

OLTP → Small operations → Fast operations → read/update on small data

DATABASE WORKLOADS

Some operation that may/may not bring change in db.

On-Line Transaction Processing (OLTP)

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

OLAP → Larger operations → Touches time → read lot of data.

On-Line Analytical Processing (OLAP)

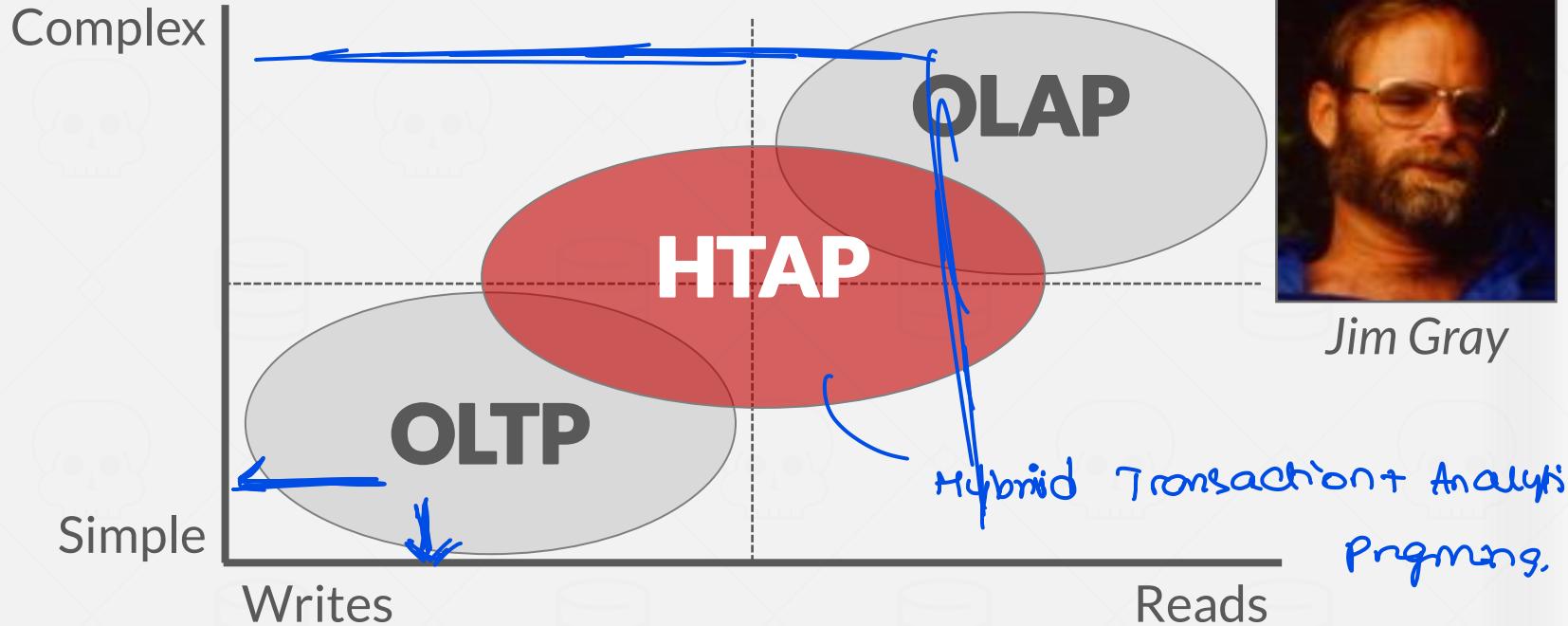
- Complex queries that read a lot of data to compute aggregates.

In Hybrid Transaction + Analytical Processing
we use OLTP + OLAP together on the same database instance

Operation Complexity

OLAP does read on huge sizes of data

DATABASE WORKLOADS



Jim Gray

Hybrid Transaction+Analytic Processing.

If we see OLTP does more writes on small amount of data.
Workload Focus

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
);
```

OBSERVATION

The relational model does not specify that the DBMS must store all a tuple's attributes together in a single page.

*This may not actually be the best layout for some workloads...

Relational model of DB does not specify that all attributes of a tuple need not be stored together in single page.

OLTP

most of the time this happens from above graph

On-line Transaction Processing:

→ Simple queries that read/update a small amount of data that is related to a single entity in the database.

