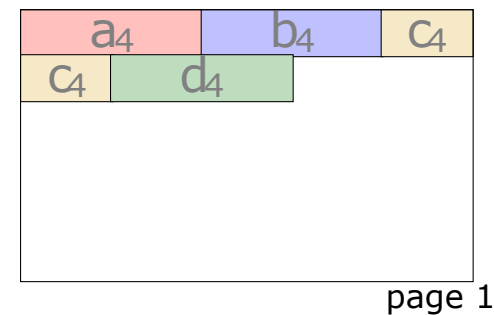
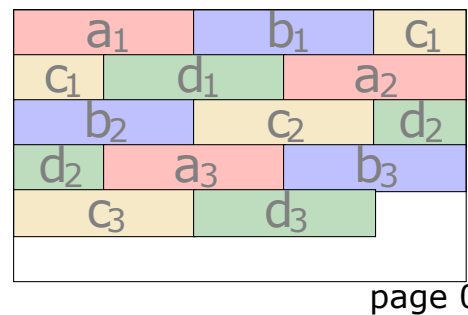
Rows, Columns, and the Page Layout

10

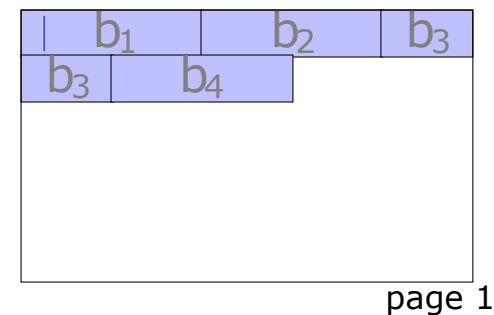
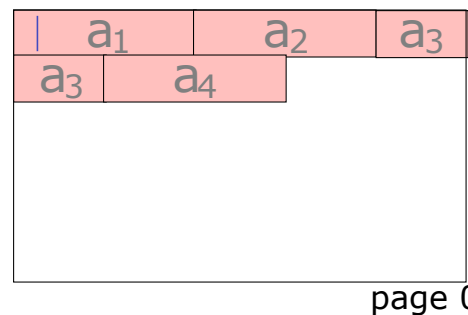
a ₁	b ₁	c ₁	d ₁
a ₂	b ₂	c ₂	d ₂
a ₃	b ₃	c ₃	d ₃
a ₄	b ₄	c ₄	d ₄



- **Row-oriented page layout** (n-ary storage model)



- **Column-oriented page layout** (decomposed storage model)

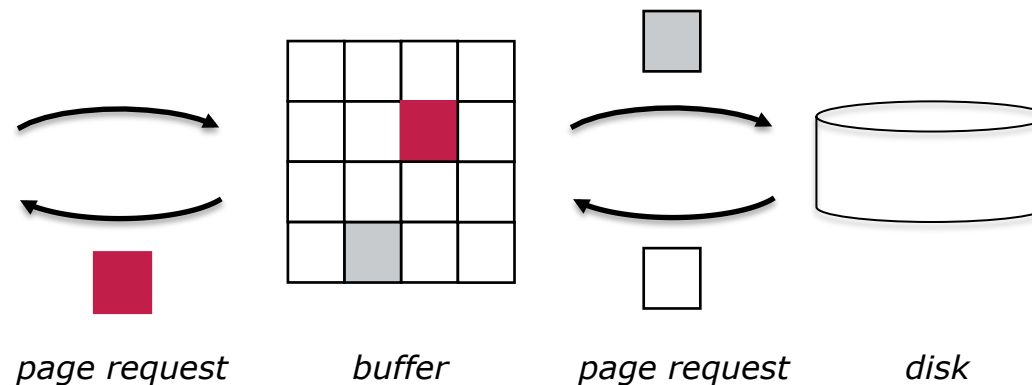


...

Buffer Management

11

- **Buffer** caches copies of pages in main memory
- Buffer Manager **maintains** these pages
 - Hit: requested page in buffer
 - Miss: page on disk
 - Allocate page frame
 - Read page
 - Page replacement
 - Dirty flag for write back
 - Least Recently Used (LRU)
 - Most Recently Used (MRU)



In a Nutshell

12

- Optimizations
 - Sequential Access
 - Buffering and scheduling algorithms
 - In-memory indices to pages
 - Pre-calculation and materialization
 - Etc.
- Page structure leads to
 - Good write performance
 - Efficient single tuple access
 - **Overhead** if single attributes scanned
 - regardless of disk throughput -

Question + Answer

13

Why is an in-memory database faster than a fully cached disk-based database?

- Disk access

- Low throughput
- Slow random access

- Buffer Management

- Disk-oriented data structures
(even in main memory)

- Page layout
- Indices

Question

14

Does this mean to keep data in main memory to achieve performance while the physical data representation can be neglected?

Why?

Memory Access

Jens Krueger

Enterprise Platform and Integration Concepts
Hasso Plattner Intitute

Capacity vs. Speed (latency)

16

Memory hierarchy:

- Capacity restricted by price/performance
- SRAM vs. DRAM (refreshing needed every 64ms)
- SRAM is very fast but very expensive

➔ Memory is organized in hierarchies

- Fast but small memory on the top
- Slow but lots of memory at the bottom

	technology	latency	size
CPU	SRAM	< 1 ns	bytes
L1 Cache	SRAM	~ 1 ns	KB
L2 Cache	SRAM	< 10 ns	MB
Main Memory	DRAM	100 ns	GB

Capacity vs. Speed (latency)

17

	latency	size
CPU	< 1 ns	bytes
L1 Cache	~ 1 ns	KB
L2 Cache	< 10 ns	MB
Main Memory	100 ns	GB
Magnetic Disk	~ 10 000 000 ns (10 ms)	TB

Data Processing

18

In DBMS, on disk as well as in memory, data processing is often:

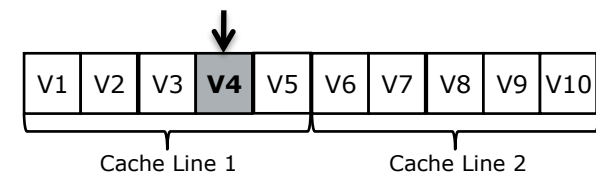
- Not CPU bound
- **But** bandwidth bound
- "I/O Bottleneck"

➔ CPU could process data faster

Memory Access:

- **Not** truly random (in the sense of constant latency)
- Data is read in **blocks**/cache lines
- Even if only parts of a block are requested

➔ Potential **waste** of bandwidth



Memory Basics I

19

- **Cache**

Small but fast memory, which keeps data from main memory for fast access.

→ Cache performance is **crucial**

- Similar to disk cache (e.g. buffer pool)

But: Caches are controlled by hardware.

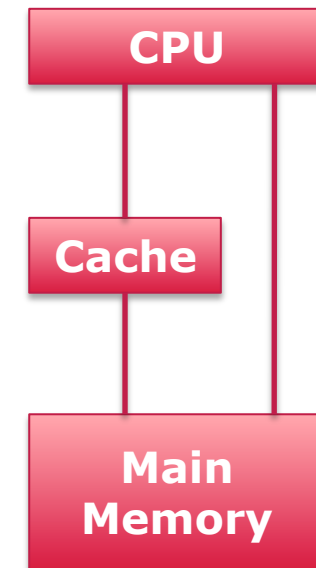
- **Cache hit**

Data was found in the cache.

Fastest data access since no lower level is involved.

- **Cache miss**

Data was **not** found in the cache. CPU has to load data from main memory into cache (**miss penalty**).



