



香港中文大學(深圳)
The Chinese University of Hong Kong, Shenzhen



Ack: Prof. Jignesh Patel @ CMU
Prof. Andy Pavlo @CMU

CSC3170

5: Storage Model

Chenhao Ma

School of Data Science

The Chinese University of Hong Kong, Shenzhen

Database Workloads

- **On-Line Transaction Processing (OLTP)**

→ Fast operations that only read/update a small amount of data each time.

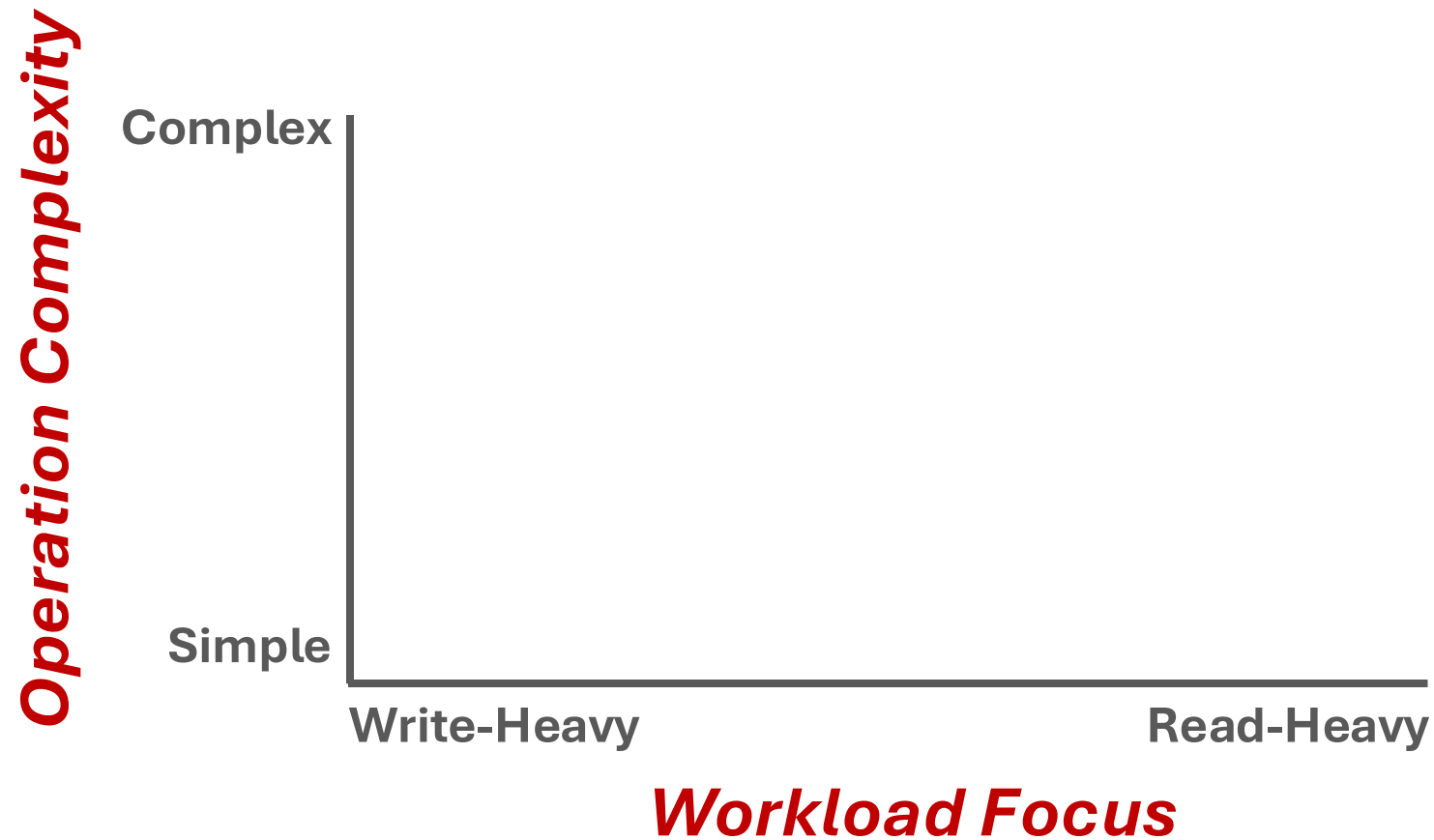
- **On-Line Analytical Processing (OLAP)**

→ Complex queries that read a lot of data to compute aggregates.

- **Hybrid Transaction + Analytical Processing**

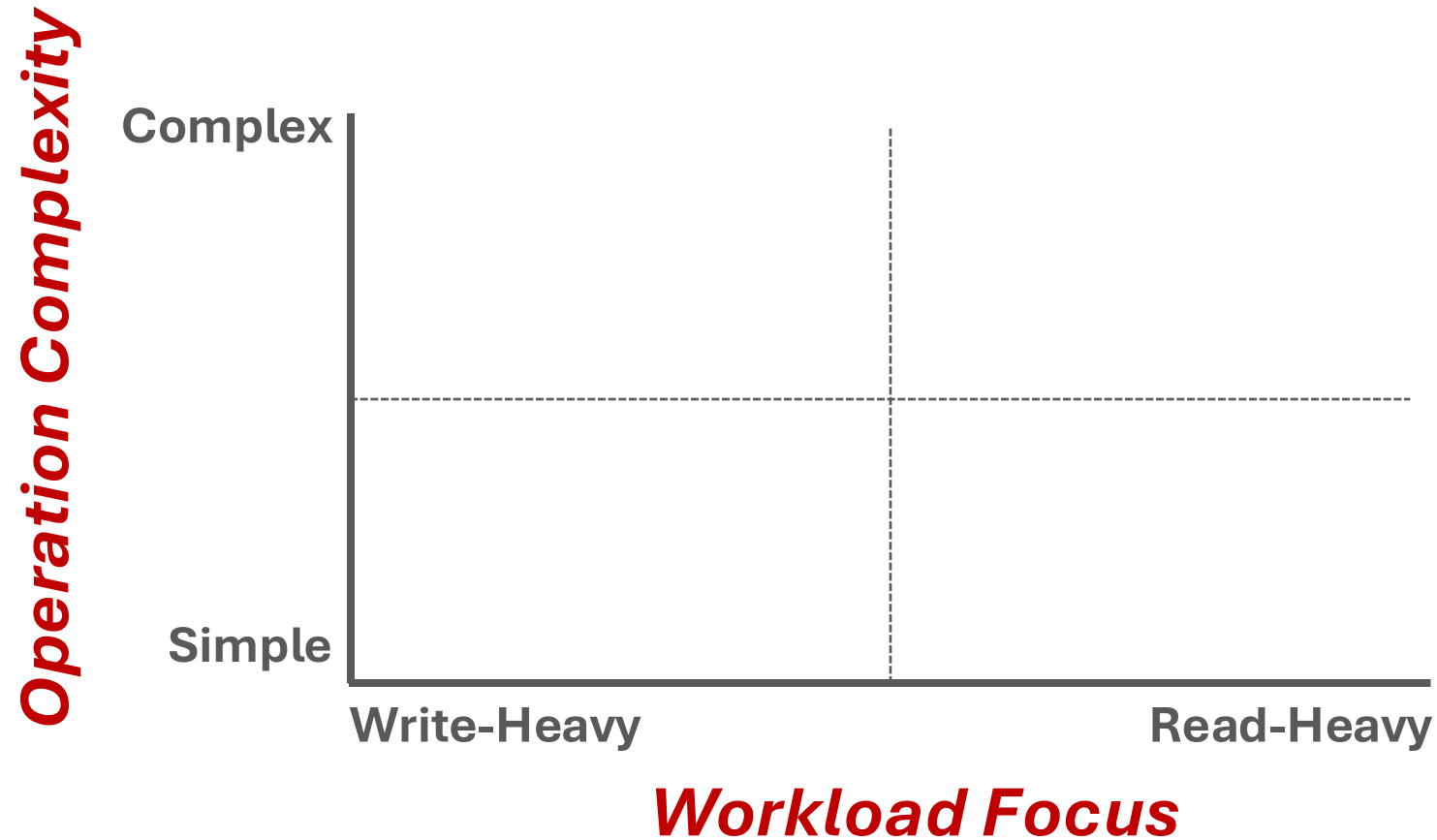
→ OLTP + OLAP together on the same database instance

Database Workloads



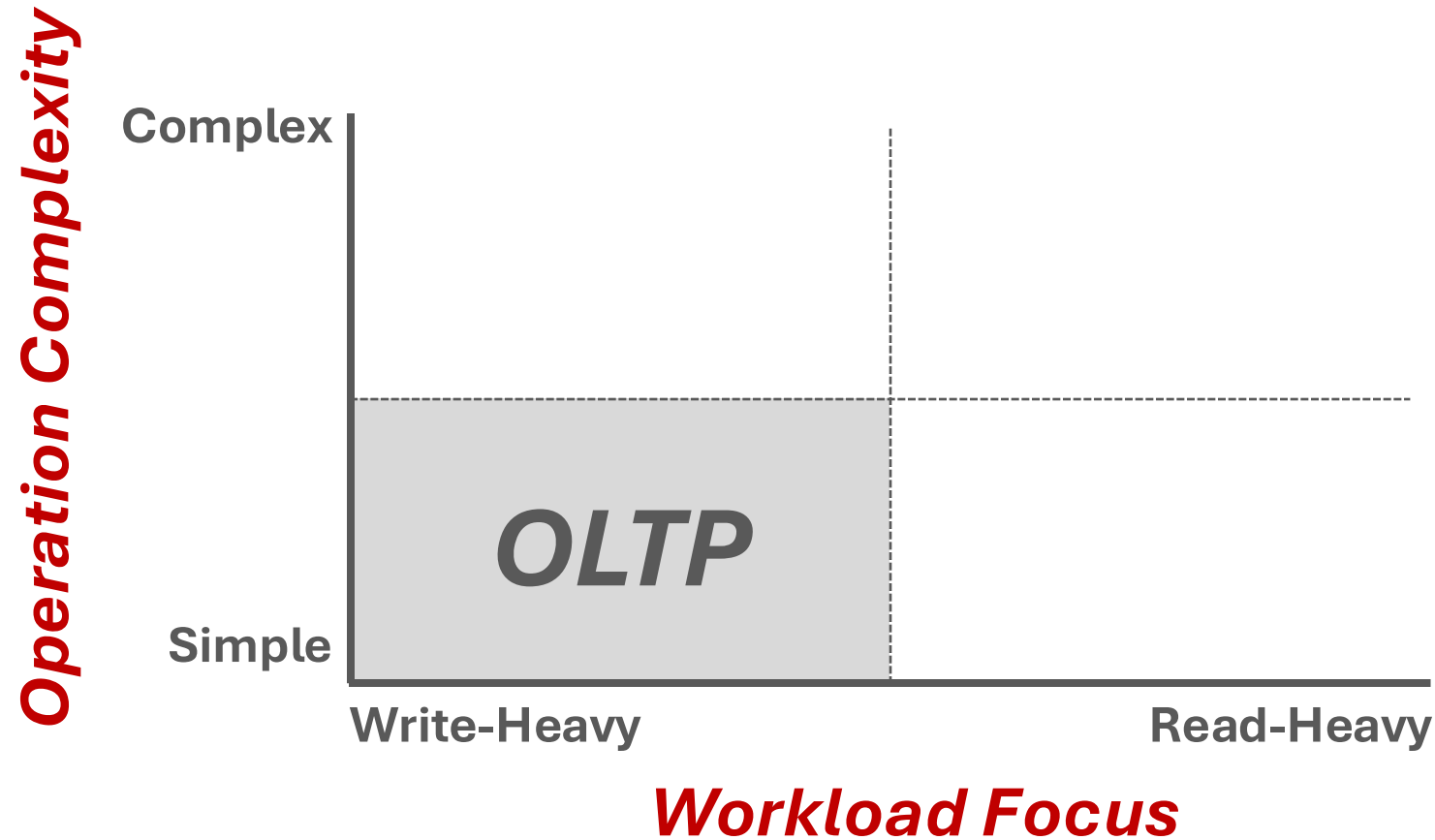
Source: [Mike Stonebraker](#)

Database Workloads



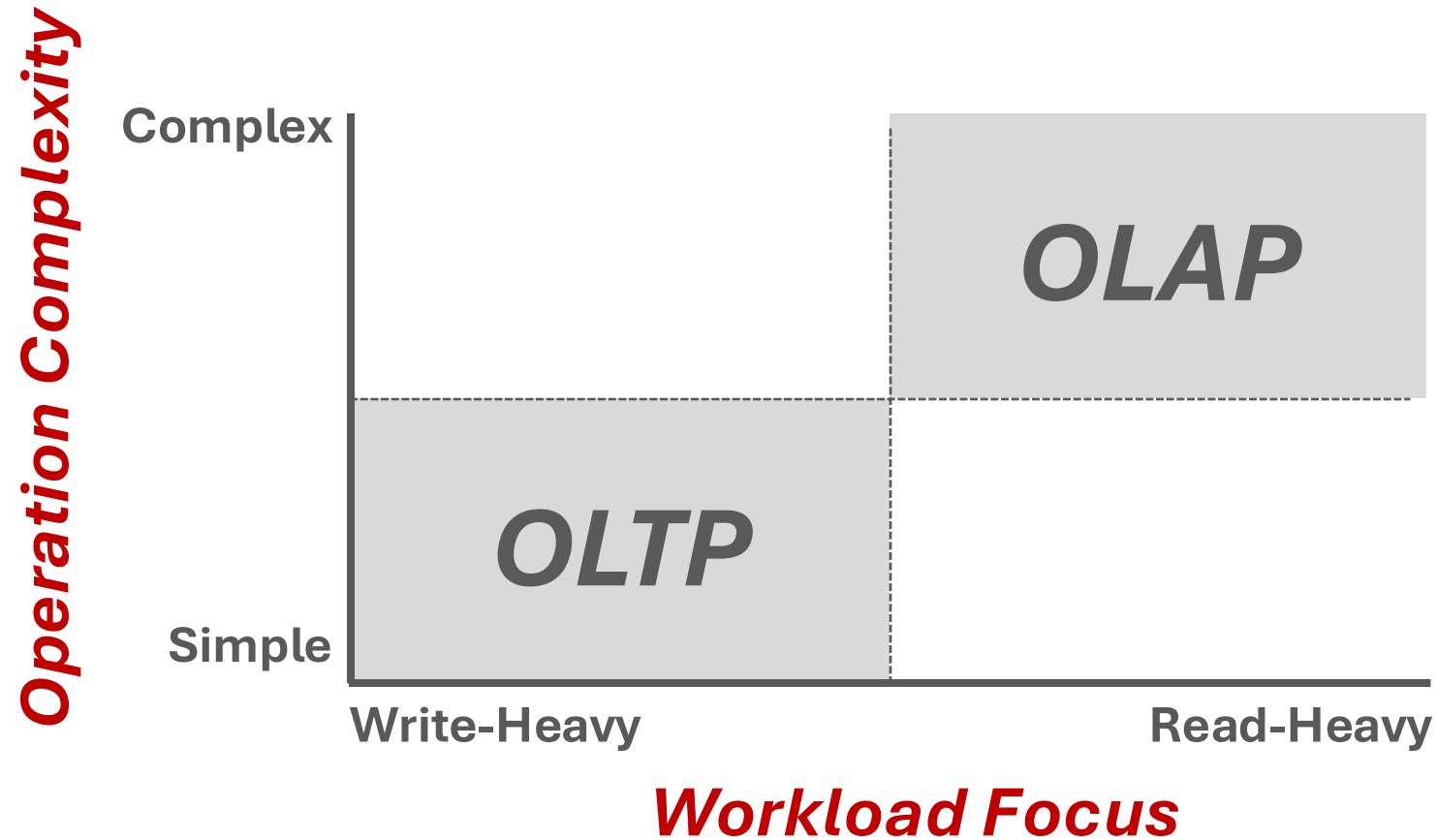
Source: [Mike Stonebraker](#)

Database Workloads



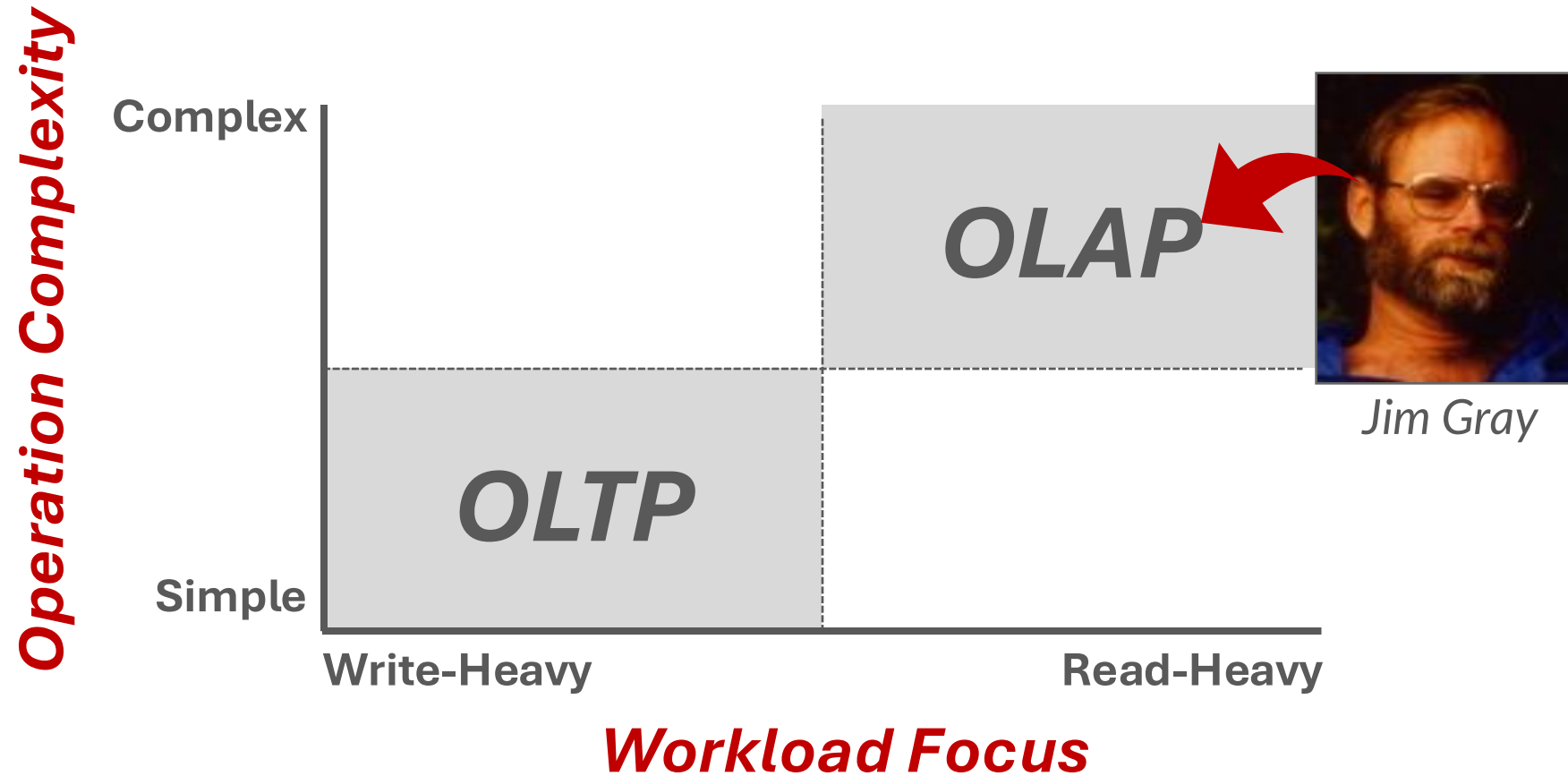
Source: [Mike Stonebraker](#)