```
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 you have 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)
```

OLAP is executed on the data collected

from OLTP applications

DATA STORAGE MODELS

The DBMS can store tuples in different ways that are better for either OLTP or OLAP workloads.

We have been assuming the n-ary storage model (aka "row storage") so far this semester.

n-ary storage model \Rightarrow attributes for a tuple are stored continuously in row in single page.

N-ARY STORAGE MODEL (NSM)

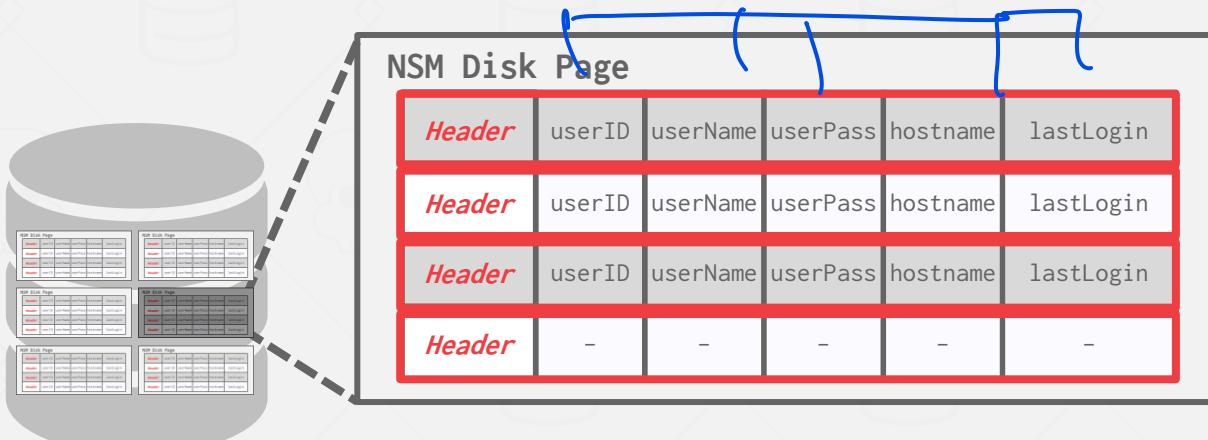
The DBMS stores all attributes for a single tuple contiguously in a page.

Ideal for OLTP workloads where queries tend to operate only on an individual entity and insert-heavy workloads.

array or
row storage

N-ARY STORAGE MODEL (NSM)

The DBMS stores all attributes for a single tuple contiguously in a page.



ALL

N-ARY STORAGE MODEL (NSM)

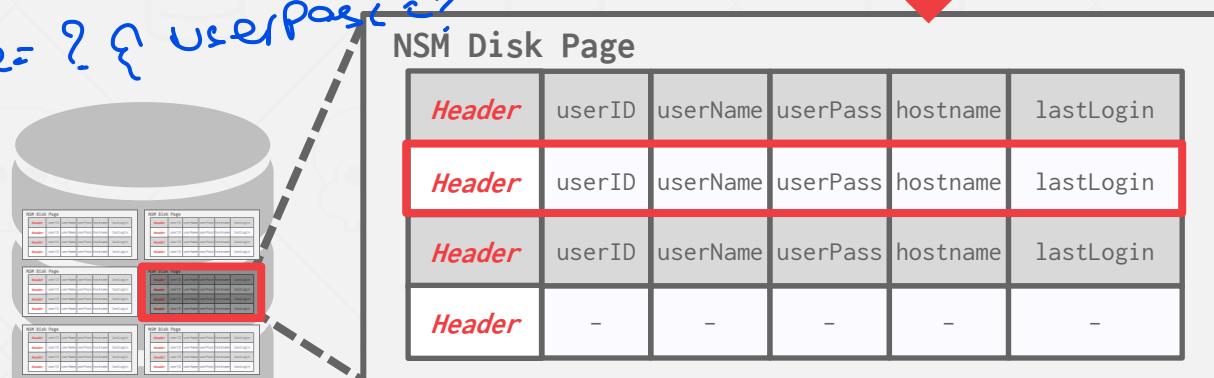
`SELECT * FROM useracct
WHERE userName = ?
AND userPass = ?`

Table name:

Lecture #8

Index

Select all tuples with
userName= ? & userPass = ?



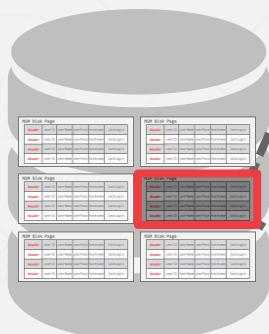
N-ARY STORAGE MODEL (NSM)

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