Memory Basics II

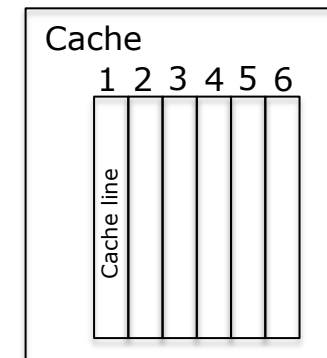
20

■ Cache lines

The cache is partitioned into lines.

- Data is read or written as whole line
- Size: 4-64 bytes

➔ Due to unnecessary data in cache lines the cache gets **polluted**.



Locality is King!

21

To improve cache behavior

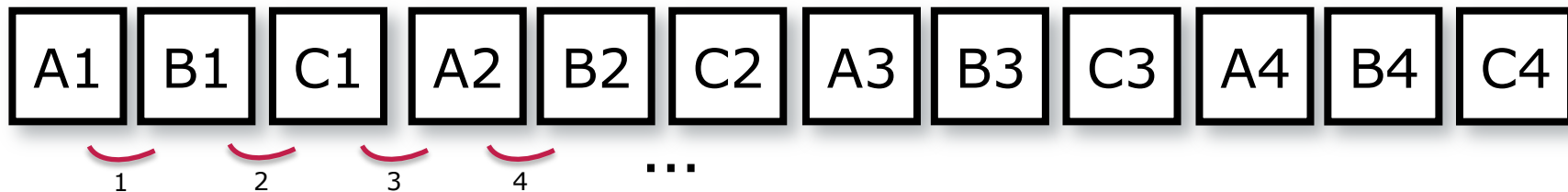
- Increase cache capacity
- Exploit locality
 - Spatial: related data is close (nearby references are likely)
 - Temporal: Re-use of data (repeat reference is likely)

To improve locality

- Non random access (e.g. scan, index traversal):
 - Leverage sequential access patterns
 - Clustering data to a cache lines
 - Partition to avoid cache line pollution (e.g. vertical decomposition)
 - Squeeze more operations into a cache line
- Random access (hash join):
 - Partition to fit in cache

Example for Sequential Access

23



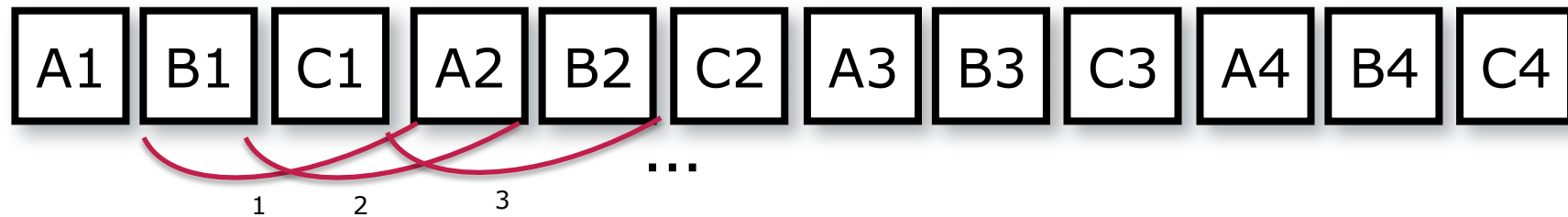
```
for (r = 0; r < rows; r++)  
    for (c = 0; c < columns; c++)  
        sum += table[r * columns + c];
```

Simulates sequential access

- All data in a cache line is read
- Prefetching and pipelining further **improve** performance

Example for Traversal Sequential Access

24



```
for (c = 0; c < columns; c++)
    for (r = 0; r < rows; r++)
        sum += table[c * columns + r];
```

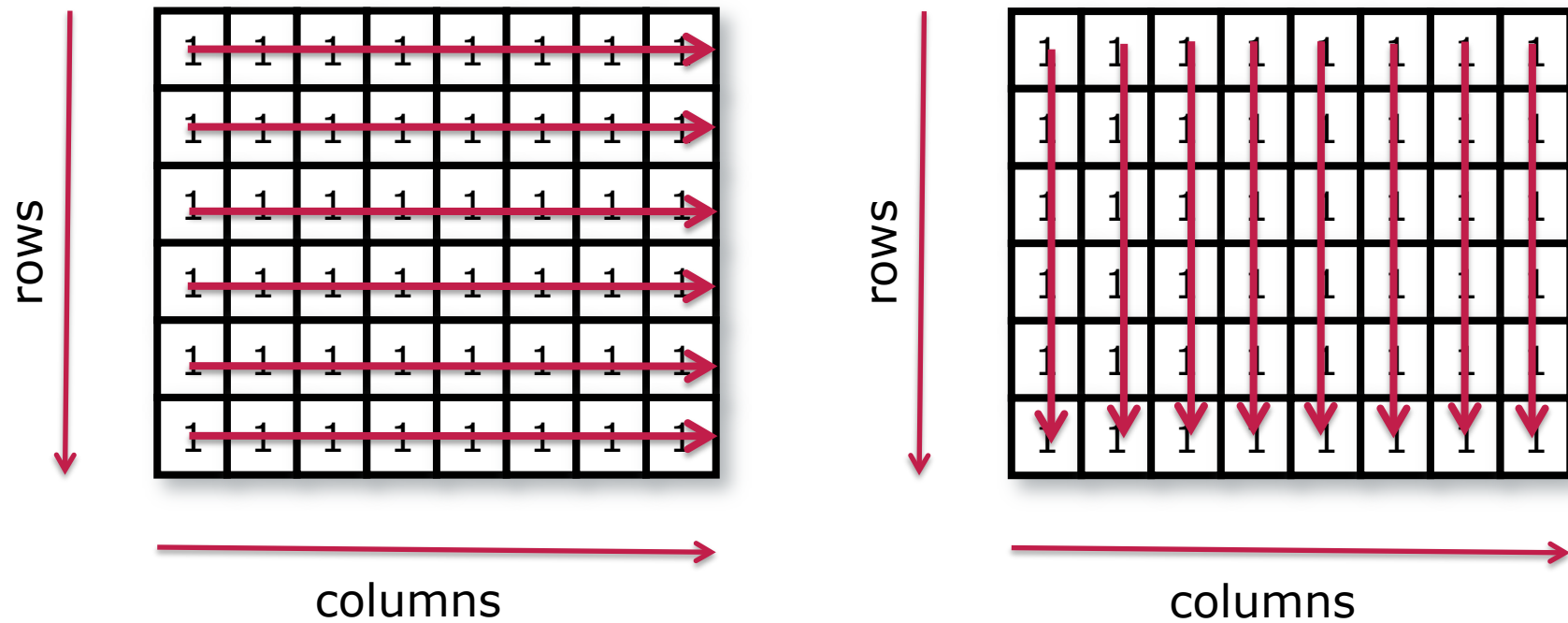
Simulates traversal sequential access

- Fixed stride (access offset) leads to cache misses
- Cache size / performance can be measured by varying the stride