Database Workloads



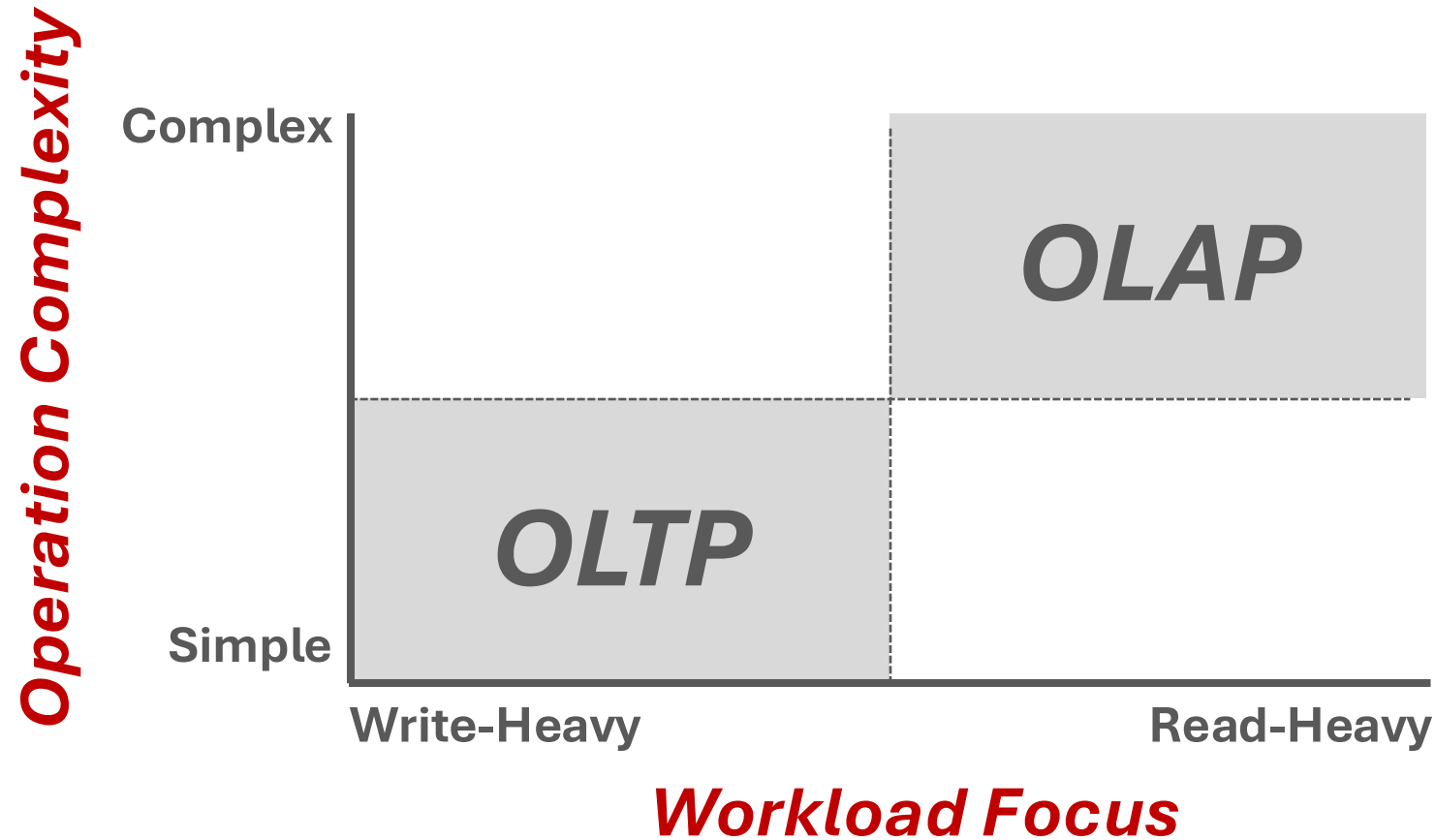
Source: [Mike Stonebraker](#)

Database Workloads



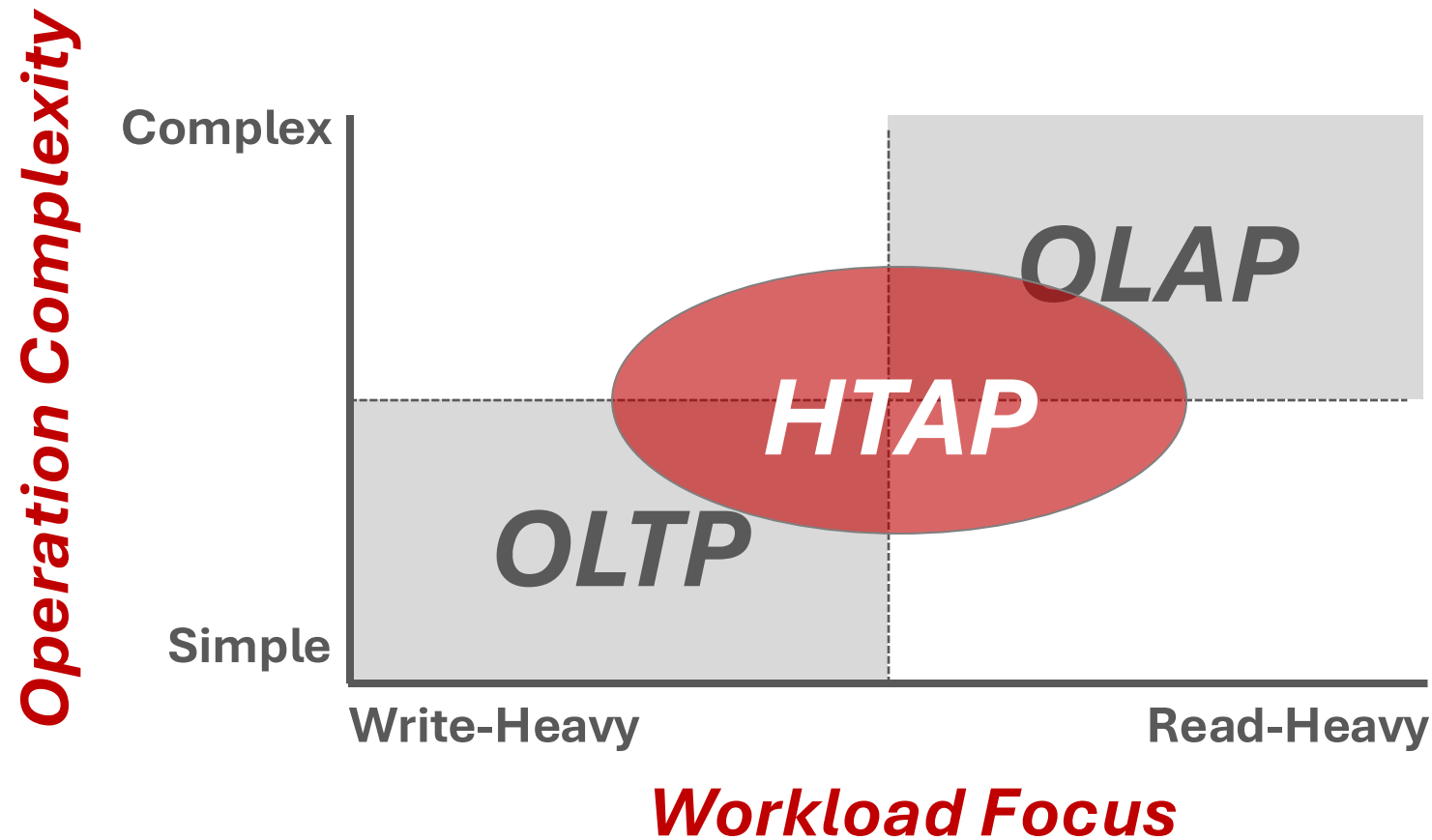
Source: [Mike Stonebraker](#)

Database Workloads



Source: [Mike Stonebraker](#)

Database Workloads



Source: [Mike Stonebraker](#)

Wikipedia Example

```
CREATE TABLE useracct (  
    userID INT PRIMARY KEY,  
    userName VARCHAR UNIQUE,  
    :  
);
```

```
CREATE TABLE pages (  
    pageID INT PRIMARY KEY,  
    title VARCHAR UNIQUE,  
    latest INT  
    ↳ REFERENCES revisions (revID),  
);
```

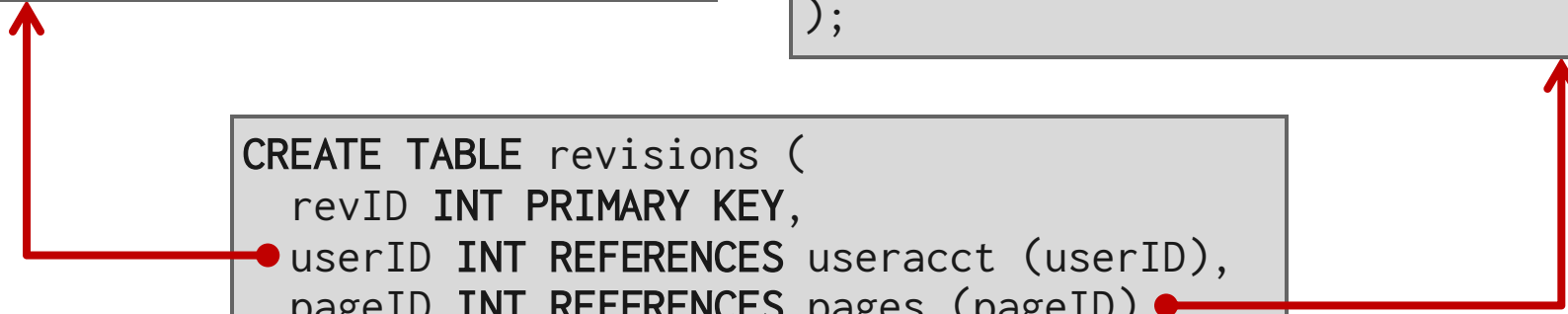
```
CREATE TABLE revisions (  
    revID INT PRIMARY KEY,  
    userID INT REFERENCES useracct (userID),  
    pageID INT REFERENCES pages (pageID),  
    content TEXT,  
    updated DATETIME  
);
```

Wikipedia Example

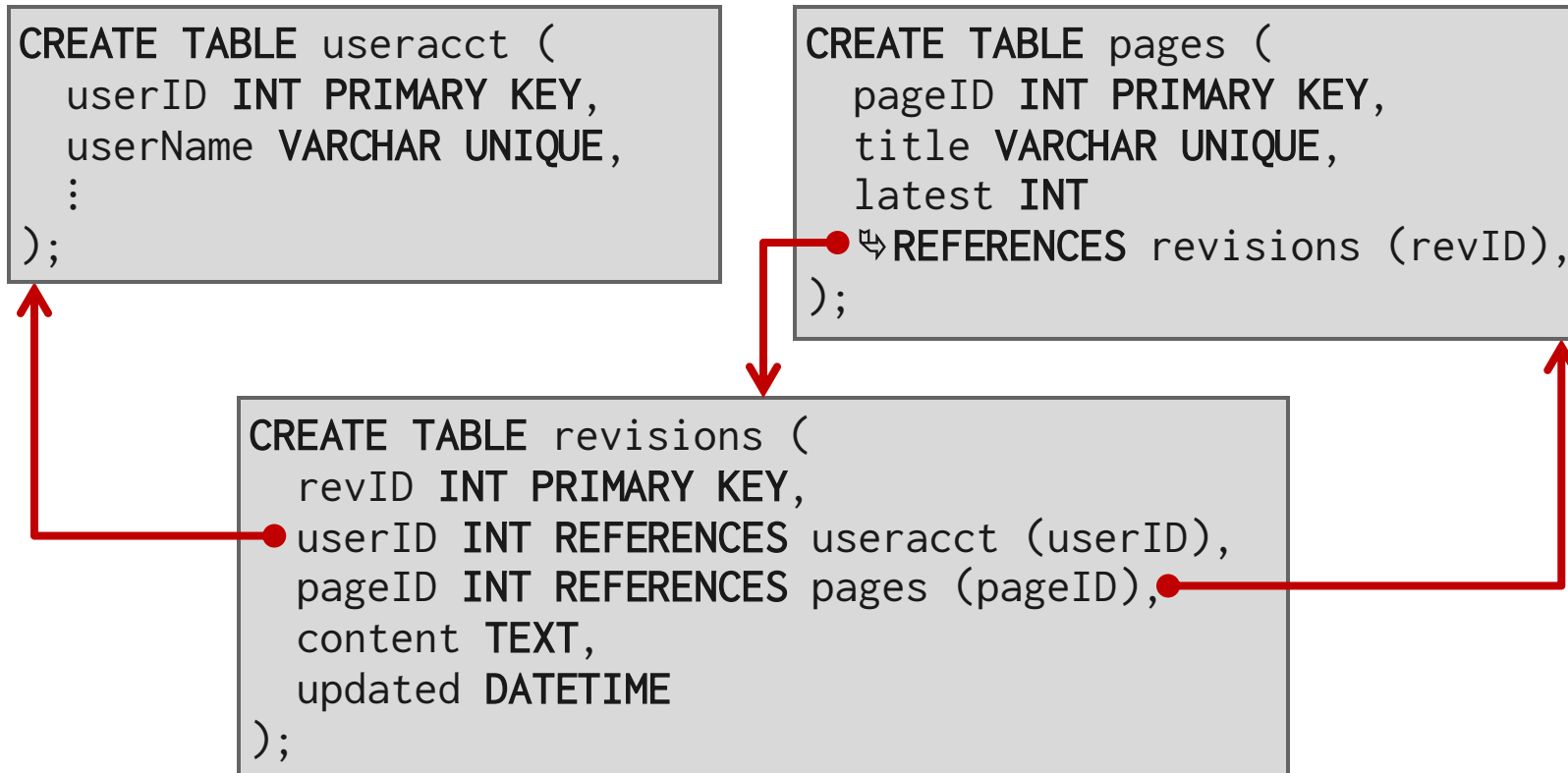
```
CREATE TABLE useracct (  
  userID INT PRIMARY KEY,  
  userName VARCHAR UNIQUE,  
  :  
);
```

```
CREATE TABLE pages (  
  pageID INT PRIMARY KEY,  
  title VARCHAR UNIQUE,  
  latest INT  
  ↳ REFERENCES revisions (revID),  
);
```

```
CREATE TABLE revisions (  
  revID INT PRIMARY KEY,  
  userID INT REFERENCES useracct (userID),  
  pageID INT REFERENCES pages (pageID),  
  content TEXT,  
  updated DATETIME  
);
```



Wikipedia Example



Observation

- The relational model does not specify that the DBMS must store all of a tuple's attributes together on a single page.
- This may not actually be the best layout for some workloads...

OLTP

- On-line Transaction Processing:
 - Simple queries that read/update a small amount of data related to a single entity in the database.
- This is usually the kind of application that people build first.

```
SELECT P.*, R.*  
  FROM pages AS P  
 INNER JOIN revisions AS R  
    ON P.latest = R.revID  
 WHERE P.pageID = ?
```

```
UPDATE useracct  
  SET lastLogin = NOW(),  
      hostname = ?  
 WHERE userID = ?
```

```
INSERT INTO revisions VALUES  
(?, ?, ?)
```

OLAP

- On-line Analytical Processing:
 - Complex queries that read large portions of the database spanning multiple entities.
- You execute these workloads on the data collected from your OLTP application(s).

```
SELECT COUNT(U.lastLogin),  
        EXTRACT(month FROM  
                U.lastLogin) AS month  
FROM useracct AS U  
WHERE U.hostname LIKE '%.gov'  
GROUP BY  
        EXTRACT(month FROM U.lastLogin)
```

Storage Models

- A DBMS's storage model specifies how it physically organizes tuples on disk and in memory.
 - Can have different performance characteristics based on the target workload (OLTP vs. OLAP).
 - Influences the design choices of the rest of the DBMS.
- **Choice #1: N-ary Storage Model (NSM)**
- **Choice #2: Decomposition Storage Model (DSM)**
- **Choice #3: Hybrid Storage Model (PAX)**

N-ary Storage Model (NSM)

N-ary Storage Model (NSM)

- The DBMS stores (almost) all attributes for a single tuple contiguously in a single page.
→ Also known as a “**row store**”.
- Ideal for OLTP workloads where queries are more likely to access individual entities and execute write-heavy workloads.
- NSM database page sizes are typically some constant multiple of 4 KB hardware pages.
→ Oracle (4 KB), Postgres (8 KB), MySQL (16 KB)

NSM: Physical Organization

- A disk-oriented NSM system stores a tuple's fixed-length and variable-length attributes contiguously in a single slotted page.
- The tuple's **record id** (page#, slot#) is how the DBMS uniquely identifies a physical tuple.

	Col A	Col B	Col C
Row #0	a0	b0	c0
Row #1	a1	b1	c1
Row #2	a2	b2	c2
Row #3	a3	b3	c3
Row #4	a4	b4	c4
Row #5	a5	b5	c5

NSM: Physical Organization

- A disk-oriented NSM system stores a tuple's fixed-length and variable-length attributes contiguously in a single slotted page.
- The tuple's **record id** (page#, slot#) is how the DBMS uniquely identifies a physical tuple.

	Col A	Col B	Col C
Row #0	a0	b0	c0
Row #1	a1	b1	c1
Row #2	a2	b2	c2
Row #3	a3	b3	c3
Row #4	a4	b4	c4
Row #5	a5	b5	c5

Database Page



NSM: Physical Organization

- A disk-oriented NSM system stores a tuple's fixed-length and variable-length attributes contiguously in a single slotted page.
- The tuple's **record id** (page#, slot#) is how the DBMS uniquely identifies a physical tuple.

	Col A	Col B	Col C
Row #0	a0	b0	c0
Row #1	a1	b1	c1
Row #2	a2	b2	c2
Row #3	a3	b3	c3
Row #4	a4	b4	c4
Row #5	a5	b5	c5

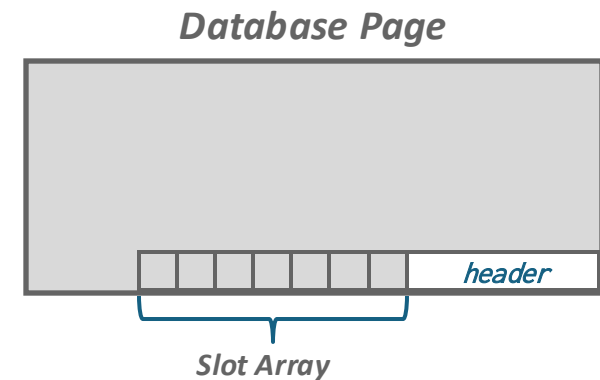
Database Page



NSM: Physical Organization

- A disk-oriented NSM system stores a tuple's fixed-length and variable-length attributes contiguously in a single slotted page.
- The tuple's **record id** (page#, slot#) is how the DBMS uniquely identifies a physical tuple.