```
INSERT INTO useracct  
VALUES (?, ?, ...?)
```

Lecture #8

Index



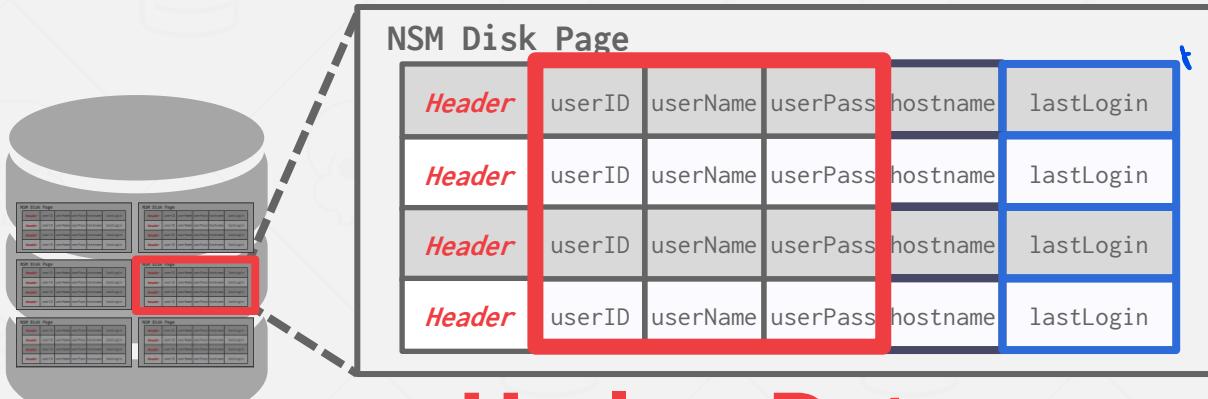
NSM Disk Page

| | | | | | |
|---------------|--------|----------|----------|----------|-----------|
| <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 |

N-ARY STORAGE MODEL (NSM)

This is
a OLAP
transaction