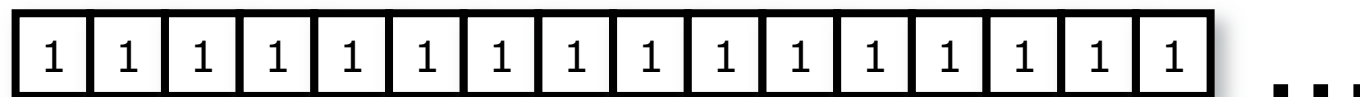
A Simple C++

25

- Logical



- Physical `int *table = (int*) calloc((rows * columns), sizeof(int));`



Demo

Demo C++ for Copy and Paste:

```

//=====
// Name      : msp02.cpp
// Author    : Jens
// Description: Aggregation
//=====
#include <sys/time.h>
#include <vector>
#include <iostream>

using namespace std;
#define C_NUMRUNS 1

typedef unsigned int uint;

void seq_read(unsigned int rows, unsigned int columns) {
    struct timeval start1, end1, start2, end2;
    long time;
    unsigned int r, c, table_size;
    int w;
    unsigned int seq_sum, seq2_sum, stride_sum;

    ///////////////////////////////////////////////////////////////////
    cout << "Fill table" << endl;

    int *table = (int*) malloc((rows * columns), sizeof(int));
    int* read = (int*) malloc(columns * sizeof(int));
    // füllen mit Random Int's
    for (r = 0; r < rows; r++)
        for (c = 0; c < columns; c++)
            table[r * columns + c] = (unsigned int) random() % 99999999;
    table_size = ((rows * columns) * sizeof(int)) / 1024 / 1024;
    cout << "Table: " << table_size
        << "MB" << endl;

    cout << "\nPress Key: ";
    cin >> w;

    ///////////////////////////////////////////////////////////////////
    cout << "Sequential Access" << endl;

    seq_sum = 0;
    time = 0;
    //lesen
    gettimeofday(&start1, NULL);
    for (r = 0; r < rows; r++)
        for (c = 0; c < columns; c++)
            //read[c] = table[r * columns + c];
            seq_sum += table[r * columns + c];
    gettimeofday(&end1, NULL);
    time = (end1.tv_sec - start1.tv_sec) * 1000000 + (end1.tv_usec
        - start1.tv_usec);
    cout << "Sum: " << seq_sum << endl;

    cout << "Time: " << time << "µsec " << (time / 1000.0) << "msec " <<
        (table_size / (time / 1000.0 / 1000.0)) << "MB/s" << endl;

    ///////////////////////////////////////////////////////////////////
    cout << "Stride Access" << endl;

    stride_sum = 0; time = 0;
    //lesen
    gettimeofday(&start2, NULL);
    for (c = 0; c < columns; c++)
        for (r = 0; r < rows; r++)
            //read[c] = table[r * columns + c];
            stride_sum += table[r * columns + c];
    gettimeofday(&end2, NULL);
    time = (end2.tv_sec - start2.tv_sec) * 1000000 + (end2.tv_usec
        - start2.tv_usec);
    cout << "Sum: " << stride_sum << endl;

    cout << "Time: " << time << "µsec " << (time / 1000.0) << "msec " <<
        (table_size / (time / 1000.0 / 1000.0)) << "MB/s" << endl;

    free(table);
    free(read);
}

/////////////////////////////////////////////////////////////////
int main(int argc, char* argv[]) {
    unsigned int rows = 3000000;
    unsigned int columns = 300;

    seq_read(rows, columns);

    cout << "##### Finish" << endl;

    return 0;
}

```

In-Memory Databases

Jens Krueger

Enterprise Platform and Integration Concepts
Hasso Plattner Intitute

In-Memory Database

28

In an In-Memory Database (IMDB)

- Data resides **permanently** in main memory
- Main Memory is the **primary** "*persistence*"
- Still: logging to **disk**/recovery from **disk**
- Main memory access is the new **bottleneck**
- Cache-conscious algorithms/data structures are **crucial**
(locality is king)

Today's Main Memory Technology

- Increased size: up to 2 TB of main memory on one main board as of today
- Increased bandwidth: up 30GB/s
- Latency is hidden by caches (memory hierarchy)

In-Memory Database

29

In an In-Memory Database (IMDB)

- Data resides **permanently** in main memory
- Main Memory is the **primary** "*persistence*"
- Still: logging to **disk**/recovery from **disk**
- Main memory access is the new **bottleneck**
- Cache-conscious algorithms/data structures are **crucial**
(locality is king)

Differences to disk-based systems

- Volatile
- Direct access
- Access time
- Access cost

Question

30

Does an entire database fit in main memory?

Question + Answer

31

Does an entire database fit in main memory?

- Yes:
 - Limited DB size, i.e. enterprise applications
 - Due to data compression (factor 10 feasible)
 - Redundant-free data schemas
- No:
 - Data could be partitioned over nodes
 - Data aging strategies for extended memory hierarchies (e.g. SSD/disks for non-active data)

More Main Memory for Disk-based DBMS?

32

What is the difference between an IMDB and a disk-based DB with a large cache?

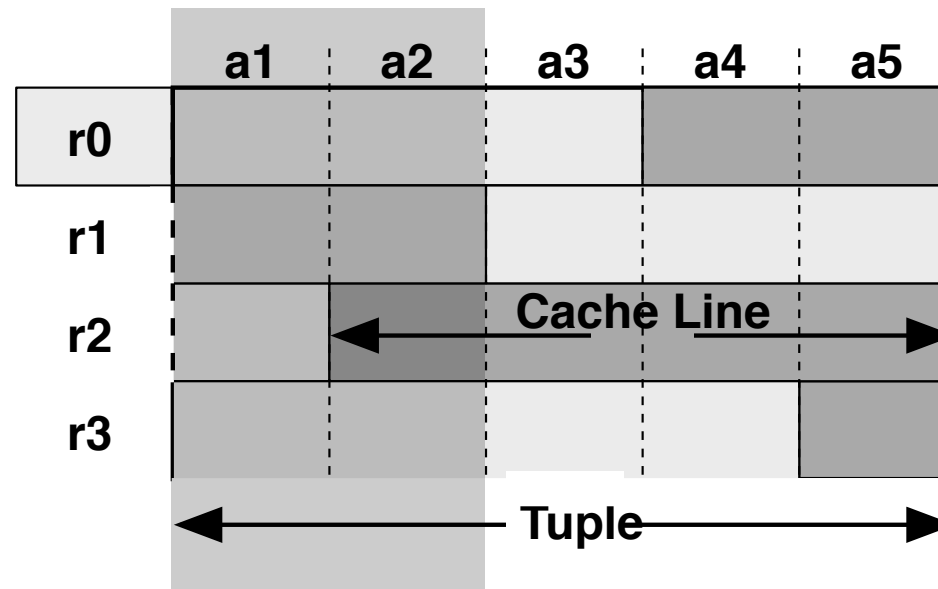
- Different optimizations for data structures, e.g.
 - Page layout
 - No access through a buffer manager
 - Index structures
 - Cache-aware data organization
 - Random access capabilities, e.g. for locking
- As disk-based DB's can have in-memory optimization, they still would have to maintain different data structures.

IMDB: Relations and Cache Lines

33

The physical data layout with regards to the workload has a significant influence on the cache behavior of the IMDB.

- Tuples are spanned over cache lines
- Wrong layout can lead to lots of (expensive) cache misses
- Row- or column-oriented can reduce cache misses if matching workload is applied



Question

34

How to optimize an IMDB?

Question + Answer

35

How to optimize an IMDB?