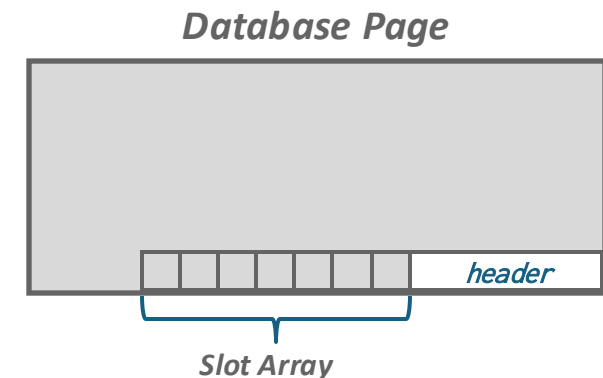
	Col A	Col B	Col C
Row #0	a0	b0	c0
Row #1	a1	b1	c1
Row #2	a2	b2	c2
Row #3	a3	b3	c3
Row #4	a4	b4	c4
Row #5	a5	b5	c5



NSM: Physical Organization

- A disk-oriented NSM system stores a tuple's fixed-length and variable-length attributes contiguously in a single slotted page.
- The tuple's **record id** (page#, slot#) is how the DBMS uniquely identifies a physical tuple.

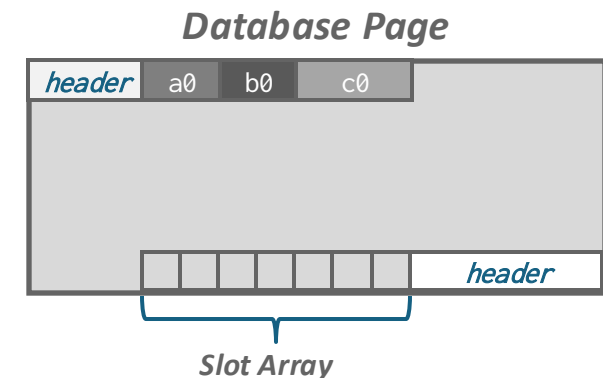
	Col A	Col B	Col C
Row #0	a0	b0	c0
Row #1	a1	b1	c1
Row #2	a2	b2	c2
Row #3	a3	b3	c3
Row #4	a4	b4	c4
Row #5	a5	b5	c5



NSM: Physical Organization

- A disk-oriented NSM system stores a tuple's fixed-length and variable-length attributes contiguously in a single slotted page.
- The tuple's **record id** (page#, slot#) is how the DBMS uniquely identifies a physical tuple.

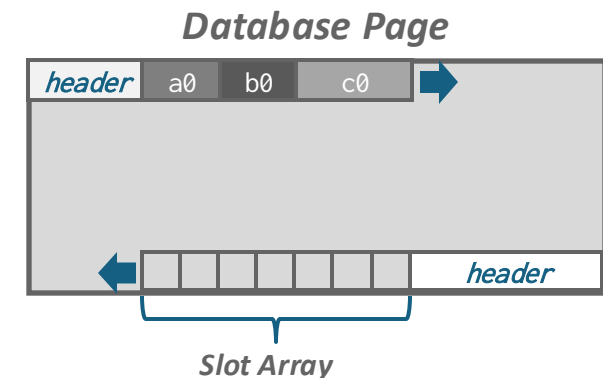
	Col A	Col B	Col C
Row #0	a0	b0	c0
Row #1	a1	b1	c1
Row #2	a2	b2	c2
Row #3	a3	b3	c3
Row #4	a4	b4	c4
Row #5	a5	b5	c5



NSM: Physical Organization

- A disk-oriented NSM system stores a tuple's fixed-length and variable-length attributes contiguously in a single slotted page.
- The tuple's **record id** (page#, slot#) is how the DBMS uniquely identifies a physical tuple.

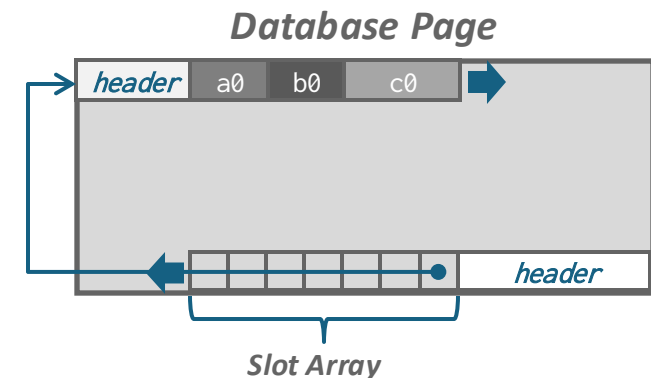
	Col A	Col B	Col C
Row #0	a0	b0	c0
Row #1	a1	b1	c1
Row #2	a2	b2	c2
Row #3	a3	b3	c3
Row #4	a4	b4	c4
Row #5	a5	b5	c5



NSM: Physical Organization

- A disk-oriented NSM system stores a tuple's fixed-length and variable-length attributes contiguously in a single slotted page.
- The tuple's **record id** (page#, slot#) is how the DBMS uniquely identifies a physical tuple.

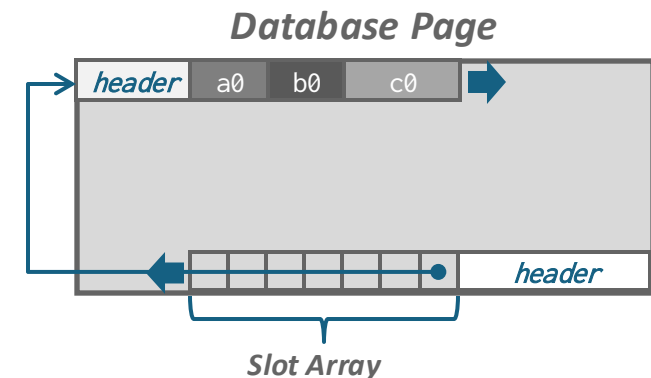
	Col A	Col B	Col C
Row #0	a0	b0	c0
Row #1	a1	b1	c1
Row #2	a2	b2	c2
Row #3	a3	b3	c3
Row #4	a4	b4	c4
Row #5	a5	b5	c5



NSM: Physical Organization

- A disk-oriented NSM system stores a tuple's fixed-length and variable-length attributes contiguously in a single slotted page.
- The tuple's **record id** (page#, slot#) is how the DBMS uniquely identifies a physical tuple.

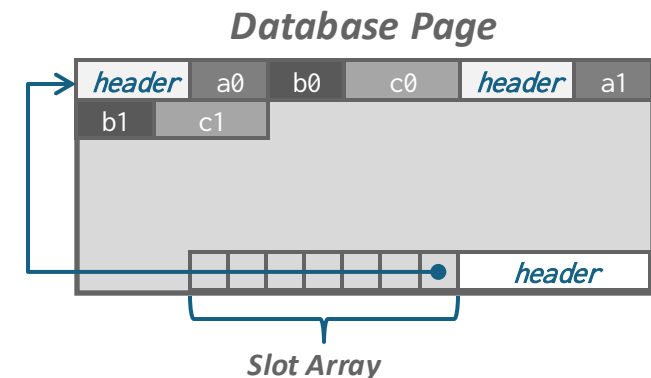
	Col A	Col B	Col C
Row #0	a0	b0	c0
Row #1	a1	b1	c1
Row #2	a2	b2	c2
Row #3	a3	b3	c3
Row #4	a4	b4	c4
Row #5	a5	b5	c5



NSM: Physical Organization

- A disk-oriented NSM system stores a tuple's fixed-length and variable-length attributes contiguously in a single slotted page.
- The tuple's **record id** (page#, slot#) is how the DBMS uniquely identifies a physical tuple.

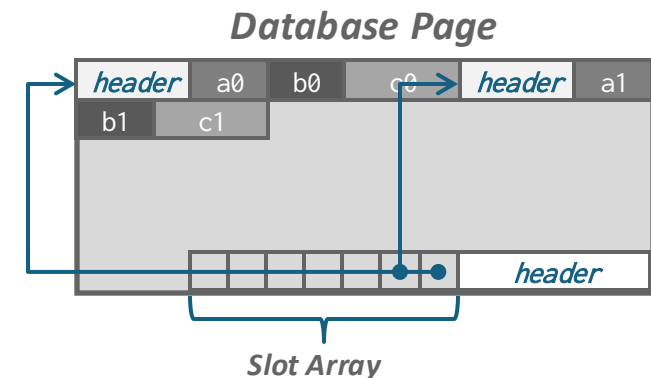
	Col A	Col B	Col C
Row #0	a0	b0	c0
Row #1	a1	b1	c1
Row #2	a2	b2	c2
Row #3	a3	b3	c3
Row #4	a4	b4	c4
Row #5	a5	b5	c5



NSM: Physical Organization

- A disk-oriented NSM system stores a tuple's fixed-length and variable-length attributes contiguously in a single slotted page.
- The tuple's **record id** (page#, slot#) is how the DBMS uniquely identifies a physical tuple.

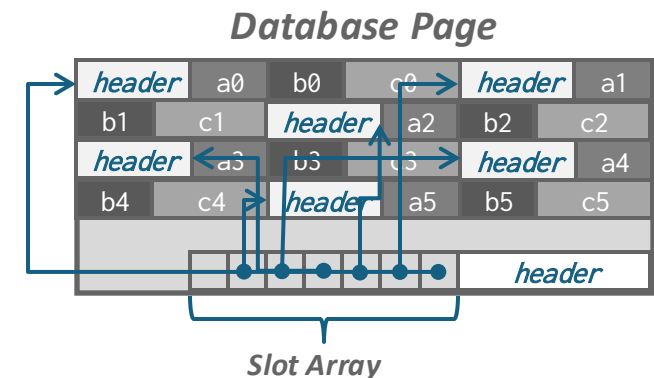
	Col A	Col B	Col C
Row #0	a0	b0	c0
Row #1	a1	b1	c1
Row #2	a2	b2	c2
Row #3	a3	b3	c3
Row #4	a4	b4	c4
Row #5	a5	b5	c5



NSM: Physical Organization

- A disk-oriented NSM system stores a tuple's fixed-length and variable-length attributes contiguously in a single slotted page.
- The tuple's **record id** (page#, slot#) is how the DBMS uniquely identifies a physical tuple.

	Col A	Col B	Col C
Row #0	a0	b0	c0
Row #1	a1	b1	c1
Row #2	a2	b2	c2
Row #3	a3	b3	c3
Row #4	a4	b4	c4
Row #5	a5	b5	c5



NSM: OLTP Example

```
SELECT * FROM useracct  
WHERE userName = ?  
AND userPass = ?
```



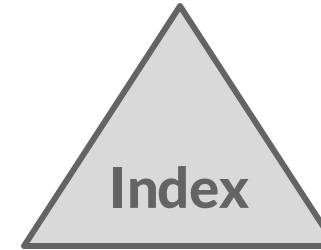
Disk

Database File



NSM: OLTP Example

```
SELECT * FROM useracct  
WHERE userName = ?  
AND userPass = ?
```



Disk

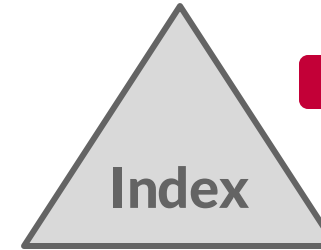
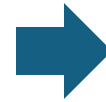
Database File



useracct	userName	userPass	...
1	John	123456	...
2	Jane	654321	...
3	Bob	987654	...
4	Alice	432109	...
5	Charlie	876543	...
6	Diana	210987	...
7	Frank	543210	...
8	Grace	098765	...
9	Henry	321098	...
10	Ivy	654321	...

NSM: OLTP Example

```
SELECT * FROM useracct  
WHERE userName = ?  
AND userPass = ?
```