```
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)
```



Useless Data

hostname &
lastLogin
↓
only these
attributes are
req for this
transaction
~ Row 3
are useless
↓
Storage waste

N-ARY STORAGE MODEL

Advantages

- Fast inserts, updates, and deletes.
- Good for queries that need the entire tuple.

Disadvantages

- Not good for scanning large portions of the table and/or a subset of the attributes.

★ OLTP transactions prefer

- NMR

now storage

↳ above example.

DECOMPOSITION STORAGE MODEL (DSM)

The DBMS stores the values of a single attribute for all tuples contiguously in a page.
→ Also known as a "column store"

⇒ ∴ Each page
contains data
about
one
attribute

Ideal for OLAP workloads where read-only
queries perform large scans over a subset of the
table's attributes.

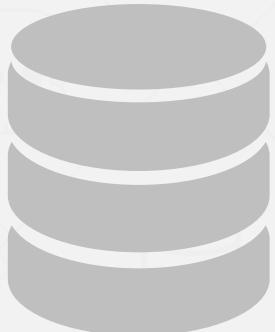
above example
where only hostname & lastlogin This page has all tuple data
one required

DECOMPOSITION STORAGE MODEL (DSM)

DSM
 ↓
 column storage.

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

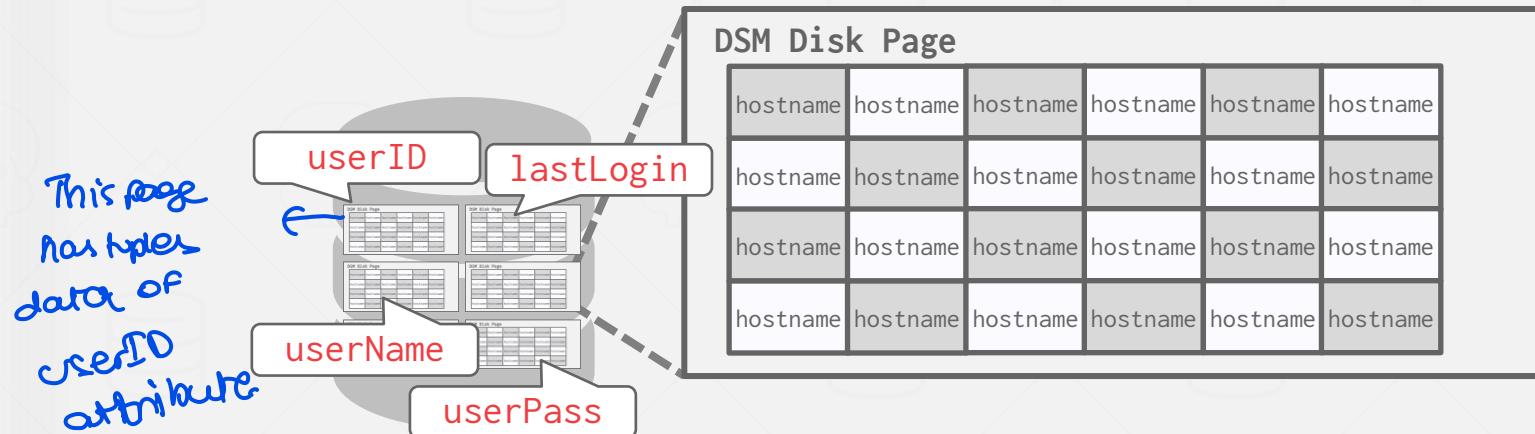
NSX|
 ↓
 rowstorage



| | | | | | |
|---------------|--------|----------|----------|----------|-----------|
| <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 |

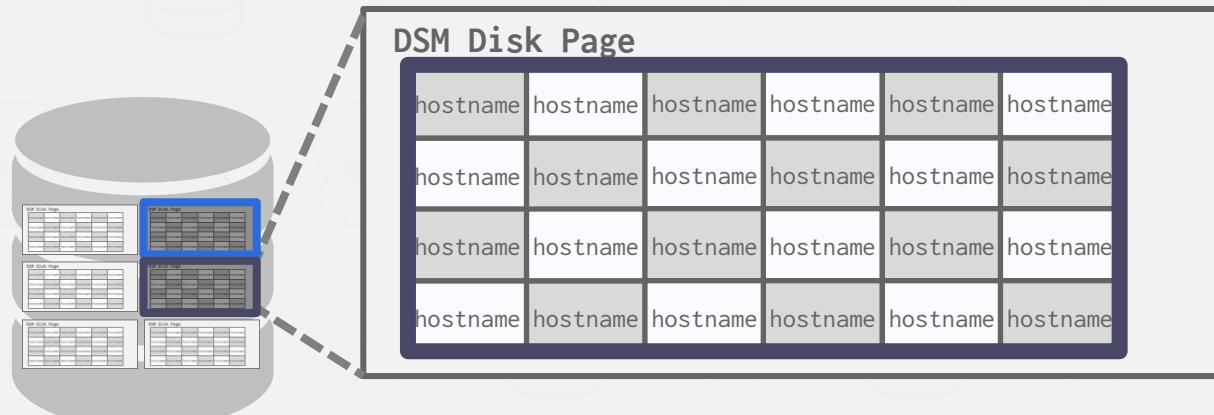
DECOMPOSITION STORAGE MODEL (DSM)

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



DECOMPOSITION STORAGE MODEL (DSM)