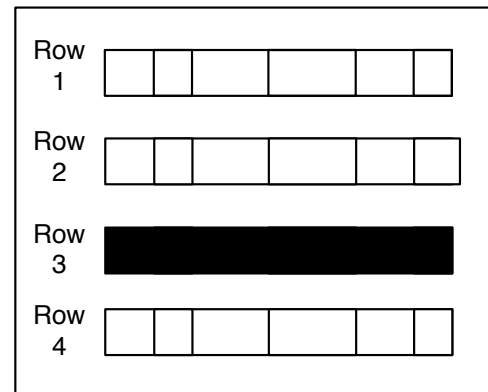
- Exploit sequential access
- Leverage locality

Row- or Column-oriented Storage

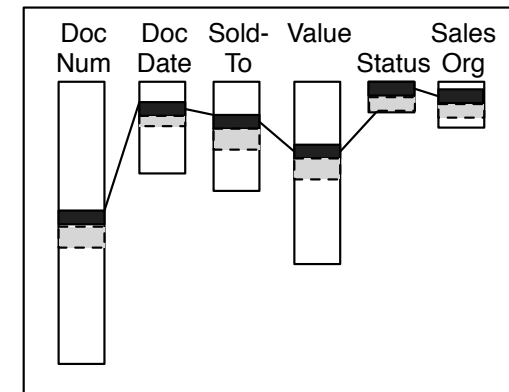
36

```
SELECT *
FROM Sales Orders
WHERE Document Number = '95779216'
```

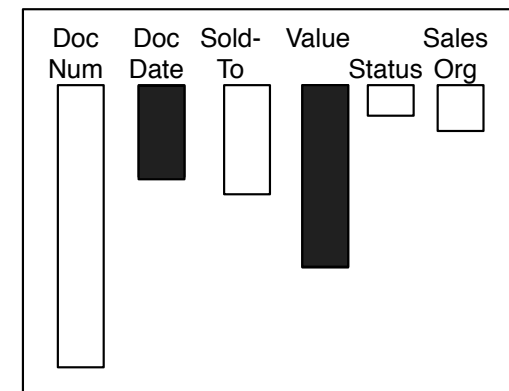
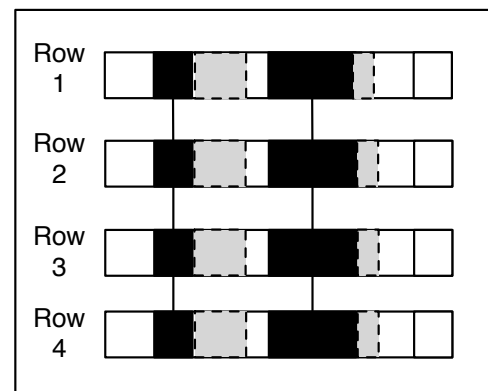
Row Store



Column Store



```
SELECT SUM(Order Value)
FROM Sales Orders
WHERE Document Date > 2009-01-20
```



Row-oriented storage

37

A1	B1	C1
A2	B2	C2
A3	B3	C3
A4	B4	C4

Row-oriented storage

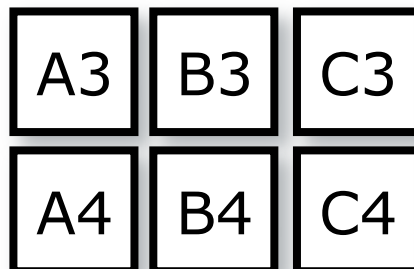
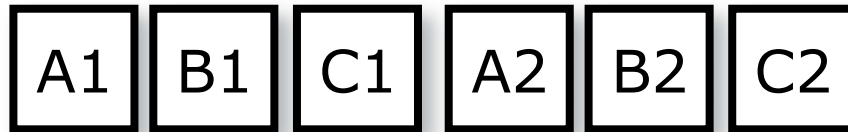
38

A1	B1	C1
----	----	----

A2	B2	C2
A3	B3	C3
A4	B4	C4

Row-oriented storage

39



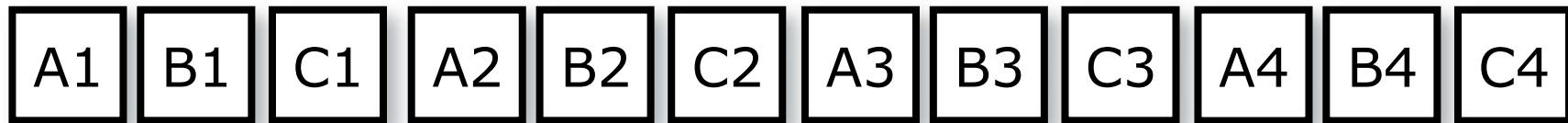
Row-oriented storage

40



Row-oriented storage

41



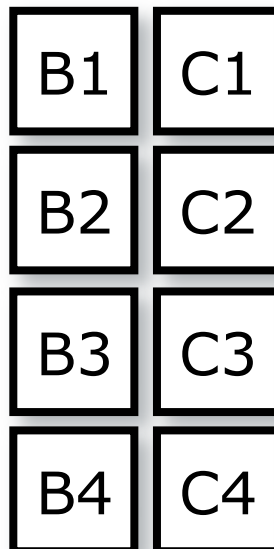
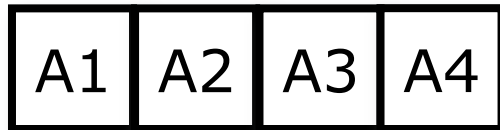
Column-oriented storage

42

A1	B1	C1
A2	B2	C2
A3	B3	C3
A4	B4	C4

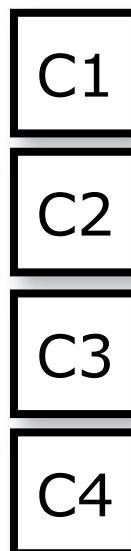
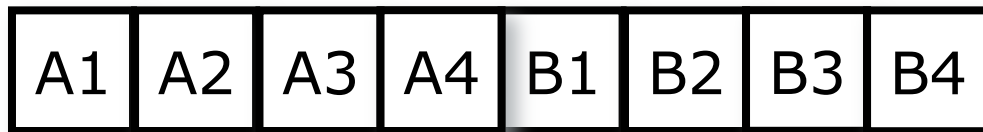
Column-oriented storage

43



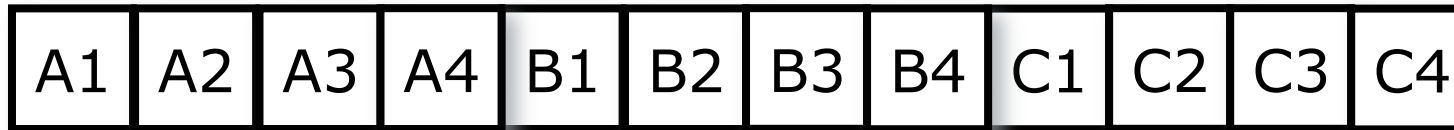
Column-oriented storage

44



Column-oriented storage

45



Example: OLTP-Style Query

46

```
struct Tuple {  
  int a,b,c;  
};
```

```
Tuple data[4];  
fill(data);
```

```
Tuple third = data[3];
```

A1	B1	C1
A2	B2	C2
A3	B3	C3
A4	B4	C4

Example: OLTP-Style Query

47

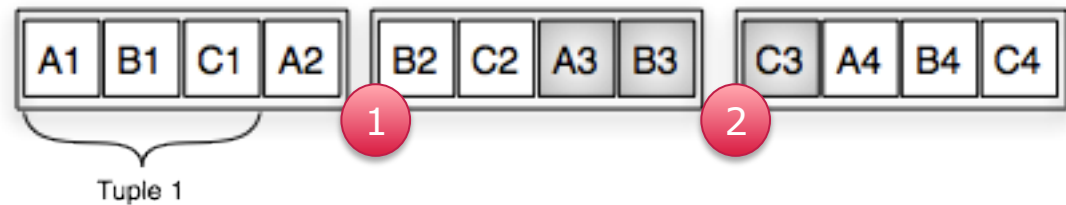
```
struct Tuple {
int a,b,c;
};
```

```
Tuple data[4];
fill(data);
```