Future Lecture

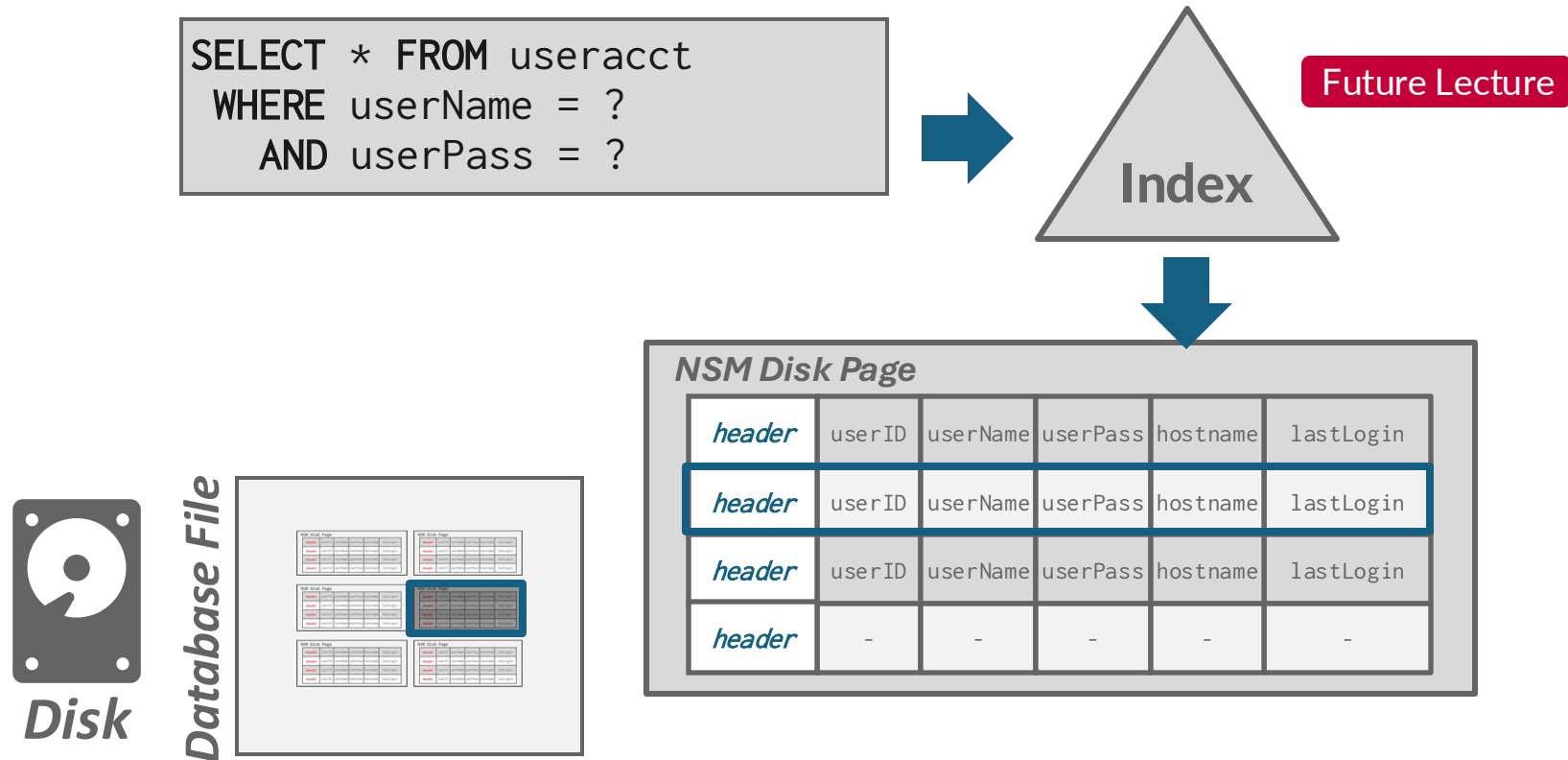


Disk

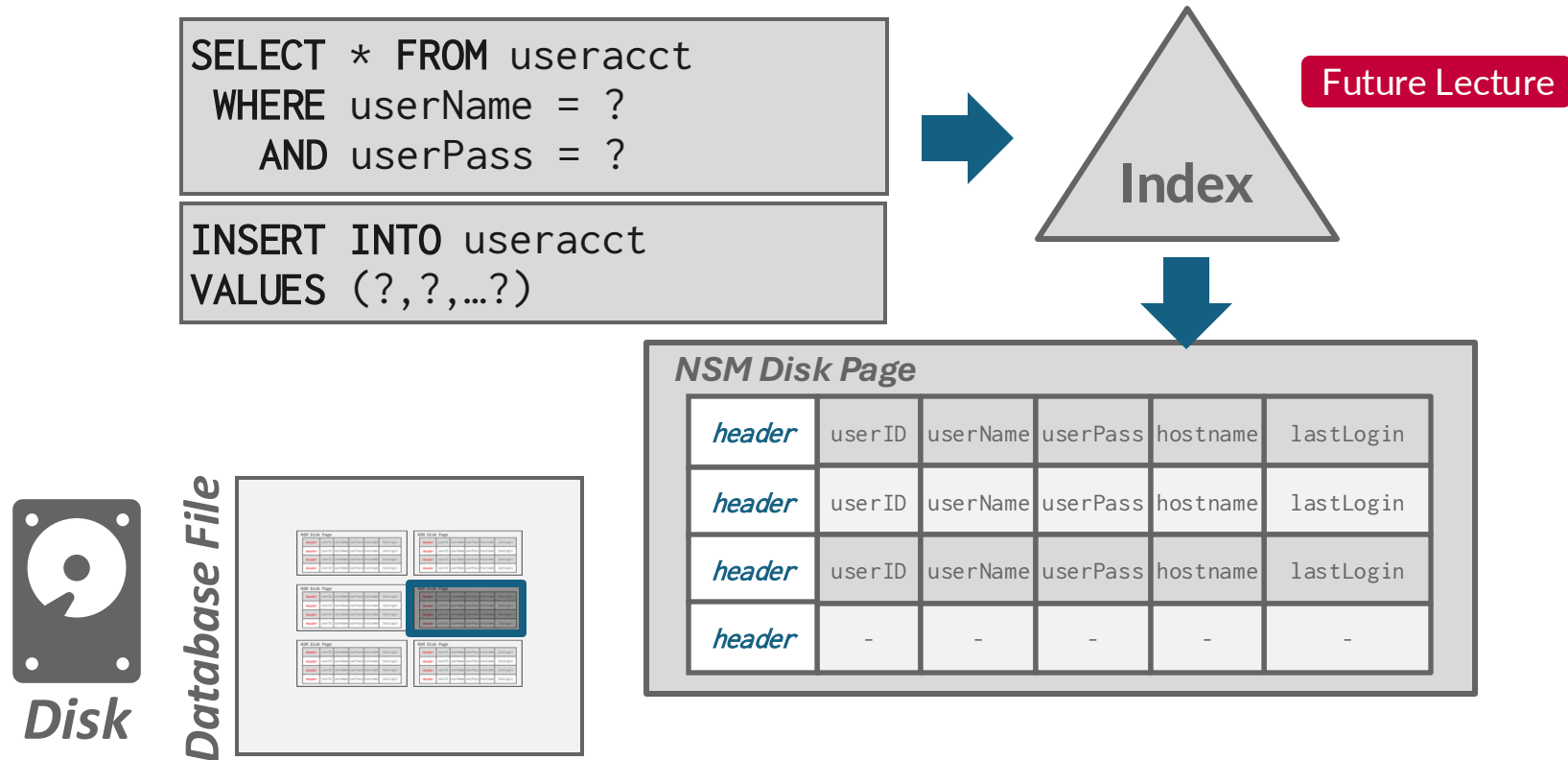
Database File



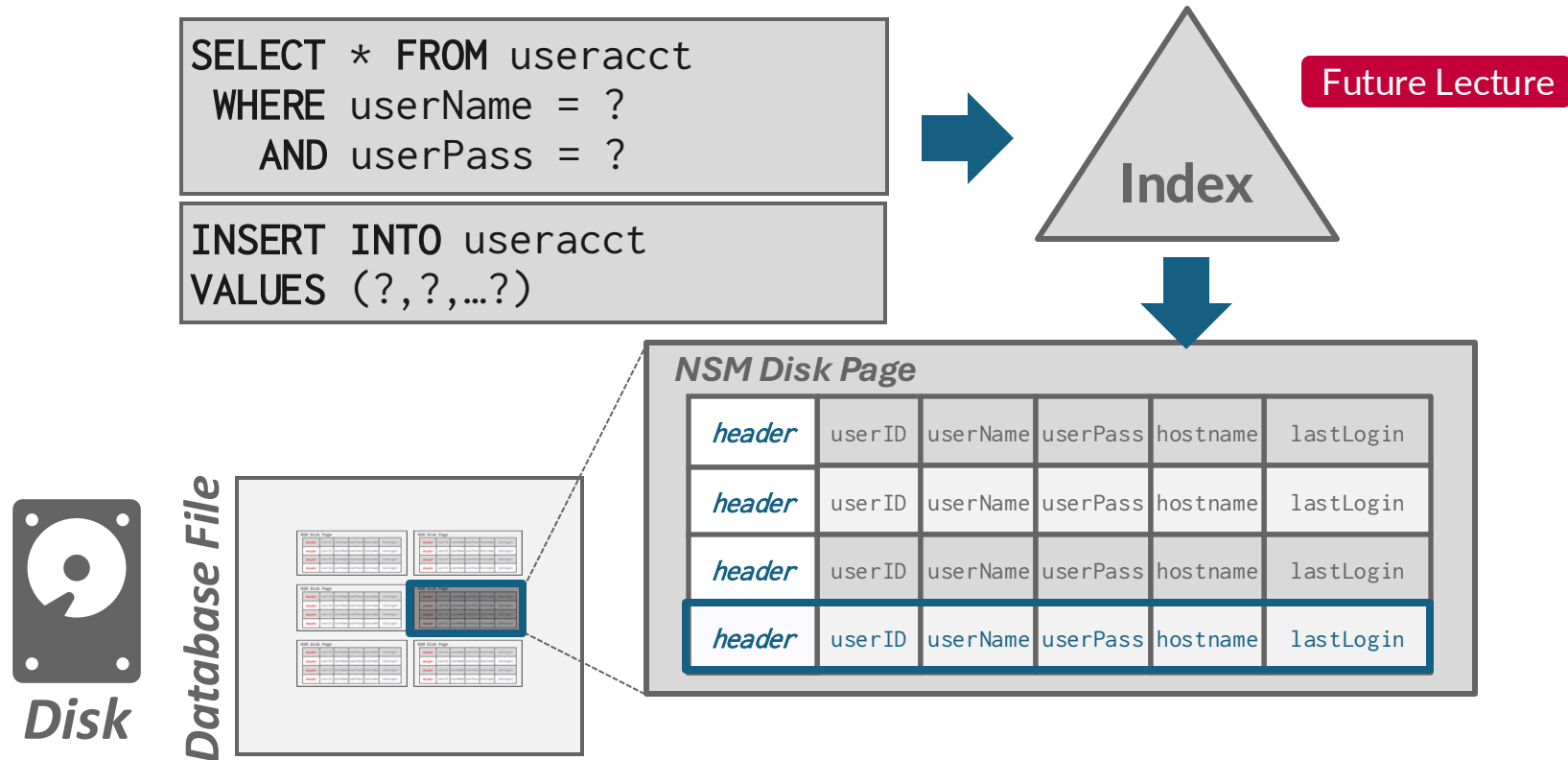
NSM: OLTP Example



NSM: OLTP Example



NSM: OLTP Example



NSM: OLAP Example

```
SELECT COUNT(U.lastLogin),  
       EXTRACT(month FROM U.lastLogin) AS month  
FROM useracct AS U  
WHERE U.hostname LIKE '%.gov'  
GROUP BY EXTRACT(month FROM U.lastLogin)
```



Database File



NSM: OLAP Example

```
SELECT COUNT(U.lastLogin),  
       EXTRACT(month FROM U.lastLogin) AS month  
FROM useracct AS U  
WHERE U.hostname LIKE '%.gov'  
GROUP BY EXTRACT(month FROM U.lastLogin)
```



Disk

Database File



NSM: OLAP Example

```
SELECT COUNT(U.lastLogin),
       EXTRACT(month FROM U.lastLogin) AS month
FROM useracct AS U
WHERE U.hostname LIKE '%.gov'
GROUP BY EXTRACT(month FROM U.lastLogin)
```



Disk

Database File



NSM Disk Page

<i>header</i>	userID	userName	userPass	hostname	lastLogin
<i>header</i>	userID	userName	userPass	hostname	lastLogin
<i>header</i>	userID	userName	userPass	hostname	lastLogin
<i>header</i>	userID	userName	userPass	hostname	lastLogin

NSM: OLAP Example

```
SELECT COUNT(U.lastLogin),
       EXTRACT(month FROM U.lastLogin) AS month
FROM useracct AS U
WHERE U.hostname LIKE '%.gov'
GROUP BY EXTRACT(month FROM U.lastLogin)
```



Disk

Database File



NSM Disk Page

<i>header</i>	userID	userName	userPass	hostname	lastLogin
<i>header</i>	userID	userName	userPass	hostname	lastLogin
<i>header</i>	userID	userName	userPass	hostname	lastLogin
<i>header</i>	userID	userName	userPass	hostname	lastLogin

NSM: OLAP Example

```
SELECT COUNT(U.lastLogin),
       EXTRACT(month FROM U.lastLogin) AS month
FROM useracct AS U
WHERE U.hostname LIKE '%.gov'
GROUP BY EXTRACT(month FROM U.lastLogin)
```



Disk

Database File



NSM Disk Page

<i>header</i>	userID	userName	userPass	hostname	lastLogin
<i>header</i>	userID	userName	userPass	hostname	lastLogin
<i>header</i>	userID	userName	userPass	hostname	lastLogin
<i>header</i>	userID	userName	userPass	hostname	lastLogin

NSM: OLAP Example

```
SELECT COUNT(U.lastLogin),  
       EXTRACT(month FROM U.lastLogin) AS month  
FROM useracct AS U  
WHERE U.hostname LIKE '%.gov'  
GROUP BY EXTRACT(month FROM U.lastLogin)
```



Disk

Database File



NSM Disk Page

<i>header</i>	userID	userName	userPass	hostname	lastLogin
<i>header</i>	userID	userName	userPass	hostname	lastLogin
<i>header</i>	userID	userName	userPass	hostname	lastLogin
<i>header</i>	userID	userName	userPass	hostname	lastLogin

NSM: OLAP Example

```
SELECT COUNT(U.lastLogin),  
       EXTRACT(month FROM U.lastLogin) AS month  
FROM useracct AS U  
WHERE U.hostname LIKE '%.gov'  
GROUP BY EXTRACT(month FROM U.lastLogin)
```



Disk

Database File



NSM Disk Page

<i>header</i>	userID	userName	userPass	hostname	lastLogin
<i>header</i>	userID	userName	userPass	hostname	lastLogin
<i>header</i>	userID	userName	userPass	hostname	lastLogin
<i>header</i>	userID	userName	userPass	hostname	lastLogin

NSM: OLAP Example

```
SELECT COUNT(U.lastLogin),  
       EXTRACT(month FROM U.lastLogin) AS month  
FROM useracct AS U  
WHERE U.hostname LIKE '%.gov'  
GROUP BY EXTRACT(month FROM U.lastLogin)
```



Disk

Database File



NSM Disk Page

<i>header</i>	userID	userName	userPass	hostname	lastLogin
<i>header</i>	userID	userName	userPass	hostname	lastLogin
<i>header</i>	userID	userName	userPass	hostname	lastLogin
<i>header</i>	userID	userName	userPass	hostname	lastLogin

NSM: OLAP Example

```
SELECT COUNT(U.lastLogin),  
       EXTRACT(month FROM U.lastLogin) AS month  
FROM useracct AS U  
WHERE U.hostname LIKE '%.gov'  
GROUP BY EXTRACT(month FROM U.lastLogin)
```



Disk

Database File



NSM Disk Page

<i>header</i>	userID	userName	userPass	hostname	lastLogin
<i>header</i>	userID	userName	userPass	hostname	lastLogin
<i>header</i>	userID	userName	userPass	hostname	lastLogin
<i>header</i>	userID	userName	userPass	hostname	lastLogin

Useless Data

NSM: Summary

- **Advantages**

- Fast inserts, updates, and deletes.
- Good for queries that need the entire tuple (OLTP).
- Can use index-oriented physical storage for clustering.

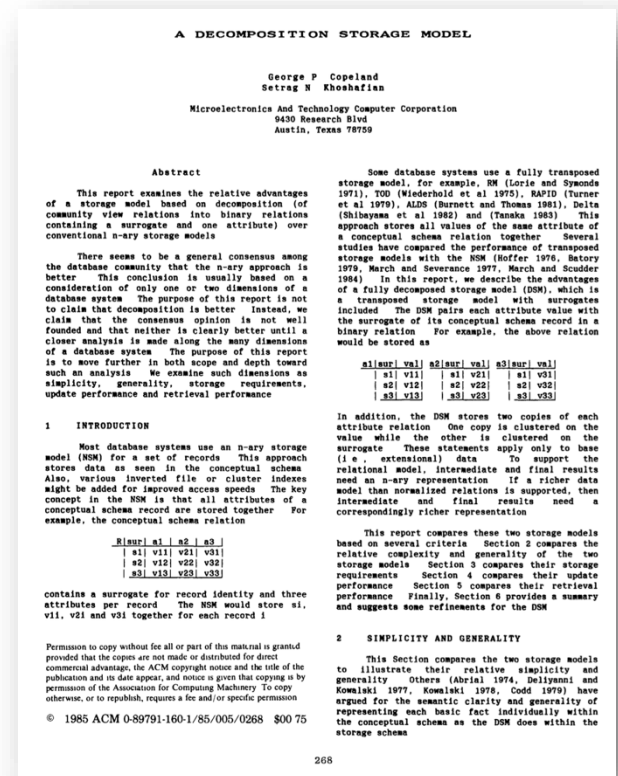
- **Disadvantages**

- Not good for scanning large portions of the table and/or a subset of the attributes.
- Terrible memory locality for OLAP access patterns.
- Not ideal for compression because of multiple value domains within a single page.

Decomposition Storage Model (DSM)

Decomposition Storage Model

- The DBMS stores a single attribute for all tuples contiguously in a block of data.
→ Also known as a “**column store**”.
- Ideal for OLAP workloads where read-only queries perform large scans over a subset of the table's attributes.
- DBMS is responsible for combining/splitting a tuple's attributes when reading/writing.



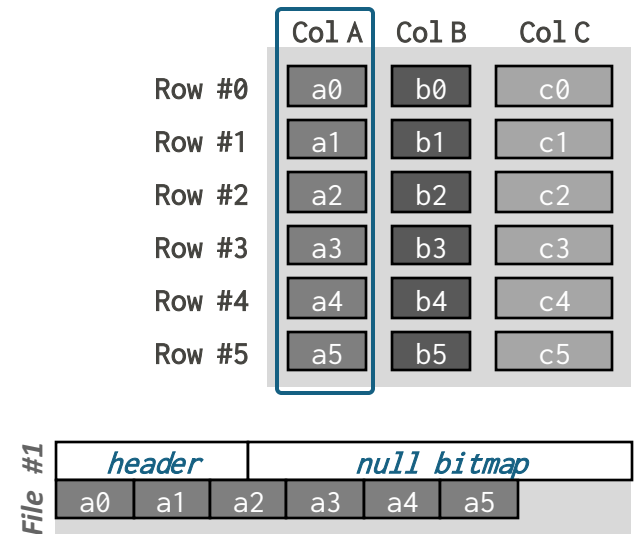
DSM: Physical Organization

- Store attributes and metadata (e.g., nulls) in separate arrays of **fixed-length** values.
 - Most systems identify unique physical tuples using offsets into these arrays.
 - Need to handle variable-length values...
- Maintain a separate file per attribute with a dedicated header area for metadata about the entire column.

	Col A	Col B	Col C
Row #0	a0	b0	c0
Row #1	a1	b1	c1
Row #2	a2	b2	c2
Row #3	a3	b3	c3
Row #4	a4	b4	c4
Row #5	a5	b5	c5

DSM: Physical Organization

- Store attributes and metadata (e.g., nulls) in separate arrays of **fixed-length** values.
 - Most systems identify unique physical tuples using offsets into these arrays.
 - Need to handle variable-length values...
- Maintain a separate file per attribute with a dedicated header area for metadata about the entire column.



DSM: Physical Organization

- Store attributes and metadata (e.g., nulls) in separate arrays of **fixed-length** values.
→ Most systems identify unique physical tuples using offsets into these arrays.
→ Need to handle variable-length values...
- Maintain a separate file per attribute with a dedicated header area for metadata about the entire column.

	Col A	Col B	Col C
Row #0	a0	b0	c0
Row #1	a1	b1	c1
Row #2	a2	b2	c2
Row #3	a3	b3	c3
Row #4	a4	b4	c4
Row #5	a5	b5	c5

File #1	header		null bitmap			
	a0	a1	a2	a3	a4	a5

File #2	header		null bitmap			
	b0	b1	b2	b3	b4	b5

DSM: Physical Organization

- Store attributes and metadata (e.g., nulls) in separate arrays of **fixed-length** values.
→ Most systems identify unique physical tuples using offsets into these arrays.
→ Need to handle variable-length values...
- Maintain a separate file per attribute with a dedicated header area for metadata about the entire column.

	Col A	Col B	Col C
Row #0	a0	b0	c0
Row #1	a1	b1	c1
Row #2	a2	b2	c2
Row #3	a3	b3	c3
Row #4	a4	b4	c4
Row #5	a5	b5	c5

File #1

header			null bitmap			
a0	a1	a2	a3	a4	a5	

File #2

header			null bitmap			
b0	b1	b2	b3	b4	b5	

File #3

header		null bitmap			
c0	c1	c2	c3	c4	
c5					

DSM: Database Example

- The DBMS stores the values of a single attribute across multiple tuples contiguously in a page.
→ Also known as a “column store”.

<i>header</i>	userID	userName	userPass	hostname	lastLogin
<i>header</i>	userID	userName	userPass	hostname	lastLogin
<i>header</i>	userID	userName	userPass	hostname	lastLogin
<i>header</i>	userID	userName	userPass	hostname	lastLogin

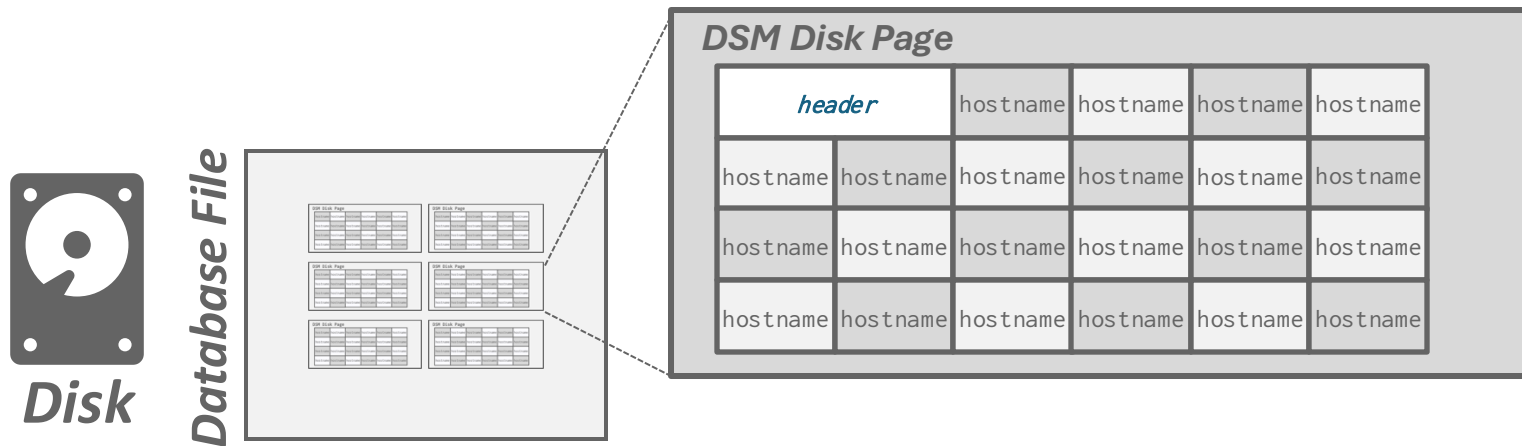
DSM: Database Example

- The DBMS stores the values of a single attribute across multiple tuples contiguously in a page.
→ Also known as a “column store”.