```
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)
```



now we can touch only 2 pages for the OLAP transaction

i.e. for a single

attribute all
values have some
length.

TUPLE IDENTIFICATION

Choice #1: Fixed-length Offsets

$\lambda = \gamma = \varphi = w$ → 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 | 2 | | | |
| 1 | 4 | | | |
| 2 | 2 | | | |
| 3 | w | | | |

Embedded Ids

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

DECOMPOSITION STORAGE MODEL (DSM)

Advantages

- Reduces the amount wasted I/O because the DBMS only reads the data that it needs.
- Better query processing and data compression (more on this later).

Disadvantages

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

∴ For insert we need to go to all pages and add ~~data~~ at each attribute.

DSM SYSTEM HISTORY

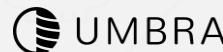
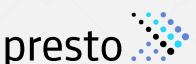
1970s: Cantor DBMS

1980s: DSM Proposal

1990s: SybaseIQ (in-memory only)

2000s: Vertica, VectorWise, MonetDB

2010s: Everyone



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.

AS we ↑ compression ratio storage space ↓
speed ↑-

REAL-WORLD DATA CHARACTERISTICS

Data sets tend to have highly skewed distributions for attribute values.

→ Example: Zipfian distribution of the [Brown Corpus](#)

Data sets tend to have high correlation between attributes of the same tuple.

→ Example: Zip Code to City, Order Date to Ship Date

DATABASE COMPRESSION

↑ for easy tuple identification

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.

∴ after compression,

when we do a query operation our

Goal #3: Must be a lossless scheme.

No data
must be lost
during
compression

Goal 3 is to postpone decompression for as long as it is possible.

LOSSLESS VS. LOSSY COMPRESSION

When a DBMS uses compression, it is always lossless because people don't like losing data.
AFA

Any kind of lossy compression must be performed at the application level.

File → Page → Block

COMPRESSION GRANULARITY

Choice #1: Block-level

- Compress a block of tuples for the same table.

→ what do u
mean by

block of
tuples

?

↓

Set of tuples .

Choice #2: Tuple-level

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

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).

NAÏVE COMPRESSION

Compress data using a general-purpose algorithm.
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.

- In MySQL
we compress
our pages

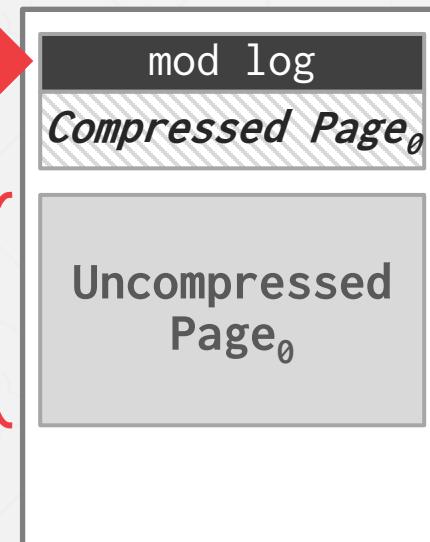
then when query comes we store all updates that are to be done in the

MYSQL INNODB COMPRESSION

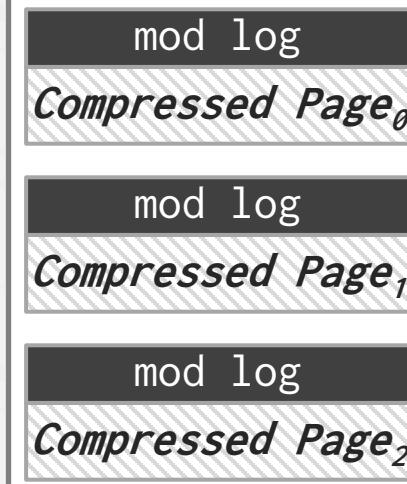
Buffer Pool

Updates

16 KB



Disk Pages



decompression we do all updates.

mod log
that is
[1,2,4,8] KB
present with
every compres-
sed page.

Then at

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.

Problem If writer then redlog can show but for a read operation
DBMS needs to decompress the page ∴ This limit the scope of
Compression Schema.

OBSERVATION

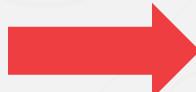