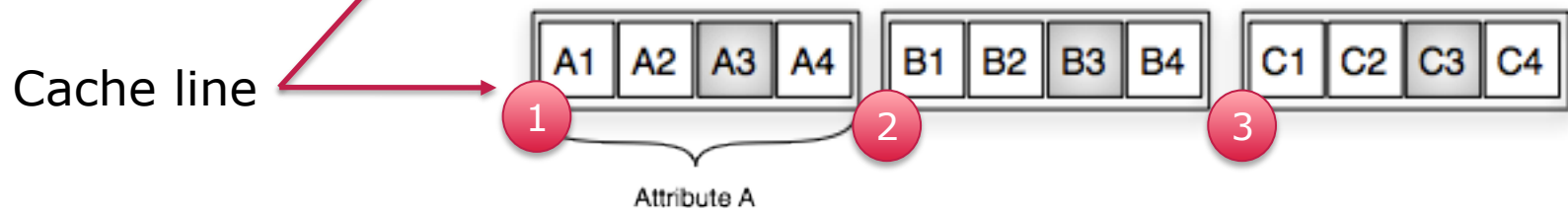
```
Tuple third = data[3];
```

A1	B1	C1
A2	B2	C2
A3	B3	C3
A4	B4	C4

Row Oriented Storage



Column Oriented Storage



Example: OLAP-Style Query

48

```
struct Tuple {  
    int a,b,c;  
};  
  
Tuple data[4];  
fill(data);  
  
int sum = 0;  
  
for(int i = 0;i<4;i++)  
  
    sum += data[i].a;
```

A1	B1	C1
A2	B2	C2
A3	B3	C3
A4	B4	C4

Example: OLAP-Style Query

49

```
struct Tuple {
int a,b,c;
};
```

```
Tuple data[4];
fill(data);
```

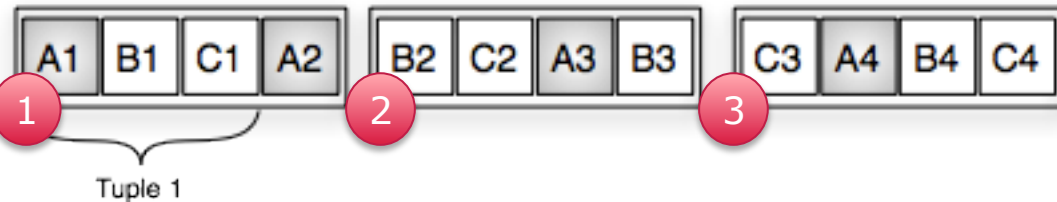
```
int sum = 0;
```

```
for(int i = 0;i<4;i++)
```

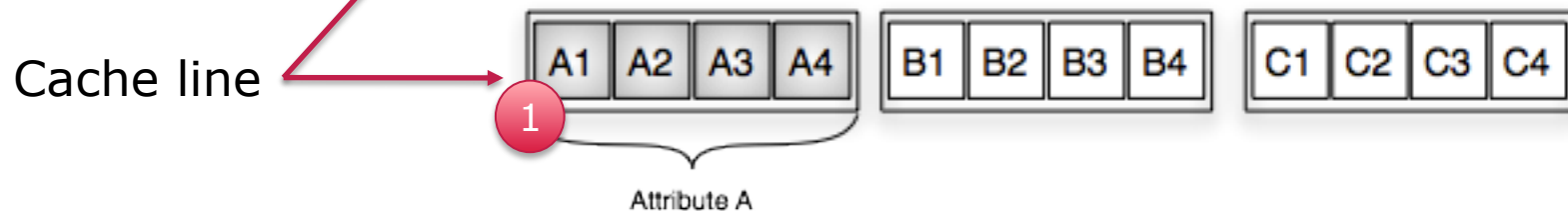
```
sum += data[i].a;
```

A1	B1	C1
A2	B2	C2
A3	B3	C3
A4	B4	C4

Row Oriented Storage



Column Oriented Storage



Mixed Workloads

50

- Mixed Workloads involve attribute- and entity-focused queries

OLTP-style queries

A1	B1	C1
A2	B2	C2
A3	B3	C3
A4	B4	C4

OLAP-style queries

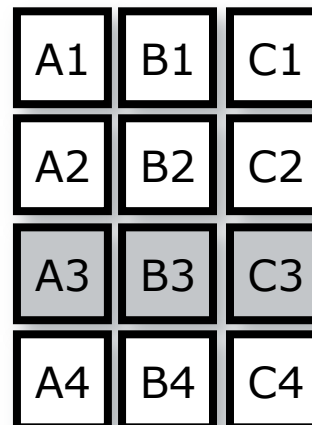
A1	B1	C1
A2	B2	C2
A3	B3	C3
A4	B4	C4

Mixed Workloads: Choosing the Layout

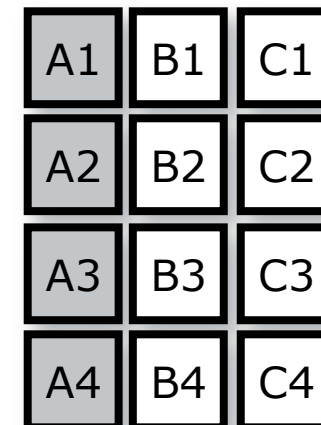
51

Layout	OLTP-Misses	OLAP-Misses	Mixed
Row	2	3	5
Column	3	1	4

OLTP-style queries



OLAP-style queries



Question

52

What is the best layout for mixed workloads?

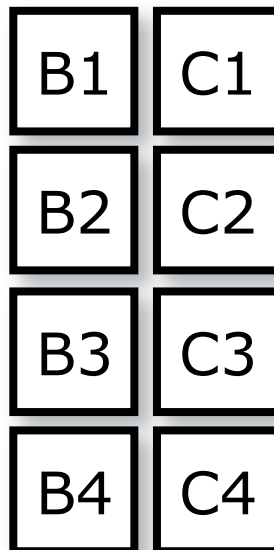
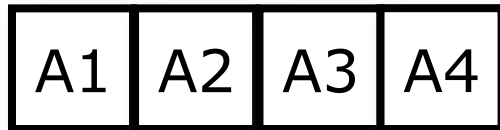
Hybrid: Grouping of Columns

53

A1	B1	C1
A2	B2	C2
A3	B3	C3
A4	B4	C4

Hybrid: Grouping of Columns

54



Hybrid: Grouping of Columns

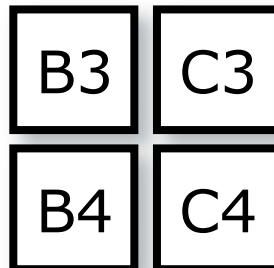
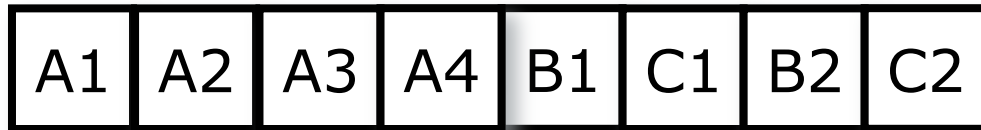
55

A1	A2	A3	A4	B1	C1
----	----	----	----	----	----

B2	C2
B3	C3
B4	C4

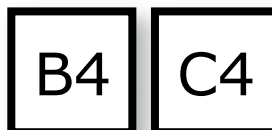
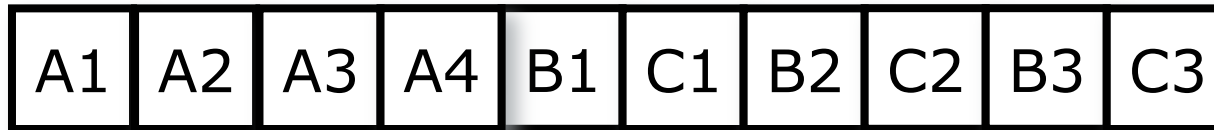
Hybrid: Grouping of Columns