<i>header</i>	userID	userName	userPass	hostname	lastLogin
<i>header</i>	userID	userName	userPass	hostname	lastLogin
<i>header</i>	userID	userName	userPass	hostname	lastLogin
<i>header</i>	userID	userName	userPass	hostname	lastLogin

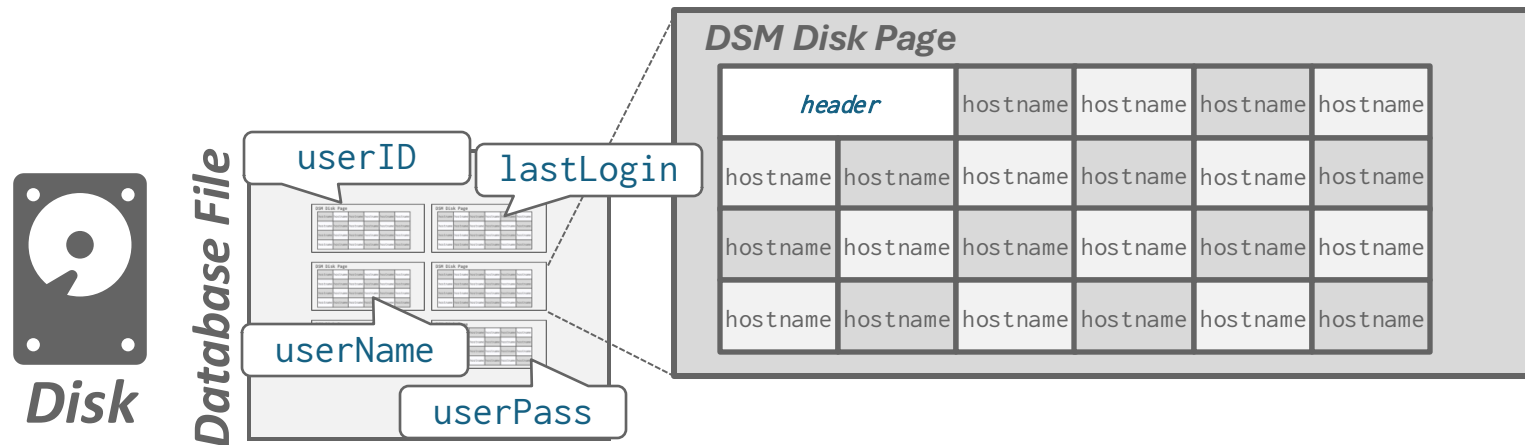
DSM: Database Example

- The DBMS stores the values of a single attribute across multiple tuples contiguously in a page.
→ Also known as a “column store”.



DSM: Database Example

- The DBMS stores the values of a single attribute across multiple tuples contiguously in a page.
→ Also known as a “column store”.



DSM: Tuple Identification

- **Choice #1: Fixed-length Offsets**

→ Each value is the same length for an attribute.

- **Choice #2: Embedded Tuple Ids**

→ Each value is stored with its tuple id in a column.

Offsets

	A	B	C	D
0				
1				
2				
3				

Embedded Ids

	A		B		C		D
0		0		0		0	
1		1		1		1	
2		2		2		2	
3		3		3		3	

DSM: Variable-Length Data

- **Padding** variable-length fields to ensure they are fixed-length is wasteful, especially for large attributes.
 - A better approach is to use *dictionary compression* to convert repetitive variable-length data into fixed-length values (typically 32-bit integers).
- More on this in a few slides.

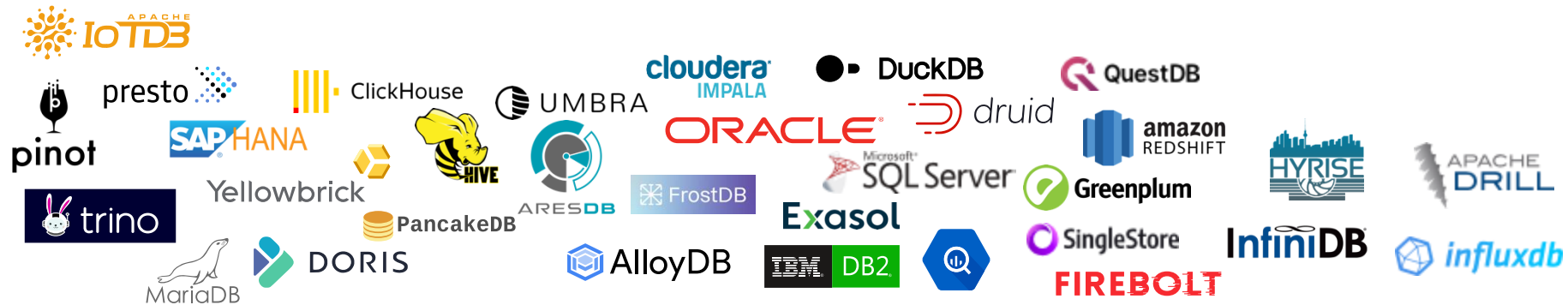
DSM: System History

- 1970s: Cantor DBMS
- 1980s: DSM Proposal
- 1990s: SybaseIQ (in-memory only)
- 2000s: Vertica, Vectorwise, MonetDB
- 2010s: Everyone + Parquet / ORC



DSM: System History

- 1970s: Cantor DBMS
- 1980s: DSM Proposal
- 1990s: SybaseIQ (in-memory only)
- 2000s: Vertica, Vectorwise, MonetDB
- 2010s: Everyone + Parquet / ORC



Decomposition Storage Model

- **Advantages**

- Reduces the amount wasted I/O per query because the DBMS only reads the data that it needs.
- Faster query processing because of increased locality and cached data reuse.
- Better data compression (more on this in a few slides).

- **Disadvantages**

- Slow for point queries, inserts, updates, and deletes because of tuple splitting/stitching/reorganization.

Hybrid Storage Model (PAX)

Observation

- OLAP queries rarely access a single column in a table by itself.
→ At some point during query execution, the DBMS must get other columns and stitch the original tuple back together.
- But we still need to store data in a columnar format to get the storage + execution benefits.
- We need a columnar scheme that still stores attributes separately but keeps the data for each tuple physically close to each other...

PAX Storage Model

- Partition Attributes Across (PAX) is a hybrid storage model that vertically partitions attributes within a database page.
→ This is what Parquet and Orc use.
- The goal is to get the benefit of faster processing on columnar storage while retaining the spatial locality benefits of row storage.



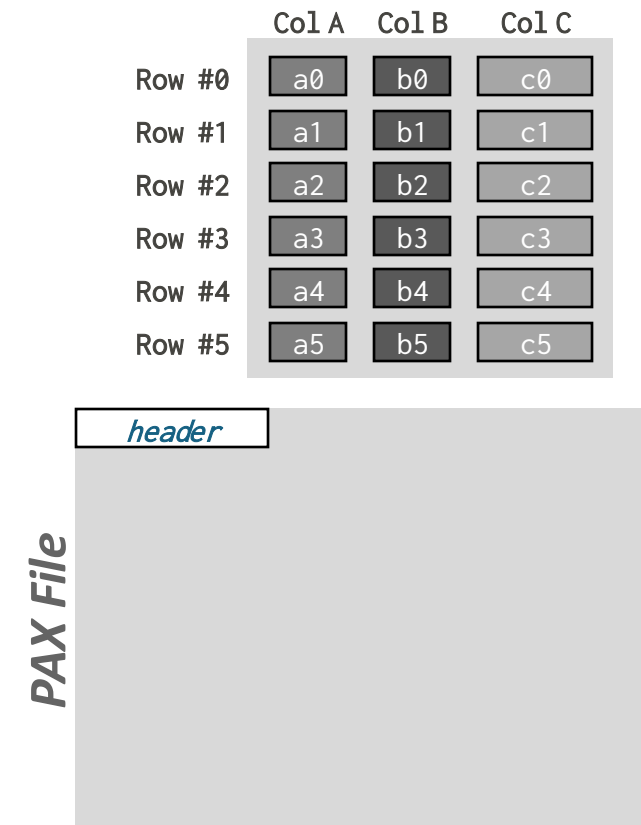
PAX: Physical Organization

- Horizontally partition rows into groups. Then vertically partition their attributes into columns.
- Global header contains directory with the offsets to the file's row groups.
 - This is stored in the footer if the file is immutable (Parquet, Orc).
- Each row group contains its own metadata header about its contents.

	Col A	Col B	Col C
Row #0	a0	b0	c0
Row #1	a1	b1	c1
Row #2	a2	b2	c2
Row #3	a3	b3	c3
Row #4	a4	b4	c4
Row #5	a5	b5	c5

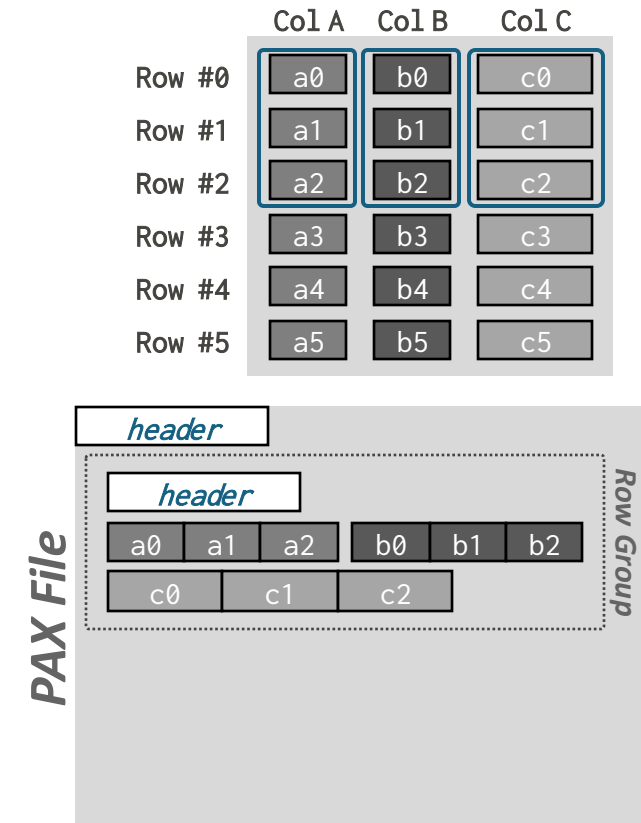
PAX: Physical Organization

- Horizontally partition rows into groups. Then vertically partition their attributes into columns.
- Global header contains directory with the offsets to the file's row groups.
 - This is stored in the footer if the file is immutable (Parquet, Orc).
- Each row group contains its own metadata header about its contents.



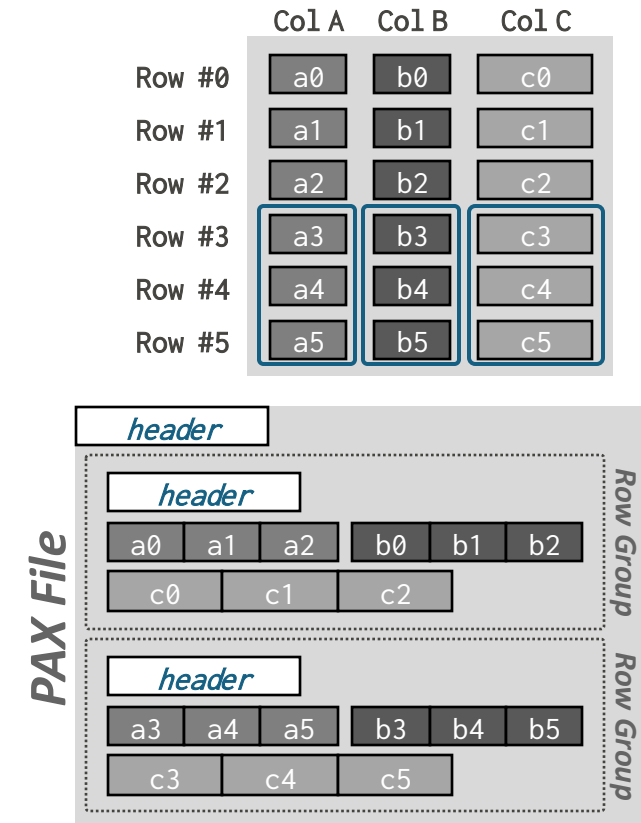
PAX: Physical Organization

- Horizontally partition rows into groups. Then vertically partition their attributes into columns.
- Global header contains directory with the offsets to the file's row groups.
→ This is stored in the footer if the file is immutable (Parquet, Orc).
- Each row group contains its own metadata header about its contents.



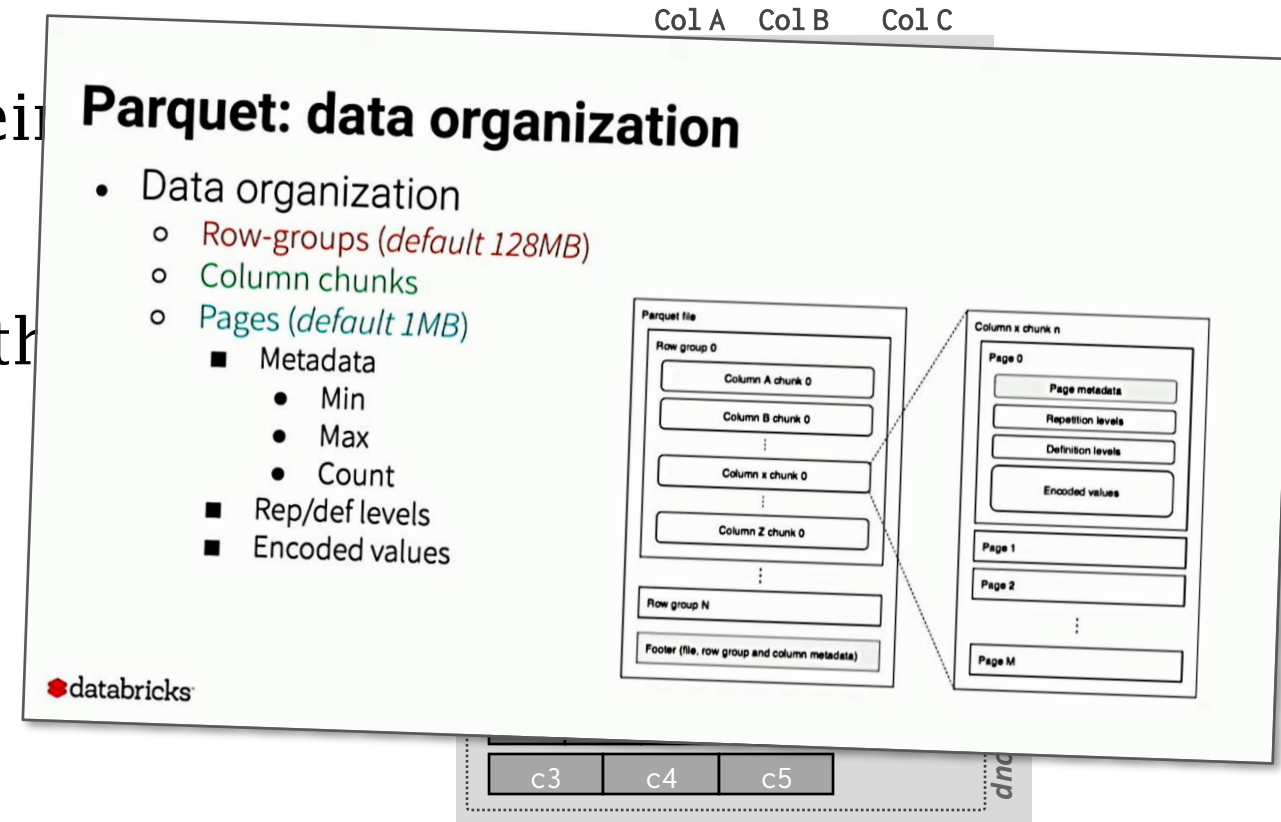
PAX: Physical Organization

- Horizontally partition rows into groups. Then vertically partition their attributes into columns.
- Global header contains directory with the offsets to the file's row groups.
→ This is stored in the footer if the file is immutable (Parquet, Orc).
- Each row group contains its own metadata header about its contents.



PAX: Physical Organization

- Horizontally partition rows into groups. Then vertically partition their attributes into columns.
- Global header contains directory with the offsets to the file's row groups.
→ This is stored in the footer if the file is immutable (Parquet, Orc).
- Each row group contains its own metadata header about its contents.



Compression

Observation

- I/O is the main bottleneck if the DBMS fetches data from disk during query execution.
- The DBMS can **compress** pages to increase the utility of the data moved per I/O operation.
- Key trade-off is speed vs. compression ratio
 - Compressing the database reduces DRAM requirements.
 - It may decrease CPU costs during query execution.

Database Compression

- **Goal #1:** Must produce fixed-length values.
→ Only exception is var-length data stored in separate pool.
- **Goal #2:** Postpone decompression for as long as possible during query execution.
→ Also known as late materialization.
- **Goal #3:** Must be a lossless scheme.

Lossless v.s. Lossy Compression

- When a DBMS uses compression, it is always lossless because people don't like losing data.
- Any kind of lossy compression must be performed at the application level.

Compression Granularity

- **Choice #1: Block-level**

→ Compress a block of tuples for the same table.

- **Choice #2: Tuple-level**

→ Compress the contents of the entire tuple (NSM-only, row-store).

- **Choice #3: Attribute-level**

→ Compress a single attribute within one tuple (overflow).

→ Can target multiple attributes for the same tuple.

- **Choice #4: Column-level**

→ Compress multiple values for one or more attributes stored for multiple tuples (DSM-only, column store).

Compression Granularity

- **Choice #1: Block-level**

→ Compress a block of tuples for the same table.

- **Choice #2: Tuple-level**

→ Compress the contents of the entire tuple (NSM-only, row-store).

- **Choice #3: Attribute-level**

→ Compress a single attribute within one tuple (overflow).

→ Can target multiple attributes for the same tuple.

- **Choice #4: Column-level**

→ Compress multiple values for one or more attributes stored for multiple tuples (DSM-only, column store).

Naïve Compression

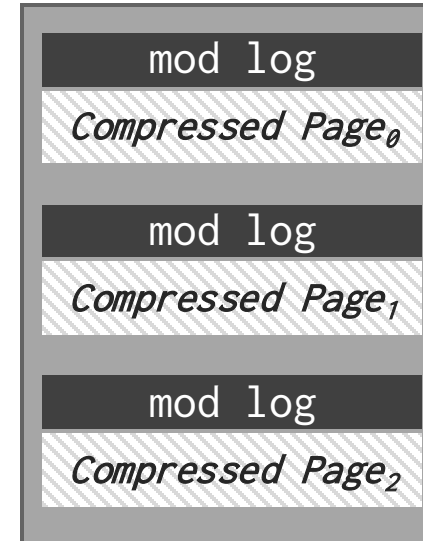
- Compress data using a general-purpose algorithm. The scope of compression is only based on the data provided as input.
 - [LZO](#) (1996), [LZ4](#) (2011), [Snappy](#) (2011),
[Oracle OZIP](#) (2014), [Zstd](#) (2015)
- Considerations
 - Computational overhead
 - Compress vs. decompress speed.