Ideally, we want the DBMS to operate on compressed data without decompressing it first.

1. If we can transfer our query like this then no need to take compressed data

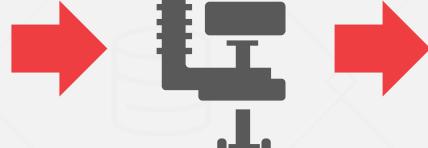
Database Magic!

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



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

| NAME | SALARY |
|------|--------|
| Andy | 99999 |
| Matt | 88888 |



| NAME | SALARY |
|------|--------|
| XX | AA |
| YY | BB |

Compression.

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).

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).

∴ we compress data about 1/more attributes

COLUMNAR COMPRESSION

→ Done for

only DSM

∴ we operate on

subset of

attributes.

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

RUN-LENGTH ENCODING

↑ tuples

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.

we group into triplet. It has ① The value ② Offset => Start point

RUN-LENGTH ENCODING

③ Size or length.

Original Data

For an attribute what data the tuple has

↓
value-

Tuple ↪

↳

| id | sex |
|----|-----|
| 1 | M |
| 2 | M |
| 3 | M |
| 4 | F |
| 6 | M |
| 7 | F |
| 8 | M |
| 9 | M |

These 3 are contiguous
∴ can be grouped

These 2 can also be grouped

On attribute

sex it has M → value.

Compressed Data

| id | sex |
|----|-----------|
| 1 | (M, 0, 3) |
| 2 | (F, 3, 1) |
| 3 | (M, 4, 1) |
| 4 | (F, 5, 1) |
| 6 | (M, 6, 2) |
| 7 | |
| 8 | |
| 9 | |

RLE Triplet

- Value
- Offset
- Length

RUN-LENGTH ENCODING

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



Compressed Data

| id | sex |
|----|-----------|
| 1 | (M, 0, 3) |
| 2 | (F, 3, 1) |
| 3 | (M, 4, 1) |
| 4 | (F, 5, 1) |
| 6 | (M, 6, 2) |
| 7 | |
| 8 | |
| 9 | |

RLE Triplet
- Value
- Offset
- Length

But instead
instead of
we now store
3 values in
Original Data

RUN-LENGTH ENCODING

which is
worse than
original.

| id | sex |
|----|-----|
| 1 | M |
| 2 | M |
| 3 | M |
| 4 | F |
| 6 | M |
| 7 | F |
| 8 | M |
| 9 | M |

There are not
continuous runs.

Compressed Data

| id | sex |
|----|-----------|
| 1 | (M, 0, 3) |
| 2 | (F, 3, 1) |
| 3 | (M, 4, 1) |
| 4 | (F, 5, 1) |
| 6 | (M, 6, 2) |
| 7 | |
| 8 | |
| 9 | |

Some way
are
stored
but triplet
↓
Data
waste.

- RLE Triplet**
- Value
- Offset
- Length

∴ If all are discrete then the compression results in some storage-

RUN-LENGTH ENCODING

∴ On sorted data

∴ continuous runs are

Sorted Data

possible
we
prefer
this

| id | sex |
|----|-----|
| 1 | M |
| 2 | M |
| 3 | M |
| 6 | M |
| 8 | M |
| 9 | M |
| 4 | F |
| 7 | F |



Compressed Data

| id | sex |
|----|-----------|
| 1 | (M, 0, 6) |
| 2 | |
| 3 | |
| 6 | |
| 8 | |
| 9 | |
| 4 | |
| 7 | |

BIT-PACKING ENCODING

When values for an attribute are always less than the value's declared largest size, store them as smaller data type.

Original Data

| int64 |
|-------|
| 2 |
| 4 |
| 45 |
| 6 |
| 18 |

→ Here attribute is defined to store 64 bits. But all values are storing less than about 6 bits.

$$\log_2(45) = 6$$

BIT-PACKING ENCODING