56



Hybrid: Grouping of Columns

57



Hybrid: Grouping of Columns

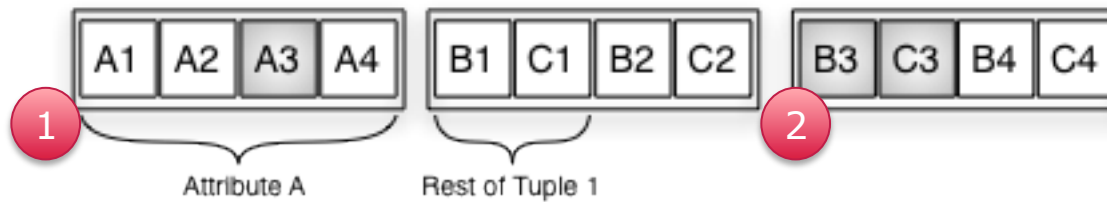
58

A1	A2	A3	A4	B1	C1	B2	C2	B3	C3	B4	C4
----	----	----	----	----	----	----	----	----	----	----	----

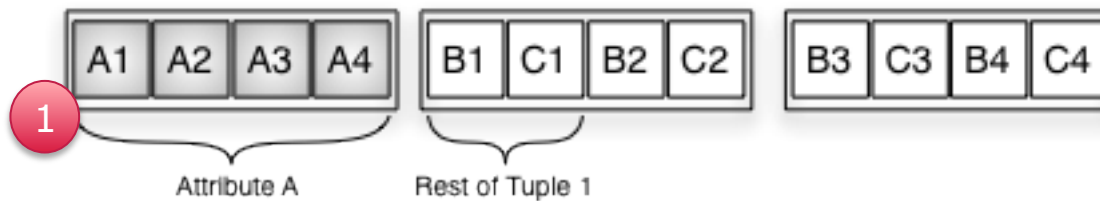
Hybrid: Grouping of Columns

59

Access tuple 3



Query attribute A



Layout	OLTP-Misses	OLAP-Misses	Mixed
Row	2	3	5
Column	3	1	4
Hybrid	2	1	3

Question

60

What other optimization for an
IMDB?

Compression in In-Memory Databases

Jens Krueger

Enterprise Platform and Integration Concepts
Hasso Plattner Intitute

Motivation

62

- Main memory is the new bottleneck
- Processor speed increases faster than memory speed
- Trade CPU time to compress and decompress data
- Compression
 - **Reduces** I/O operations to main memory
 - Leads to **less** cache misses due to more information on a cache line
 - Enables operations **directly** on compressed data
 - Allows to **offset** by the use of fixed-length data types

Compression Techniques

63

- Lightweight compression techniques:
 - **Lossless**
 - Reduce the amount of data
 - Improve query execution
 - Better utilizes cache lines
 - Techniques
 - Run Length Encoding
 - Null Suppression
 - Bit Vector Encoding
 - Dictionary Encoding

Run Length Encoding (RLE)

64

- Subsequent equal values are stored as one value with offset (value, run_length)
- Especially useful for sorted columns
- But:
 - If column store works with TupleId, only sorting by one column is possible

Null Suppression

65

- Remove leading 0's
- Most effective when encoding random sequence of small integers
 - `int x = 7;` uses 32 bits but first 29 are 0's
 - store (length, encoding) => (3, 111)
- Optimization: store byte count for next 4 values as two bits in one byte

Bit vector encoding

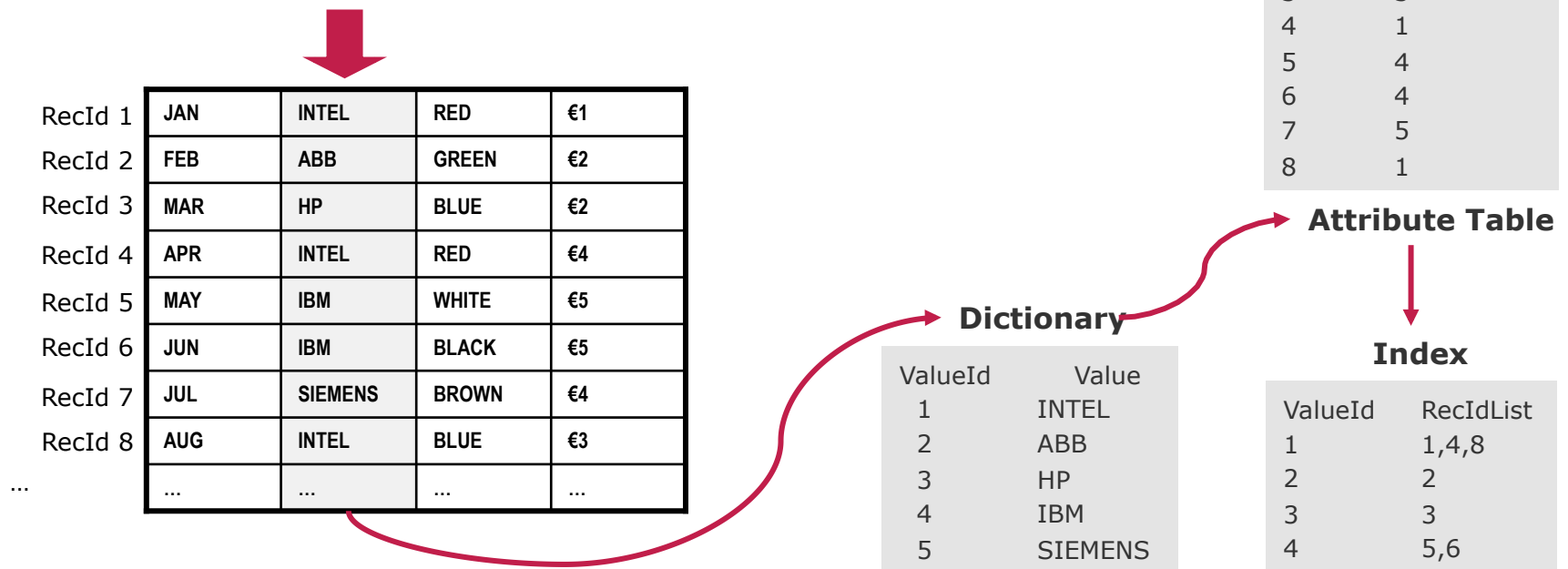
66

- Store a bitmap for each distinct value
- Values to encode: a b a a c c b
 - a => (1 0 1 1 0 0 0)
 - b => (0 1 0 0 0 0 1)
 - c => (0 0 0 0 1 1 0)
- Useful with few distinct values

Dictionary Encoding

67

- Store distinct values once in separate mapping table (the dictionary)
- Associate unique mapping key for each distinct value
- Store mapping key instead of value in value table



Example (1)

68

- Store fixed length strings of 32 characters
 - SQL-Speak: CHAR(32) - 32 Bytes
 - 1 Million entries consume $32 * 10^6$ Bytes
 - \sim 32 Megabytes

Example (2)

69

- Associate 4 byte valueID with distinct value
- Dictionary: assume 200.000 distinct values
 - Each: 1 key, 1 value => 36 Bytes
 - ~ 7.2 Megabytes
 - 1 million * 4 Bytes = ~ 4 Megabytes
- Overall: 11.2 Megabytes
- 64 byte cache line
 - Uncompressed: 2 values per cache line
 - Compressed: 16 valueID's per cache line

Question

70

How can this compression technique further be improved?