MySQL InnoDB Compression

 **Buffer Pool**



 **Database File**



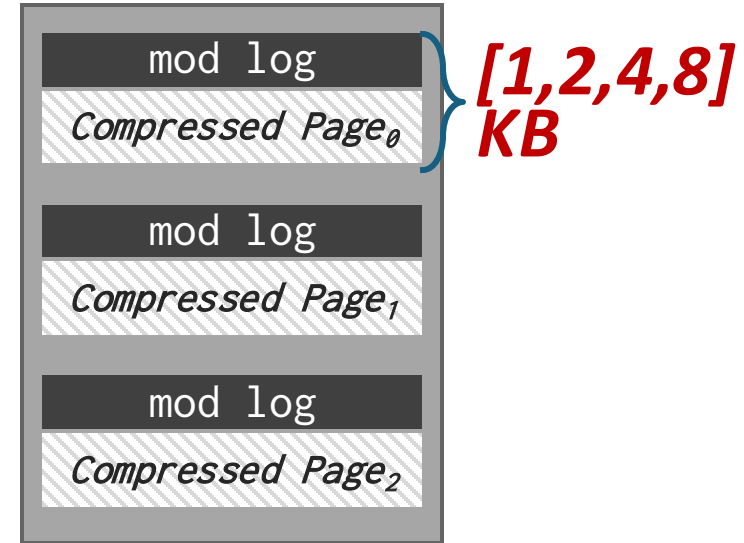
Source: [MySQL 5.7 Documentation](#)

MySQL InnoDB Compression

 *Buffer Pool*

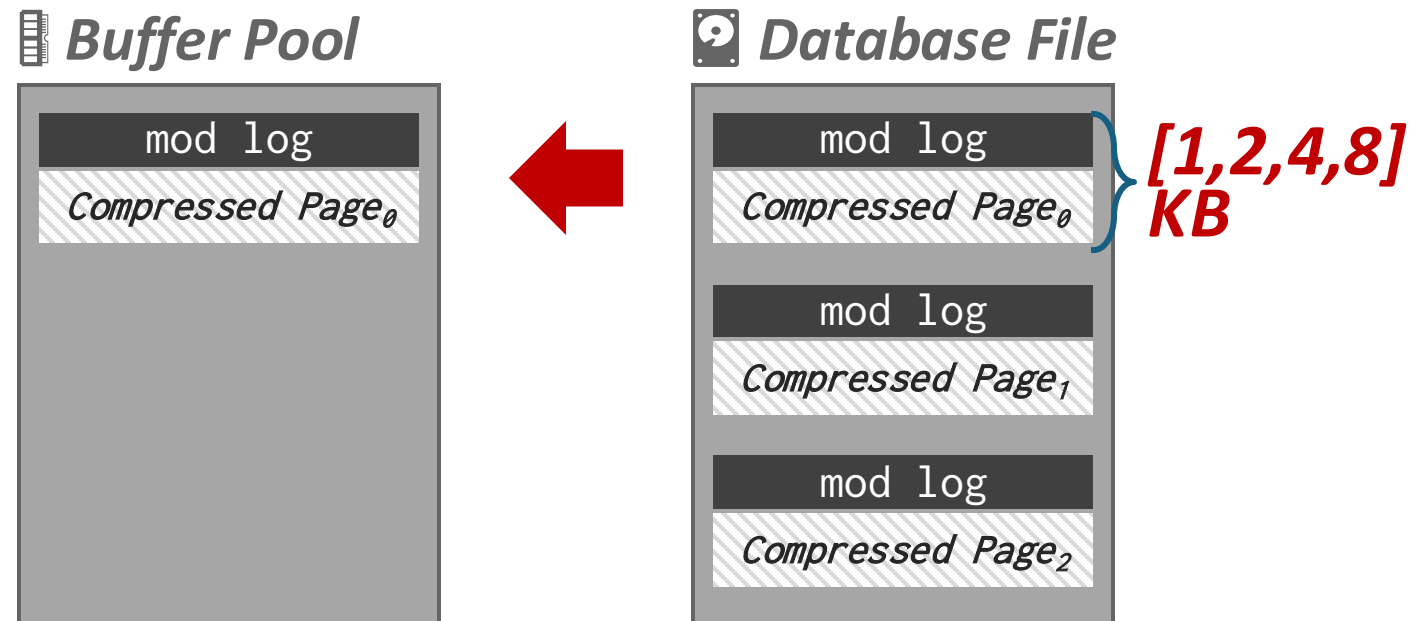


 *Database File*



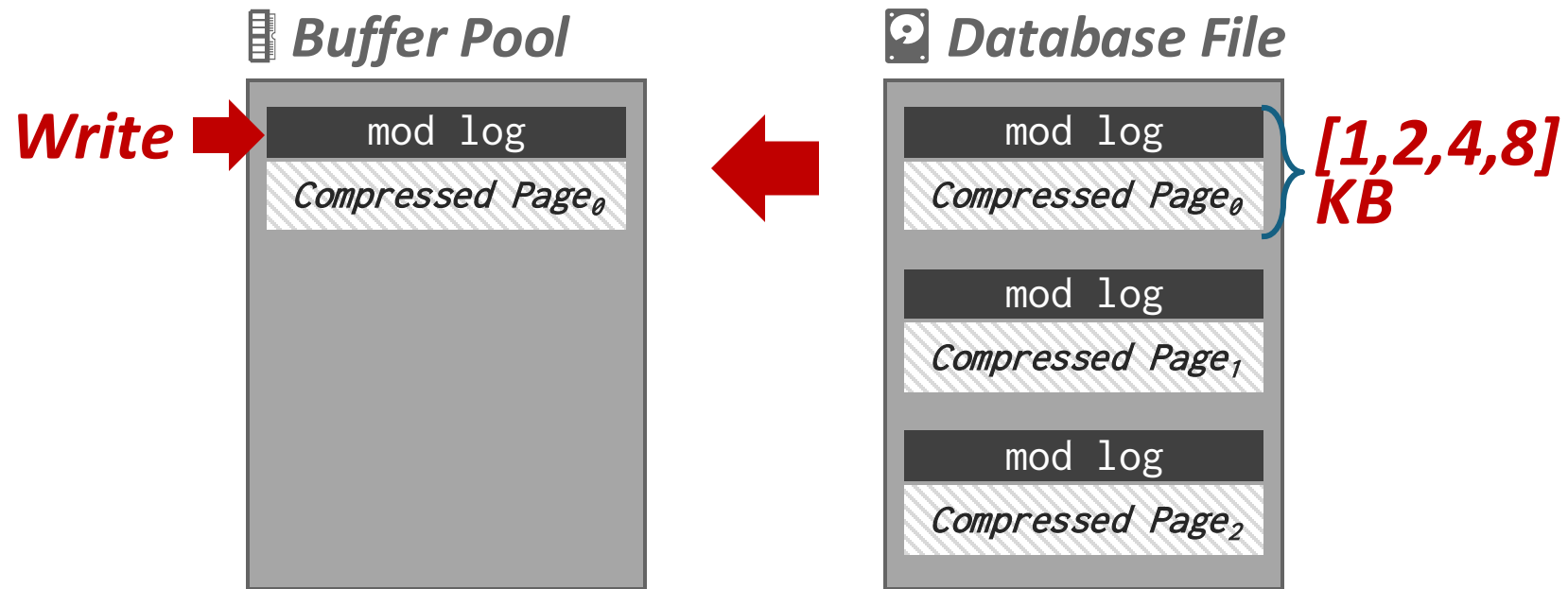
Source: [MySQL 5.7 Documentation](#)

MySQL InnoDB Compression



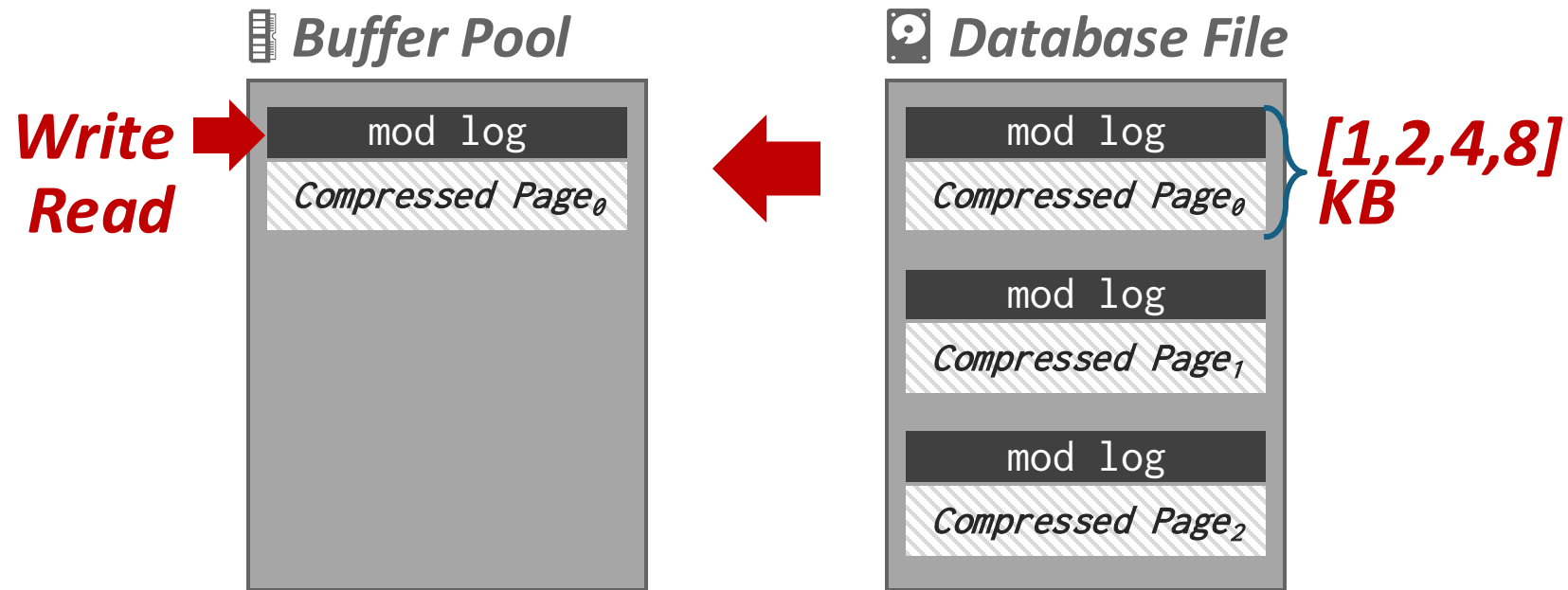
Source: [MySQL 5.7 Documentation](#)

MySQL InnoDB Compression



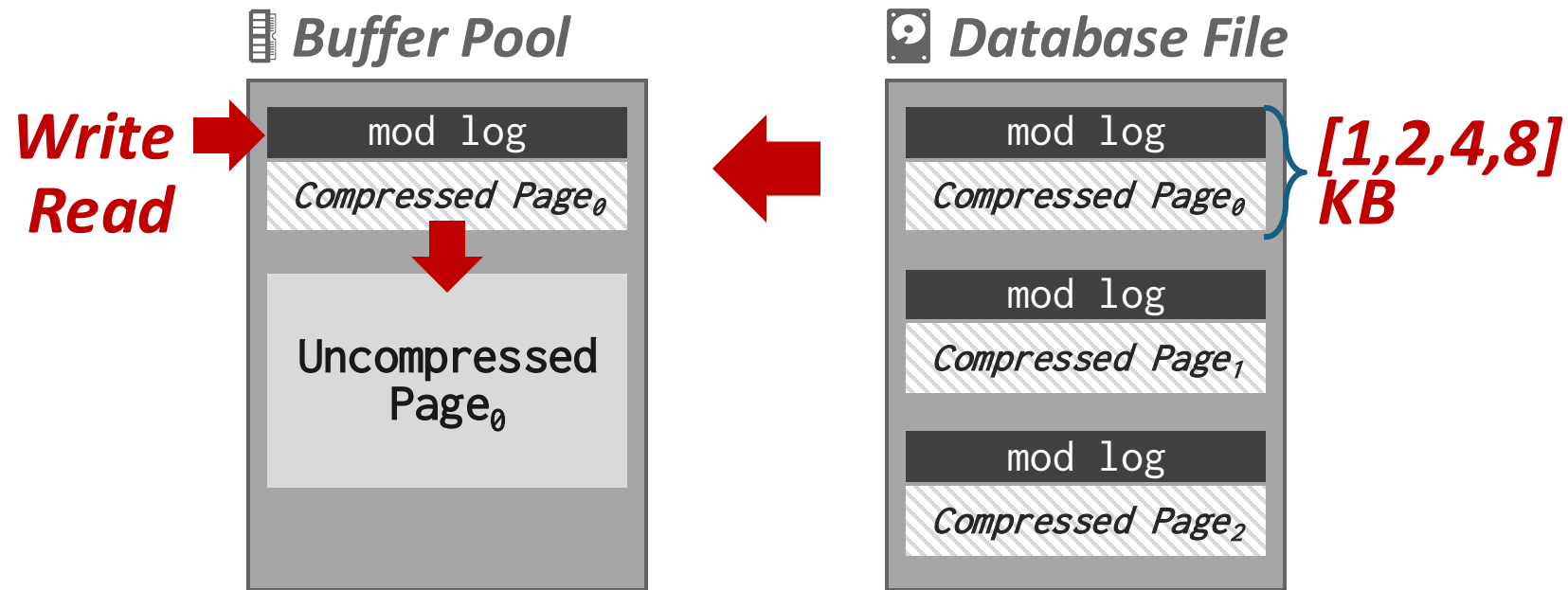
Source: [MySQL 5.7 Documentation](#)

MySQL InnoDB Compression



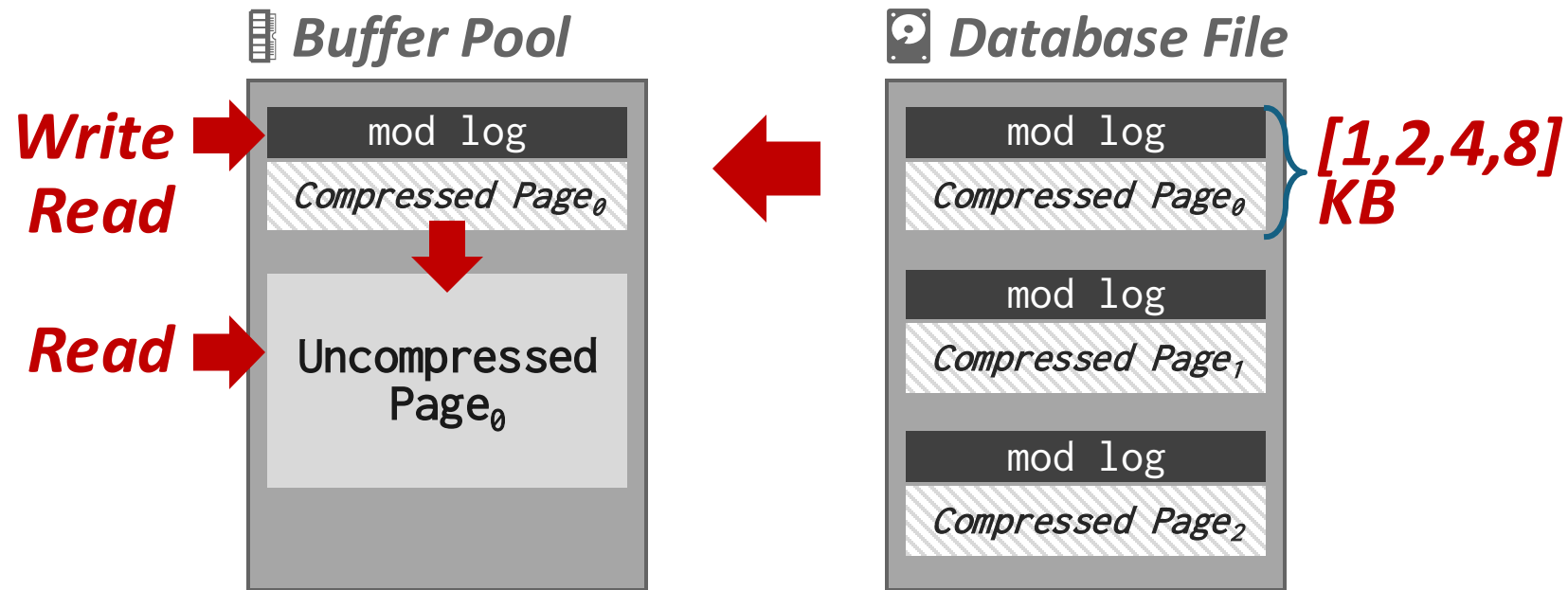
Source: [MySQL 5.7 Documentation](#)

MySQL InnoDB Compression



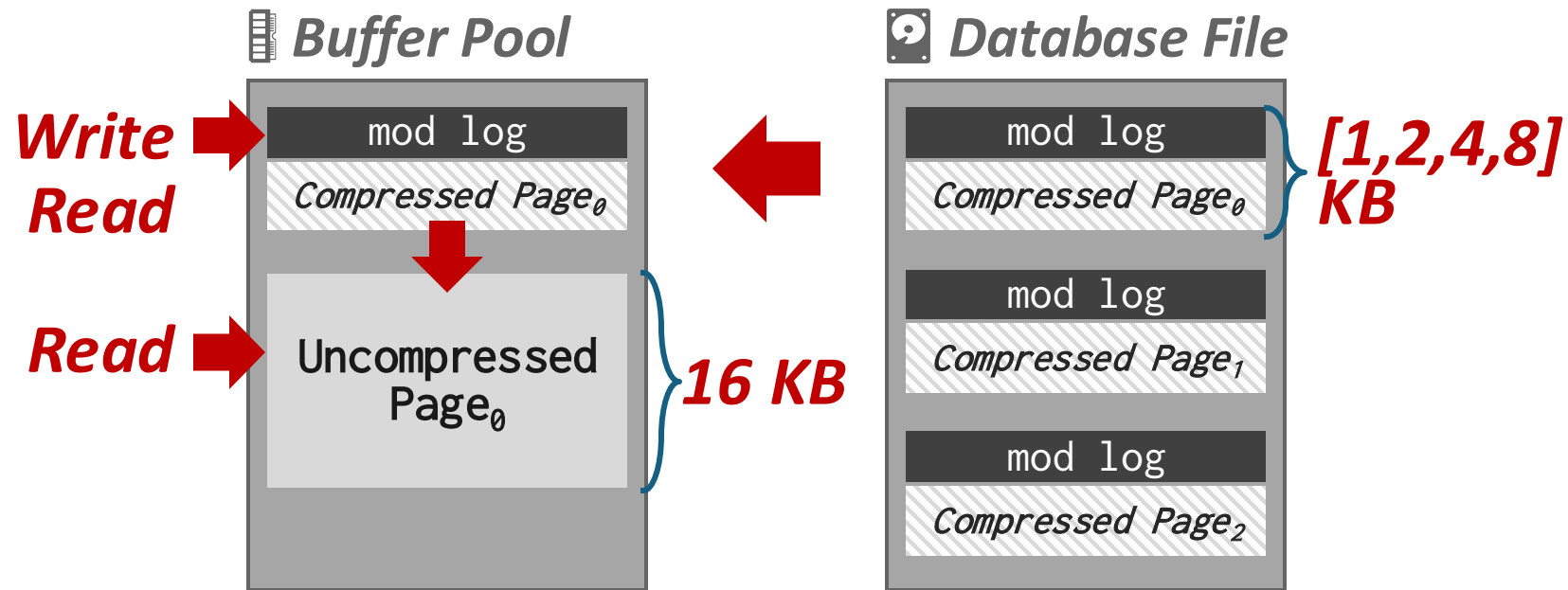
Source: [MySQL 5.7 Documentation](#)

MySQL InnoDB Compression



Source: [MySQL 5.7 Documentation](#)

MySQL InnoDB Compression



Source: [MySQL 5.7 Documentation](#)

Naïve Compression

- The DBMS must decompress data first before it can be read and (potentially) modified.
→ This limits the “scope” of the compression scheme.
- These schemes also do not consider the high-level meaning or semantics of the data.

Observation

- Ideally, we want the DBMS to operate on compressed data without decompressing it first.

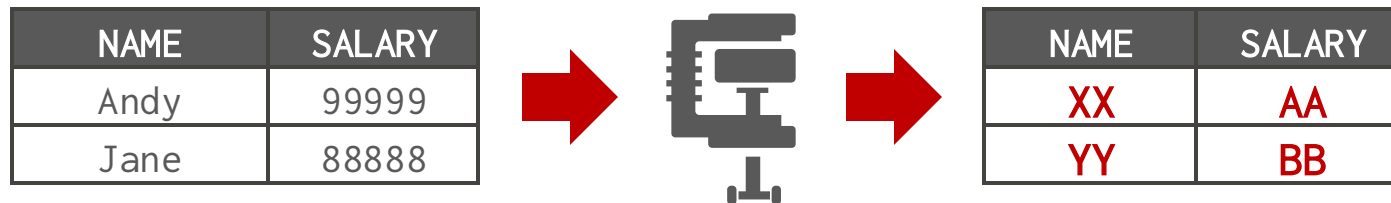
Observation

- Ideally, we want the DBMS to operate on compressed data without decompressing it first.

NAME	SALARY
Andy	99999
Jane	88888

Observation

- Ideally, we want the DBMS to operate on compressed data without decompressing it first.

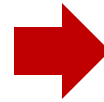
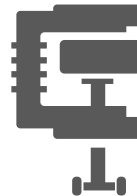
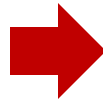


Observation

- Ideally, we want the DBMS to operate on compressed data without decompressing it first.

```
SELECT * FROM users  
WHERE name = 'Andy'
```

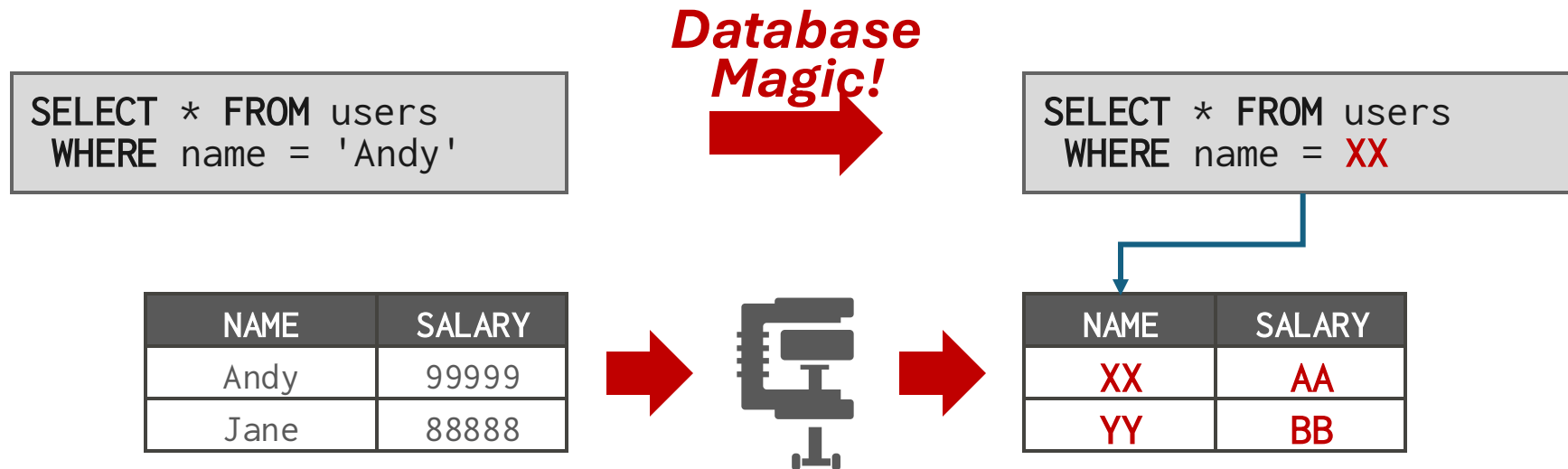
NAME	SALARY
Andy	99999
Jane	88888



NAME	SALARY
XX	AA
YY	BB

Observation

- Ideally, we want the DBMS to operate on compressed data without decompressing it first.



Compression Granularity

- **Choice #1: Block-level**

→ Compress a block of tuples for the same table.

- **Choice #2: Tuple-level**

→ Compress the contents of the entire tuple (NSM-only, row-store).

- **Choice #3: Attribute-level**

→ Compress a single attribute within one tuple (overflow).

→ Can target multiple attributes for the same tuple.

- **Choice #4: Column-level**

→ Compress multiple values for one or more attributes stored for multiple tuples (DSM-only, column store).

Compression Granularity

- **Choice #1: Block-level**

- Compress a block of tuples for the same table.

- **Choice #2: Tuple-level**

- Compress the contents of the entire tuple (NSM-only, row-store).

- **Choice #3: Attribute-level**

- Compress a single attribute within one tuple (overflow).

- Can target multiple attributes for the same tuple.

- **Choice #4: Column-level**

- Compress multiple values for one or more attributes stored for multiple tuples (DSM-only, column store).

Columnar Compression

- Run-length Encoding
- Bit-Packing Encoding
- Bitmap Encoding
- Delta Encoding
- Dictionary Encoding

Run-length Encoding

- Compress runs of the same value in a single column into triplets:
 - The value of the attribute.
 - The start position in the column segment.
 - The # of elements in the run.
- Requires the columns to be sorted intelligently to maximize compression opportunities.

Run-Length Encoding

Original Data

id	isDead
1	Y
2	Y
3	Y
4	N
6	Y
7	N
8	Y
9	Y

Run-Length Encoding

Original Data

id	isDead
1	Y
2	Y
3	Y
4	N
6	Y
7	N
8	Y
9	Y

Run-Length Encoding

Original Data

id	isDead
1	Y
2	Y
3	Y
4	N
6	Y
7	N
8	Y
9	Y



Compressed Data

id	isDead
1	(Y, 0, 3)
2	(N, 3, 1)
3	(Y, 4, 1)
4	(N, 5, 1)
6	(Y, 6, 2)
7	RLE Triplet
8	- Value
9	- Offset
	- Length

Run-Length Encoding

```
SELECT isDead, COUNT(*)  
FROM users  
GROUP BY isDead
```



Compressed Data

id	isDead
1	(Y, 0, 3)
2	(N, 3, 1)
3	(Y, 4, 1)
4	(N, 5, 1)
6	(Y, 6, 2)
7	<i>RLE Triplet</i> - Value - Offset - Length
8	
9	

Run-Length Encoding

Original Data

id	isDead
1	Y
2	Y
3	Y
4	N
6	Y
7	N
8	Y
9	Y



Compressed Data

id	isDead
1	(Y, 0, 3)
2	(N, 3, 1)
3	(Y, 4, 1)
4	(N, 5, 1)
6	(Y, 6, 2)
7	RLE Triplet
8	- Value
9	- Offset
	- Length

Run-Length Encoding

Original Data

id	isDead
1	Y
2	Y
3	Y
4	N
6	Y
7	N
8	Y
9	Y



Compressed Data

id	isDead
1	(Y, 0, 3)
2	(N, 3, 1)
3	(Y, 4, 1)
4	(N, 5, 1)
6	(Y, 6, 2)
7	RLE Triplet
8	- Value
9	- Offset
	- Length

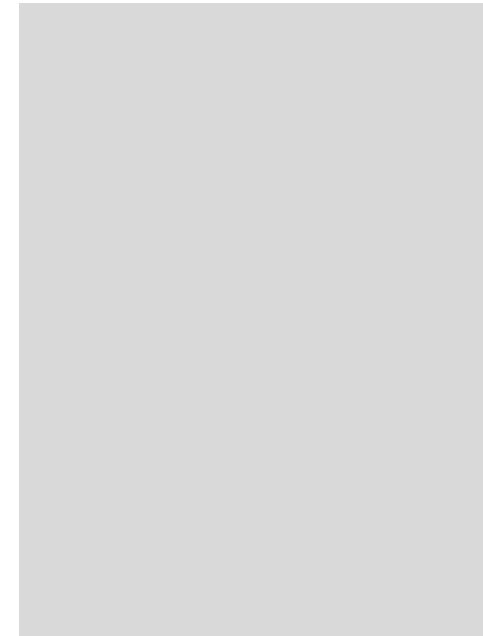
Run-Length Encoding

Sorted Data

id	isDead
1	Y
2	Y
3	Y
6	Y
8	Y
9	Y
4	N
7	N



Compressed Data



Run-Length Encoding

Sorted Data

id	isDead
1	Y
2	Y
3	Y
6	Y
8	Y
9	Y
4	N
7	N



Compressed Data

id	isDead
1	(Y, 0, 6)
2	(N, 7, 2)
3	
6	
8	
9	
4	
7	

Bit Packing

- If the values for an integer attribute is smaller than the range of its given data type size, then reduce the number of bits to represent each value.
- Use bit-shifting tricks to operate on multiple values in a single word.

Original Data

int32
13
191
56
92
81
120
231
172

Bit Packing

- If the values for an integer attribute is smaller than the range of its given data type size, then reduce the number of bits to represent each value.
- Use bit-shifting tricks to operate on multiple values in a single word.

Original Data

int32	
13	00000000 00000000 00000000 00001101
191	00000000 00000000 00000000 10111111
56	00000000 00000000 00000000 00111000
92	00000000 00000000 00000000 01011100
81	00000000 00000000 00000000 01010001
120	00000000 00000000 00000000 01111000
231	00000000 00000000 00000000 11100111
172	00000000 00000000 00000000 10101100

Bit Packing

- If the values for an integer attribute is smaller than the range of its given data type size, then reduce the number of bits to represent each value.
- Use bit-shifting tricks to operate on multiple values in a single word.

Original Data

Original:
 $8 \times 32\text{-bits} =$
256 bits

int32	
13	00000000 00000000 00000000 00001101
191	00000000 00000000 00000000 10111111
56	00000000 00000000 00000000 00111000
92	00000000 00000000 00000000 01011100
81	00000000 00000000 00000000 01010001
120	00000000 00000000 00000000 01111000
231	00000000 00000000 00000000 11100111
172	00000000 00000000 00000000 10101100

Bit Packing

- If the values for an integer attribute is smaller than the range of its given data type size, then reduce the number of bits to represent each value.
- Use bit-shifting tricks to operate on multiple values in a single word.

Original Data

Original:
 $8 \times 32\text{-bits} =$
256 bits

int32	
13	00000000 00000000 00000000 00001101
191	00000000 00000000 00000000 10111111
56	00000000 00000000 00000000 00111000
92	00000000 00000000 00000000 01011100
81	00000000 00000000 00000000 01010001
120	00000000 00000000 00000000 01111000
231	00000000 00000000 00000000 11100111
172	00000000 00000000 00000000 10101100

Bit Packing

- If the values for an integer attribute is smaller than the range of its given data type size, then reduce the number of bits to represent each value.
- Use bit-shifting tricks to operate on multiple values in a single word.

Original Data

int32	
13	00000000 00000000 00000000 00001101
191	00000000 00000000 00000000 10111111
56	00000000 00000000 00000000 00111000
92	00000000 00000000 00000000 01011100
81	00000000 00000000 00000000 01010001
120	00000000 00000000 00000000 01111000
231	00000000 00000000 00000000 11100111
172	00000000 00000000 00000000 10101100

Original:
 $8 \times 32\text{-bits} =$
256 bits

Bit Packing

- If the values for an integer attribute is smaller than the range of its given data type size, then reduce the number of bits to represent each value.
- Use bit-shifting tricks to operate on multiple values in a single word.

Original Data

int32	
13	00001101
191	10111111
56	00111000
92	01011100
81	01010001
120	01111000
231	11100111
172	10101100

Original:
 $8 \times 32\text{-bits} =$
256 bits

Compressed:
 $8 \times 8\text{-bits} =$
64 bits

Patching / Mostly Encoding

- A variation of bit packing when attribute's values are “mostly” less than the largest size, store them with the smaller data type.
→ The remaining values that cannot be compressed are stored in their raw form.

Original Data

int32
13
191
99999999
92
81
120
231
172

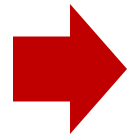
Source: [Redshift Documentation](#)

Patching / Mostly Encoding

- A variation of bit packing when attribute's values are “mostly” less than the largest size, store them with the smaller data type.
→ The remaining values that cannot be compressed are stored in their raw form.

Original Data

int32
13
191
99999999
92
81
120
231
172



Compressed Data

mostly8	offset	value
13	3	99999999
181		
XXX		
92		
81		
120		
231		
172		

Source: [Redshift Documentation](#)

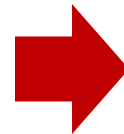
Patching / Mostly Encoding

- A variation of bit packing when attribute's values are “mostly” less than the largest size, store them with the smaller data type.
→ The remaining values that cannot be compressed are stored in their raw form.

Original:
 $8 \times 32\text{-bits} =$
256 bits

Original Data

int32
13
191
99999999
92
81
120
231
172



Compressed Data

mostly8	offset	value
13	3	99999999
181		
XXX		
92		
81		
120		
231		
172		

Compressed:
 $(8 \times 8\text{-bits}) +$
 $16\text{-bits} + 32\text{-bits}$
= 112 bits

Source: [Redshift Documentation](#)

Bitmap Encoding

- Store a separate bitmap for each unique value for an attribute where an offset in the vector corresponds to a tuple.
 - The i^{th} position in the Bitmap corresponds to the i^{th} tuple in the table.
 - Typically segmented into chunks to avoid allocating large blocks of contiguous memory.
- Only practical if the value cardinality is low.
- Some DBMSs provide bitmap indexes.

Bitmap Encoding

Original Data

id	isDead
1	Y
2	Y
3	Y
4	N
6	Y
7	N
8	Y
9	Y

Bitmap Encoding

Original Data

id	isDead
1	Y
2	Y
3	Y
4	N
6	Y
7	N
8	Y
9	Y

Bitmap Encoding

Original Data

id	isDead
1	Y
2	Y
3	Y
4	N
6	Y
7	N
8	Y
9	Y



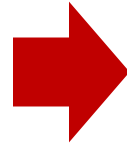
Compressed Data

id	isDead	
	Y	N
1	1	0
2	1	0
3	1	0
4	0	1
6	1	0
7	0	1
8	1	0
9	1	0

Bitmap Encoding

Original Data

id	isDead
1	Y
2	Y
3	Y
4	N
6	Y
7	N
8	Y
9	Y



Compressed Data

id	isDead	
	Y	N
1	1	0
2	1	0
3	1	0
4	0	1
6	1	0
7	0	1
8	1	0
9	1	0

Bitmap Encoding

Original Data

id	isDead
1	Y
2	Y
3	Y
4	N
6	Y
7	N
8	Y
9	Y

Original:
 $8 \times 8\text{-bits} =$
64 bits

Compressed Data

id	isDead	
	Y	N
1	1	0
2	1	0
3	1	0
4	0	1
6	1	0
7	0	1
8	1	0
9	1	0

Bitmap Encoding

Original Data

id	isDead
1	Y
2	Y
3	Y
4	N
6	Y
7	N
8	Y
9	Y

Original:
 $8 \times 8\text{-bits} =$
 64 bits

Compressed Data

id	isDead	
	Y	N
1	1	0
2	1	0
3	1	0
4	0	1
6	1	0
7	0	1
8	1	0
9	1	0

$2 \times 8\text{-bits} =$
 16 bits

Bitmap Encoding

Original Data

id	isDead
1	Y
2	Y
3	Y
4	N
6	Y
7	N
8	Y
9	Y

Original:
 $8 \times 8\text{-bits} =$
64 bits

Compressed Data

id	isDead	
	Y	N
1	1	0
2	1	0
3	1	0
4	0	1
6	1	0
7	0	1
8	1	0
9	1	0

Compressed:
16 bits + 18
bits = 34 bits

$2 \times 8\text{-bits} =$
16 bits

$9 \times 2\text{-bits} =$
18 bits

Bitmap Encoding: Example

- Assume we have 10 million tuples.
43,000 zip codes in the US.
→ $10000000 \times 32\text{-bits} = 40 \text{ MB}$
→ $10000000 \times 43000 = 53.75 \text{ GB}$
- Every time the application inserts a new tuple, the DBMS must extend 43,000 different bitmaps.

```
CREATE TABLE customer (  
    id INT PRIMARY KEY,  
    name VARCHAR(32),  
    email VARCHAR(64),  
    address VARCHAR(64),  
    zip_code INT  
);
```


Bitmap Encoding: Example

- Assume we have 10 million tuples.
43,000 zip codes in the US.
→ $10000000 \times 32\text{-bits} = 40 \text{ MB}$
→ $10000000 \times 43000 = 53.75 \text{ GB}$
- Every time the application inserts a new tuple, the DBMS must extend 43,000 different bitmaps.

```
CREATE TABLE customer (  
    id INT PRIMARY KEY,  
    name VARCHAR(32),  
    email VARCHAR(64),  
    address VARCHAR(64),  
    zip_code INT  
);
```

Delta Encoding

- Recording the difference between values that follow each other in the same column.
 - Store base value in-line or in a separate look-up table.
 - Combine with RLE to get even better compression ratios.

Original Data

time64	temp
12:00	99.5
12:01	99.4
12:02	99.5
12:03	99.6
12:04	99.4

Delta Encoding

- Recording the difference between values that follow each other in the same column.
 - Store base value in-line or in a separate look-up table.
 - Combine with RLE to get even better compression ratios.

Original Data

time64	temp
12:00	99.5
12:01	99.4
12:02	99.5
12:03	99.6
12:04	99.4

Delta Encoding

- Recording the difference between values that follow each other in the same column.
 - Store base value in-line or in a separate look-up table.
 - Combine with RLE to get even better compression ratios.

Original Data

time64	temp
12:00	99.5
12:01	99.4
12:02	99.5
12:03	99.6
12:04	99.4



Compressed Data

time64	temp
12:00	99.5
+1	-0.1
+1	+0.1
+1	+0.1
+1	-0.2

Delta Encoding

- Recording the difference between values that follow each other in the same column.
 - Store base value in-line or in a separate look-up table.
 - Combine with RLE to get even better compression ratios.