With regards to:

- **Amount** of data
- Query **execution**

Answer

71

- Amount of data
 - Idea: compress valueID's
 - Use only bits needed to represent the cardinality of distinct values - $\log_2(\text{distinct values})$
 - Optimal for only a few distinct values
 - Re-encoding if more bits to encode needed
- Query execution
 - Use order-preserving dictionaries
 - ValueID's have same order as uncompressed values
 - $\text{value1} < \text{value2} \iff \text{valueID1} < \text{valueID2}$

Materialization in Column Stores

Jens Krueger

Enterprise Platform and Integration Concepts
Hasso Plattner Intitute

Strategies for Tuple Reconstruction

73

Strategies:

- **Early** materialization
Create a row-wise data representation
at the first operator
- **Late** materialization
Operate on columns as long as possible

Example:

74

Query:

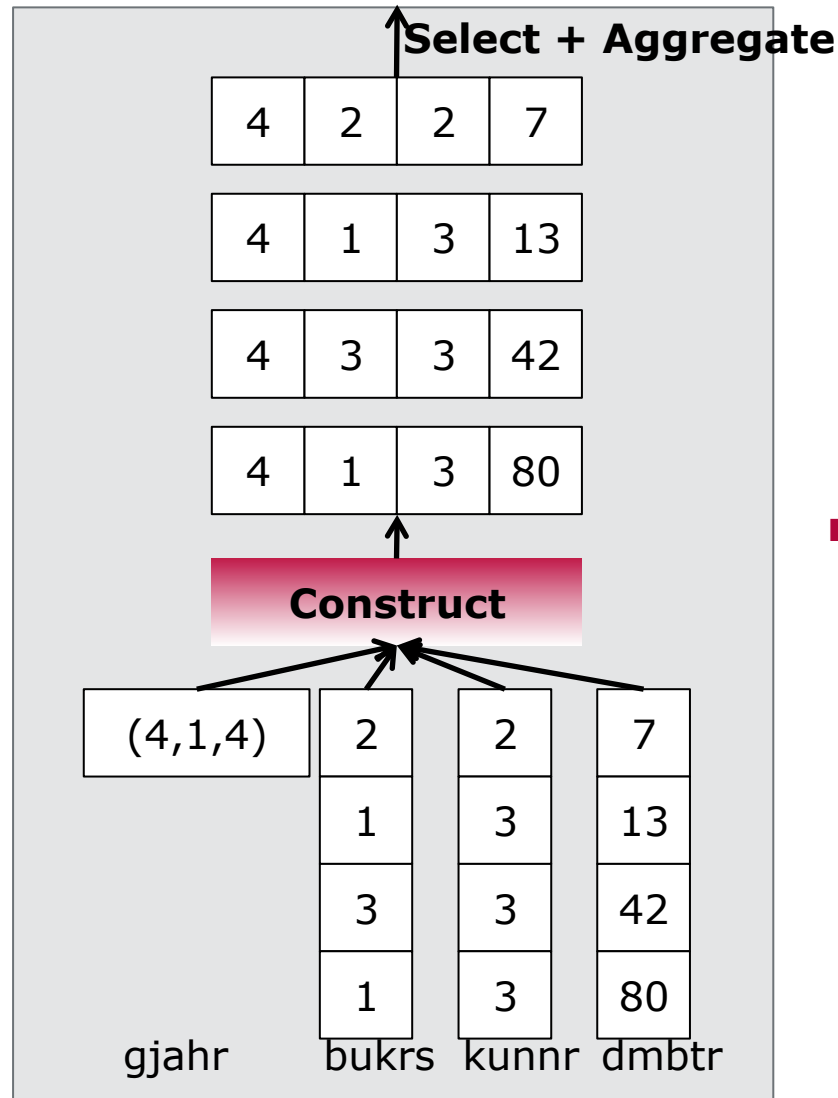
```
SELECT kunnr, sum(dmbtr)
FROM BSEG
WHERE   gjahr = 4
AND     bukrs = 1
GROUP BY kunnr
```

Table BSEG

4	2	2	7
4	1	3	13
4	3	3	42
4	1	3	80
gjahr	bukrs	kunnr	dmbtr

Early materialization

75



Query:

```
SELECT kunnr, sum(dmbtr)
FROM BSEG
WHERE   gjahr = 4
AND     bukrs = 1
GROUP BY kunnr
```

■ Create rows first

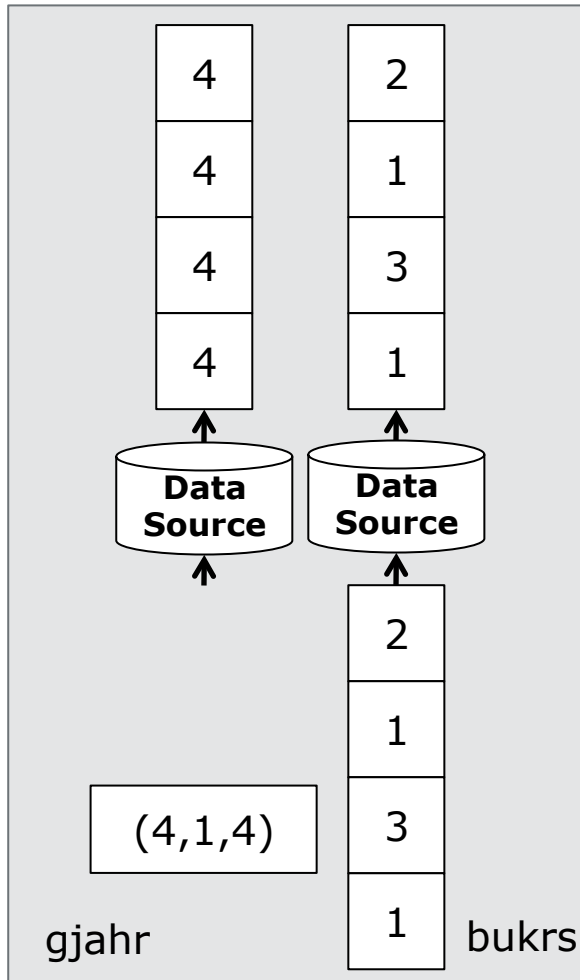
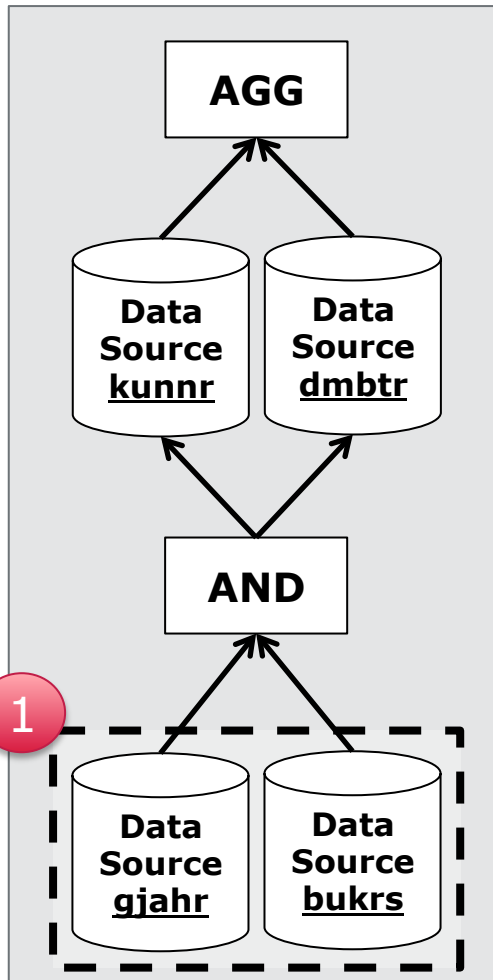
But:

- Need to construct **ALL** tuples
- Need to decompress data
- Poor memory bandwidth utilization

Late materialization I

76

Operate on columns



Query:

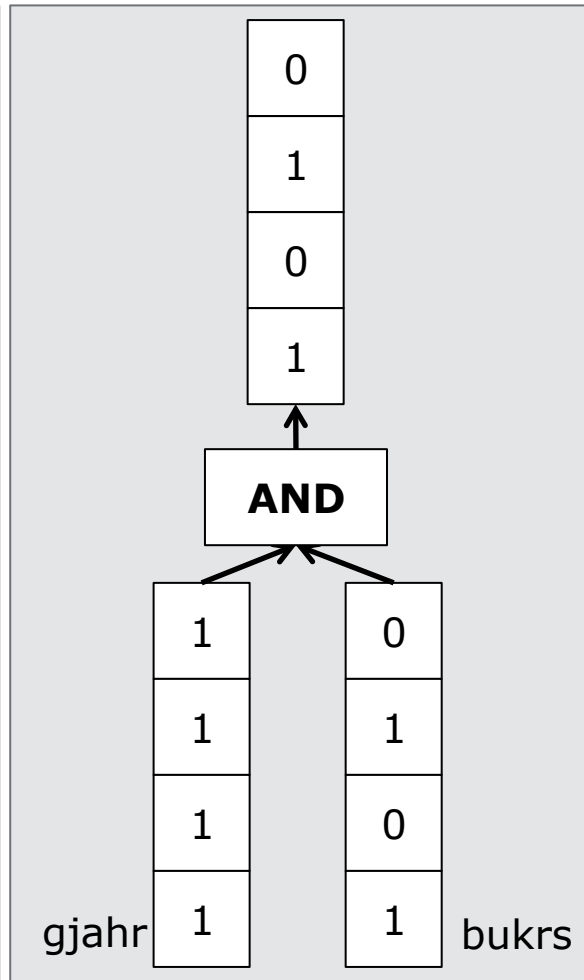
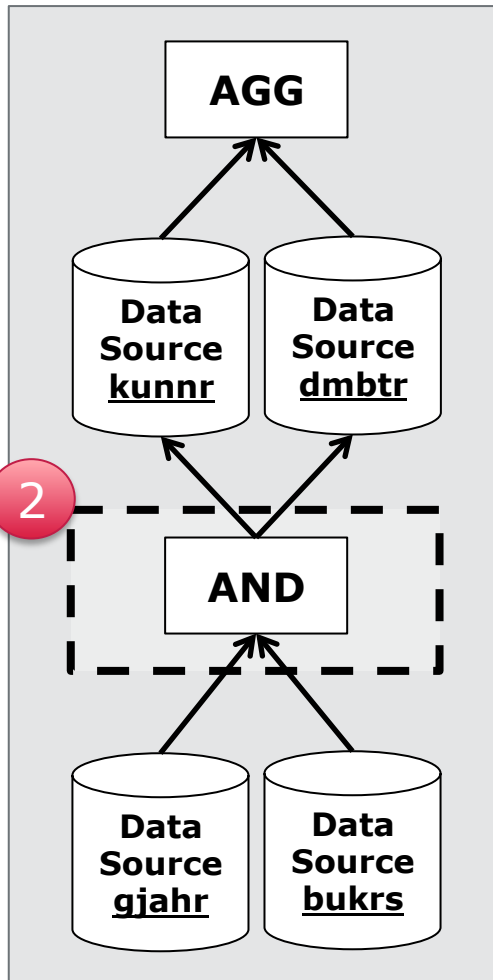
```
SELECT kunnr, sum(dmbtr)
FROM BSEG
WHERE   gjahr = 4
AND     bukrs = 1
GROUP BY kunnr
```

4	2	2	7
4	1	3	13
4	3	3	42
4	1	3	80
gjahr	bukrs	kunnr	dmbtr

Late materialization II

77

Operate on columns



Query:

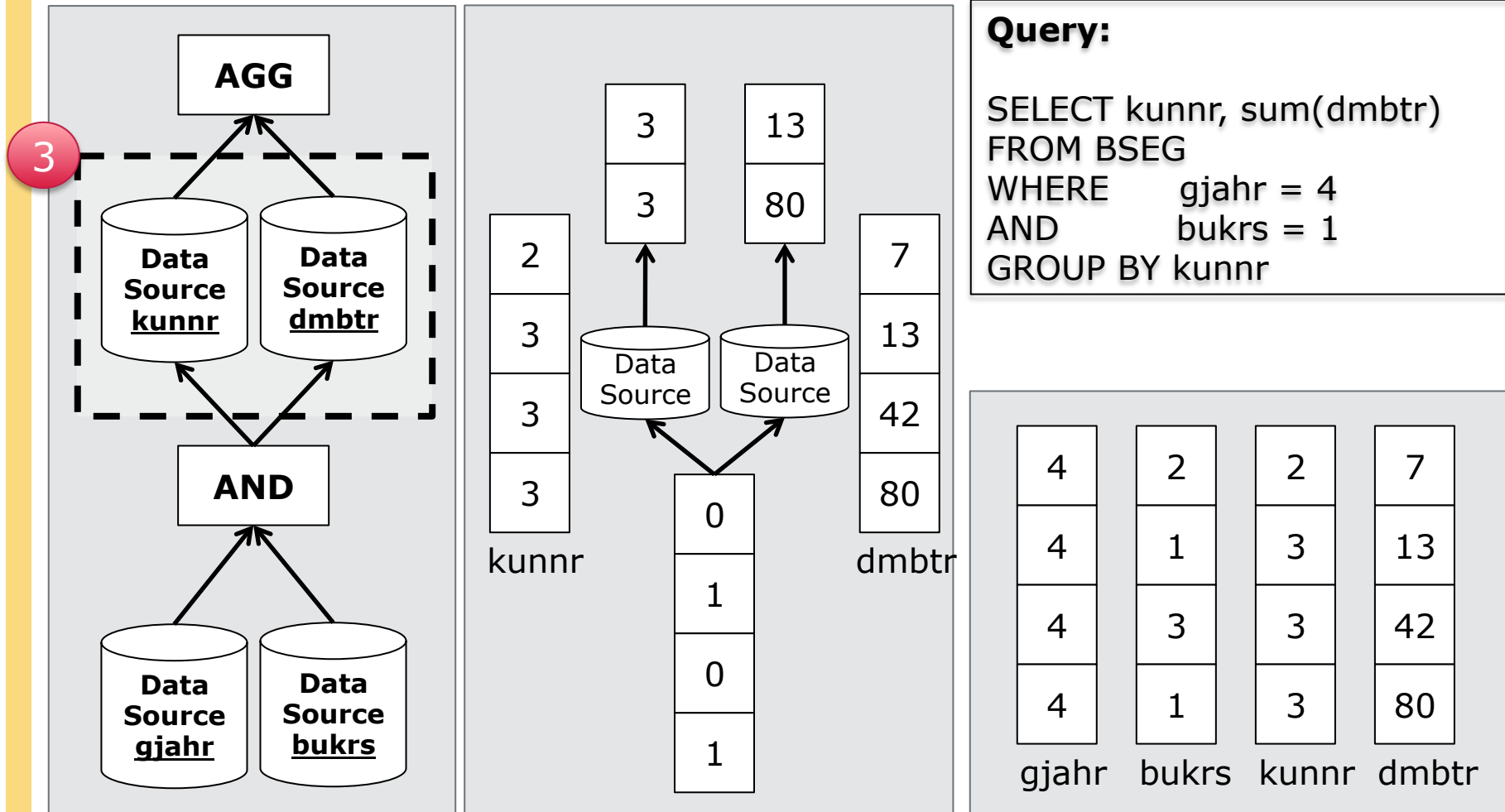
```
SELECT kunnr, sum(dmbtr)
FROM BSEG
WHERE   gjahr = 4
AND     bukrs = 1
GROUP BY kunnr
```

4	2	2	7
4	1	3	13
4	3	3	42
4	1	3	80
gjahr	bukrs	kunnr	dmbtr

Late materialization III

78

Operate on columns



Late materialization IV

79

Operate on columns