Original Data

time64	temp
12:00	99.5
12:01	99.4
12:02	99.5
12:03	99.6
12:04	99.4

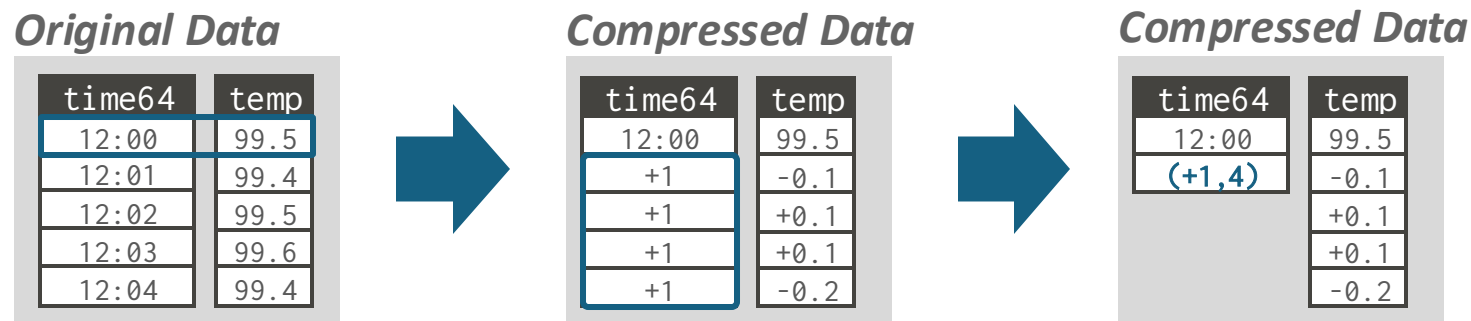


Compressed Data

time64	temp
12:00	99.5
+1	-0.1
+1	+0.1
+1	+0.1
+1	-0.2

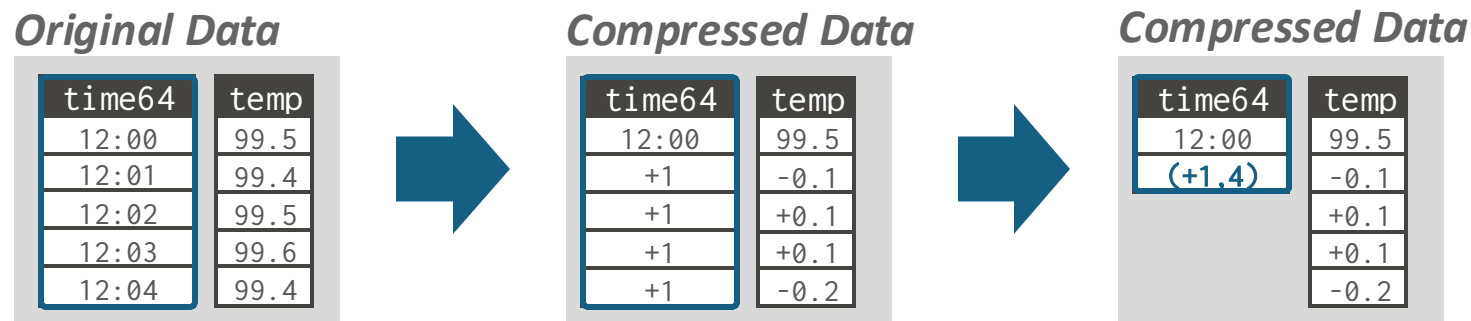
Delta Encoding

- Recording the difference between values that follow each other in the same column.
 - Store base value in-line or in a separate look-up table.
 - Combine with RLE to get even better compression ratios.



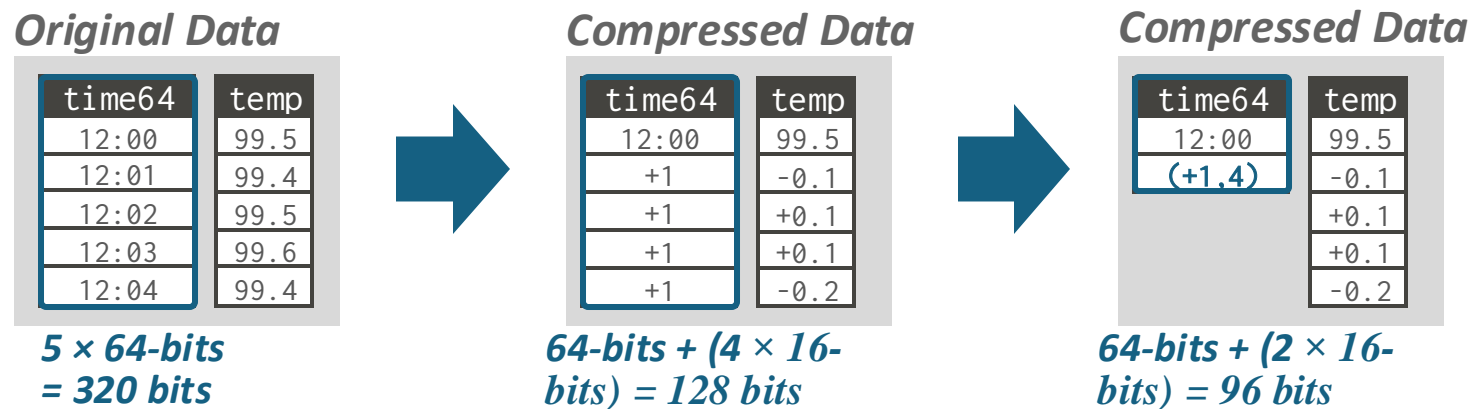
Delta Encoding

- Recording the difference between values that follow each other in the same column.
 - Store base value in-line or in a separate look-up table.
 - Combine with RLE to get even better compression ratios.



Delta Encoding

- Recording the difference between values that follow each other in the same column.
- Store base value in-line or in a separate look-up table.
- Combine with RLE to get even better compression ratios.



Dictionary Compression

- Replace frequent values with smaller fixed-length codes and then maintain a mapping (dictionary) from the codes to the original values
 - Typically, one code per attribute value.
 - Most widely used native compression scheme in DBMSs.
- The ideal dictionary scheme supports fast encoding and decoding for both point and range queries.

Dictionary: Example

Original Data

name
Andrea
Prashanth
Andy
Matt
Prashanth

Dictionary: Example

Original Data

name
Andrea
Prashanth
Andy
Matt
Prashanth



Compressed Data

name	value	code
10	Andrea	10
20	Prashanth	20
30	Andy	30
40	Matt	40
20		

Dictionary: Example

Original Data

name
Andrea
Prashanth
Andy
Matt
Prashanth



Compressed Data

name	value	code
10	Andrea	10
20	Prashanth	20
30	Andy	30
40	Matt	40
20		

Dictionary

Dictionary: Example

```
SELECT * FROM users
WHERE name = 'Andy'
```

Original Data

name
Andrea
Prashanth
Andy
Matt
Prashanth



Compressed Data

name	value	code
10	Andrea	10
20	Prashanth	20
30	Andy	30
40	Matt	40
20		

Dictionary

Dictionary: Example

```
SELECT * FROM users
WHERE name = 'Andy'
```



```
SELECT * FROM users
WHERE name = 30
```

Original Data

name
Andrea
Prashanth
Andy
Matt
Prashanth



Compressed Data

name	value	code
10	Andrea	10
20	Prashanth	20
30	Andy	30
40	Matt	40
20		

Dictionary

Dictionary: Encoding / Decoding

- A dictionary needs to support two operations:
 - **Encode/Locate:** For a given uncompressed value, convert it into its compressed form.
 - **Decode/Extract:** For a given compressed value, convert it back into its original form.
- No magic hash function will do this for us.

Dictionary: Order-Preserving

- The encoded values need to support the same collation as the original values.

Original Data

name
Andrea
Prashanth
Andy
Matt
Prashanth

Dictionary: Order-Preserving

- The encoded values need to support the same collation as the original values.

Original Data

name
Andrea
Prashanth
Andy
Matt
Prashanth



Compressed Data

name	value	code
10	Andrea	10
40	Andy	20
20	Matt	30
30	Prashanth	40
40		

Dictionary: Order-Preserving

- The encoded values need to support the same collation as the original values.

Original Data

name
Andrea
Prashanth
Andy
Matt
Prashanth



Compressed Data

name	value	code
10	Andrea	10
40	Andy	20
20	Matt	30
30	Prashanth	40
40		

*Sorted
Dictionary*

Dictionary: Order-Preserving

- The encoded values need to support the same collation as the original values.

```
SELECT * FROM users  
WHERE name LIKE 'And%'
```

Original Data

name
Andrea
Prashanth
Andy
Matt
Prashanth



Compressed Data

name	value	code
10	Andrea	10
40	Andy	20
20	Matt	30
30	Prashanth	40
40		

*Sorted
Dictionary*

Dictionary: Order-Preserving

- The encoded values need to support the same collation as the original values.

```
SELECT * FROM users  
WHERE name LIKE 'And%'
```



```
SELECT * FROM users  
WHERE name BETWEEN 10 AND 20
```

Original Data

name
Andrea
Prashanth
Andy
Matt
Prashanth



Compressed Data

name	value	code
10	Andrea	10
40	Andy	20
20	Matt	30
30	Prashanth	40
40		

*Sorted
Dictionary*

Order-Preserving Encoding

Original Data

name
Andrea
Prashanth
Andy
Matt
Prashanth



Compressed Data

<i>name</i>	<i>value</i>	<i>code</i>
10	Andrea	10
40	Andy	20
20	Matt	30
30	Prashanth	40
40		

*Sorted
Dictionary*

Order-Preserving Encoding

```
SELECT name FROM users  
WHERE name LIKE 'And%'
```



Still must perform scan on column

Original Data

name
Andrea
Prashanth
Andy
Matt
Prashanth



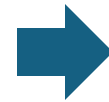
Compressed Data

name	value	code
10	Andrea	10
40	Andy	20
20	Matt	30
30	Prashanth	40
40		

*Sorted
Dictionary*

Order-Preserving Encoding

```
SELECT name FROM users
WHERE name LIKE 'And%'
```



Still must perform scan on column

```
SELECT DISTINCT name
FROM users
WHERE name LIKE 'And%'
```



Only need to access dictionary

Original Data

name
Andrea
Prashanth
Andy
Matt
Prashanth



Compressed Data

name	value	code
10	Andrea	10
40	Andy	20
20	Matt	30
30	Prashanth	40
40		

*Sorted
Dictionary*

Dictionary: Data Structures

- **Choice #1: Array**

- One array of variable length strings and another array with pointers that maps to string offsets.
- Expensive to update so only usable in immutable files.

- **Choice #2: Hash Table**

- Fast and compact.
- Unable to support range and prefix queries.

- **Choice #3: B+Tree**

- Slower than a hash table and takes more memory.
- Can support range and prefix queries.

Dictionary: Array

- First sort the values and then store them sequentially in a byte array.
 - Need to also store the size of the value if they are variable-length.
- Replace the original data with dictionary codes that are the (byte) offset into this array.

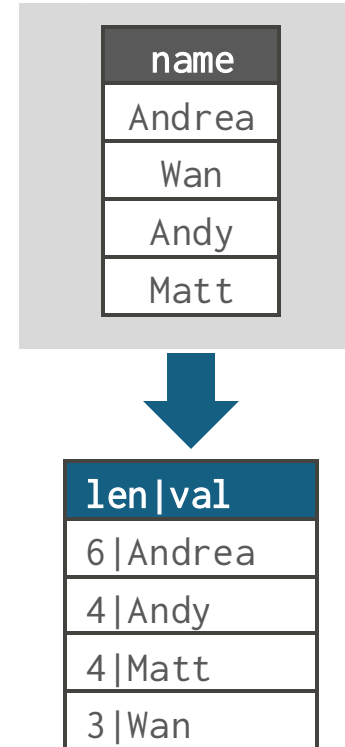
Original Data

name
Andrea
Wan
Andy
Matt

Dictionary: Array

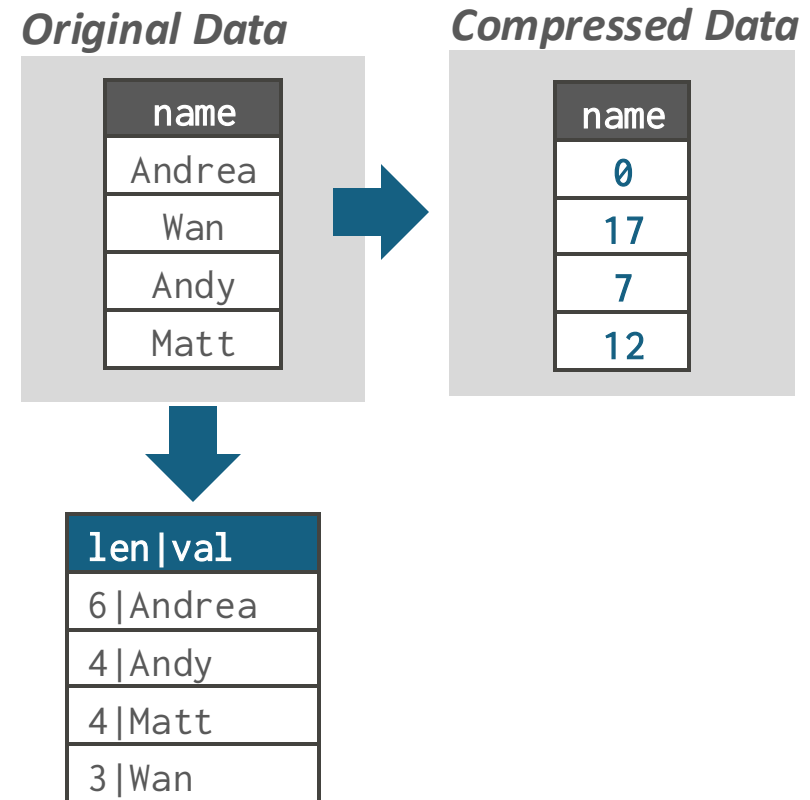
- First sort the values and then store them sequentially in a byte array.
 - Need to also store the size of the value if they are variable-length.
- Replace the original data with dictionary codes that are the (byte) offset into this array.

Original Data



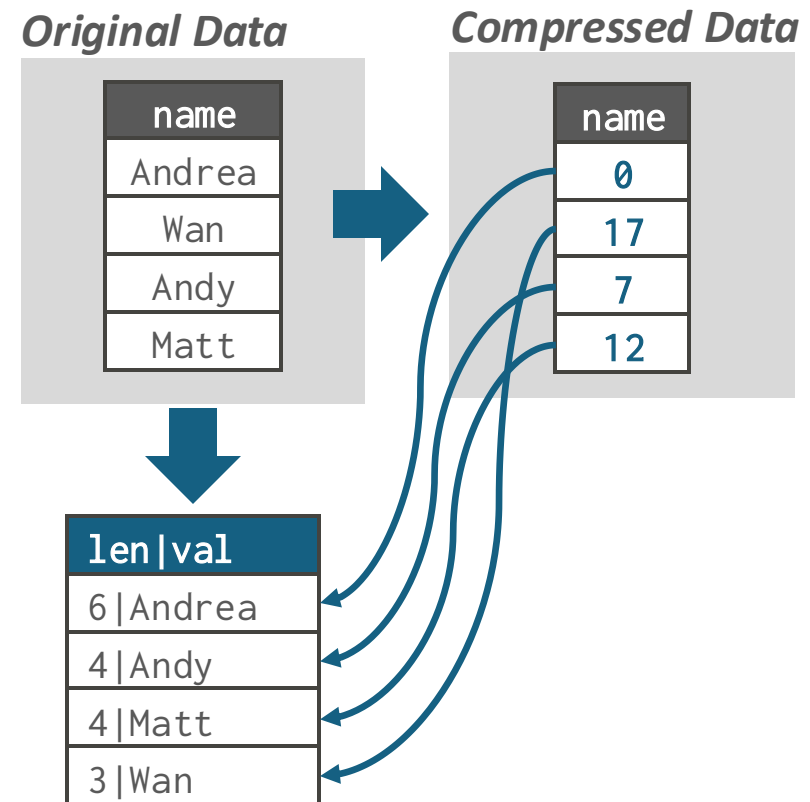
Dictionary: Array

- First sort the values and then store them sequentially in a byte array.
 - Need to also store the size of the value if they are variable-length.
- Replace the original data with dictionary codes that are the (byte) offset into this array.



Dictionary: Array

- First sort the values and then store them sequentially in a byte array.
 - Need to also store the size of the value if they are variable-length.
- Replace the original data with dictionary codes that are the (byte) offset into this array.



Conclusion

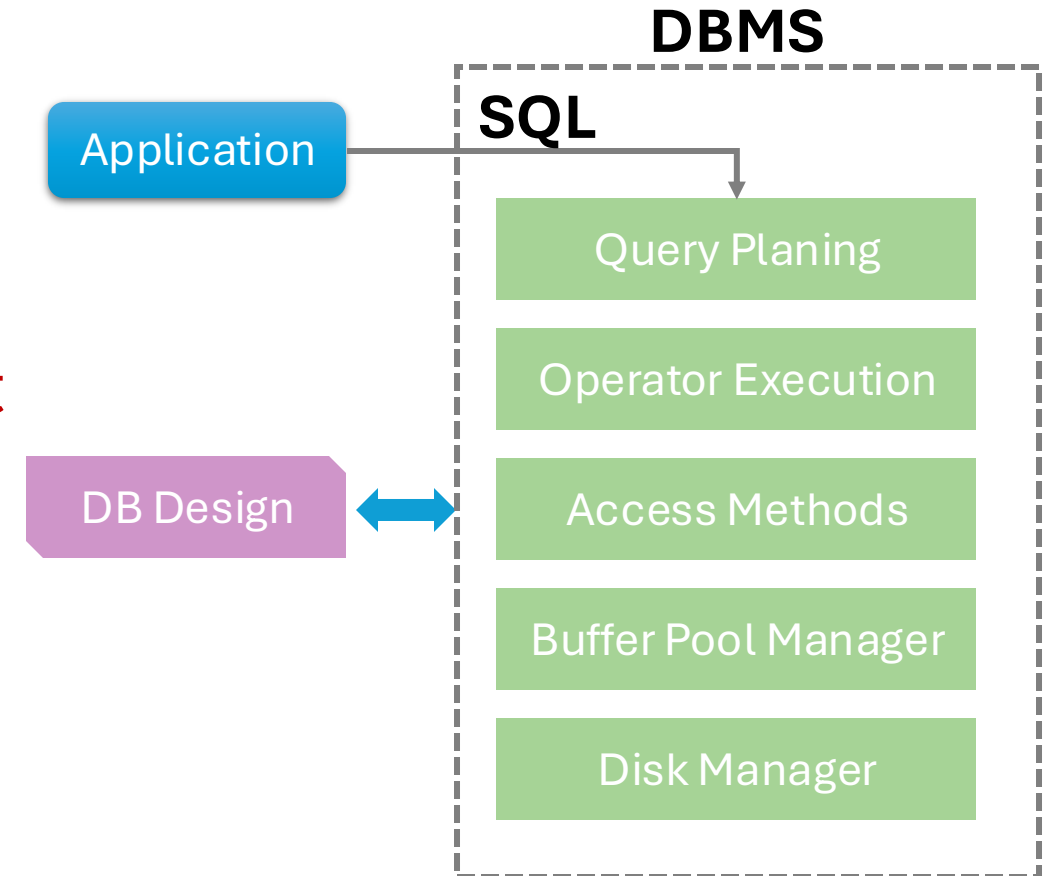
- It is important to choose the right storage model for the target workload:
 - OLTP = Row Store
 - OLAP = Column Store
- DBMSs can combine different approaches for even better compression.
- Dictionary encoding is probably the most useful scheme because it does not require pre-sorting.

Database Storage

- **Problem #1:** How the DBMS represents the database in files on disk.

- **Problem #2:** How the DBMS manages its memory and moves data back-and-forth from disk.

← Next



Database Storage

- **Problem #1:** How the DBMS represents the database in files on disk.

- **Problem #2:** How the DBMS manages its memory and moves data back-and-forth from disk.

← Next