When values for an attribute are always less than the value's declared largest size, store them as smaller data type.

- Here we compress the data type? Store them -

Original Data

| int64 |
|--|
| 0010 |
| 0011 |
| 00 |
| 00 |
| 00011001 |
| 00 |
| 0010 |
| 00 |
| 00010010 |

5 × 64-bits = 320 bits



Compressed Data

| packed-int8 |
|-------------|
| 00000010 |
| 00000011 |
| 00011101 |
| 00000110 |
| 00010010 |

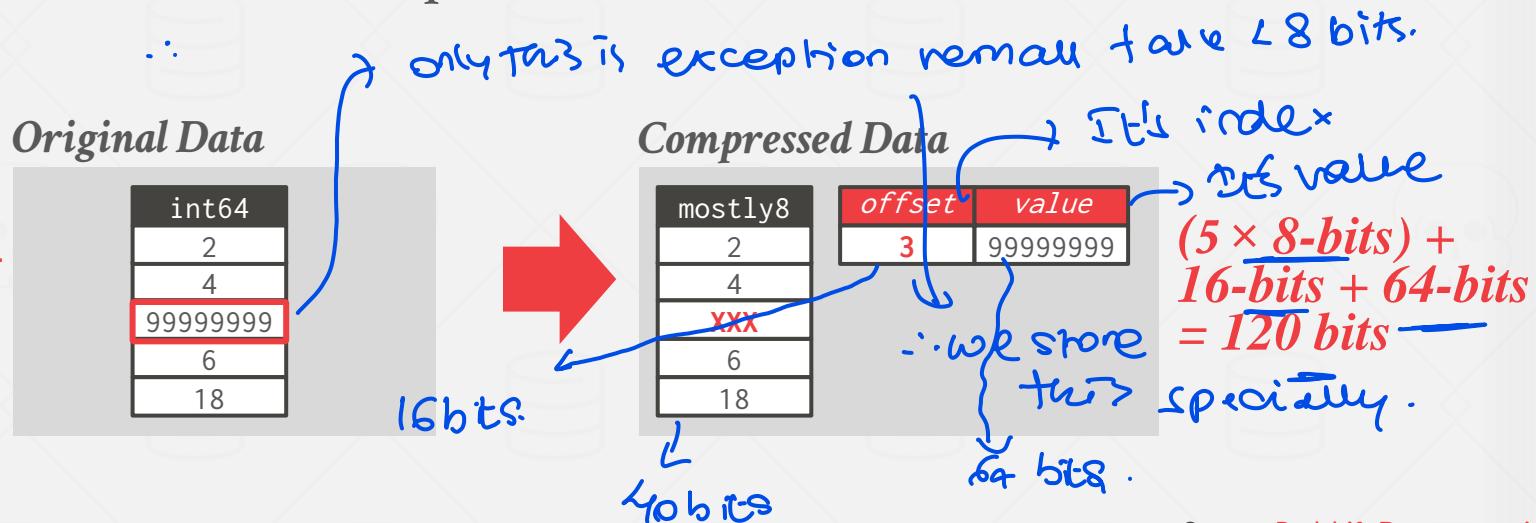
(5 × 8-bits) = 40 bits

- now only 8 bits are required to store each value -

Why the
1 & bits
here ??

MOSTLY ENCODING

Bit-packing variant that uses a special marker to indicate when a value exceeds largest size and then maintain a look-up table to store them.



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.

types of values that attribute can store are low

Some DBMSs provide bitmap indexes.

Ex sex male or female

zip code → 6 digits → diff for each place.

we don't
prefer to
apply Bitmap encoding to

BITMAP ENCODING

Original Data

| id | sex |
|----|-----|
| 1 | M |
| 2 | M |
| 3 | M |
| 4 | F |
| 6 | M |
| 7 | F |
| 8 | M |
| 9 | M |

attributes with

Compressed Data

primary key
constraint
because each
value
→ different



| id | M | F |
|----|---|---|
| 1 | 1 | 0 |
| 2 | 1 | 0 |
| 3 | 1 | 0 |
| 4 | 0 | 1 |
| 6 | 1 | 0 |
| 7 | 0 | 1 |
| 8 | 1 | 0 |
| 9 | 1 | 0 |

If n tuples
→ n columns
must be
created in bitmap.

1 char → 1 byte → 8 bits.

BITMAP ENCODING

Original Data

| id | sex |
|----|-----|
| 1 | M |
| 2 | M |
| 3 | M |
| 4 | F |
| 6 | M |
| 7 | F |
| 8 | M |
| 9 | M |

→ 8 bits.
 $9 \times 8\text{-bits} = 72\text{ bits}$

Compressed Data

| id | sex | M | F |
|----|-----|---|---|
| 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 chars M and F
each 8 bits

8 Columns

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

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

16

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

Each zipcode req. 31 bit ∵ 40 million tuples req

$10 \text{ million} \times 32 = 40 \text{ mb}$

but for each zipcode we give a bitmap size very huge.

the difference 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.

Better for
continuous
value

Original Data

| time | temp |
|-------|------|
| 12:00 | 99.5 |
| 12:01 | 99.4 |
| 12:02 | 99.5 |
| 12:03 | 99.6 |
| 12:04 | 99.4 |

Compressed Data

| time | temp |
|-------|------|
| 12:00 | 99.5 |
| +1 | -0.1 |
| +1 | +0.1 |
| +1 | +0.1 |
| +1 | -0.2 |

Delta

Compressed Data

| time | temp |
|---------|------|
| 12:00 | 99.5 |
| (+1, 4) | -0.1 |
| +0.1 | +0.1 |
| +0.1 | -0.2 |

Runlength.

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

| time | temp |
|-------|------|
| 12:00 | 99.5 |
| 12:01 | 99.4 |
| 12:02 | 99.5 |
| 12:03 | 99.6 |
| 12:04 | 99.4 |

$5 \times 32\text{-bits}$
 $= 160 \text{ bits}$

Compressed Data

| time | temp |
|-------|------|
| 12:00 | 99.5 |
| +1 | -0.1 |
| +1 | +0.1 |
| +1 | +0.1 |
| +1 | -0.2 |

$32\text{-bits} + (4 \times 16\text{-bits})$
 $= 96 \text{ bits}$

Compressed Data

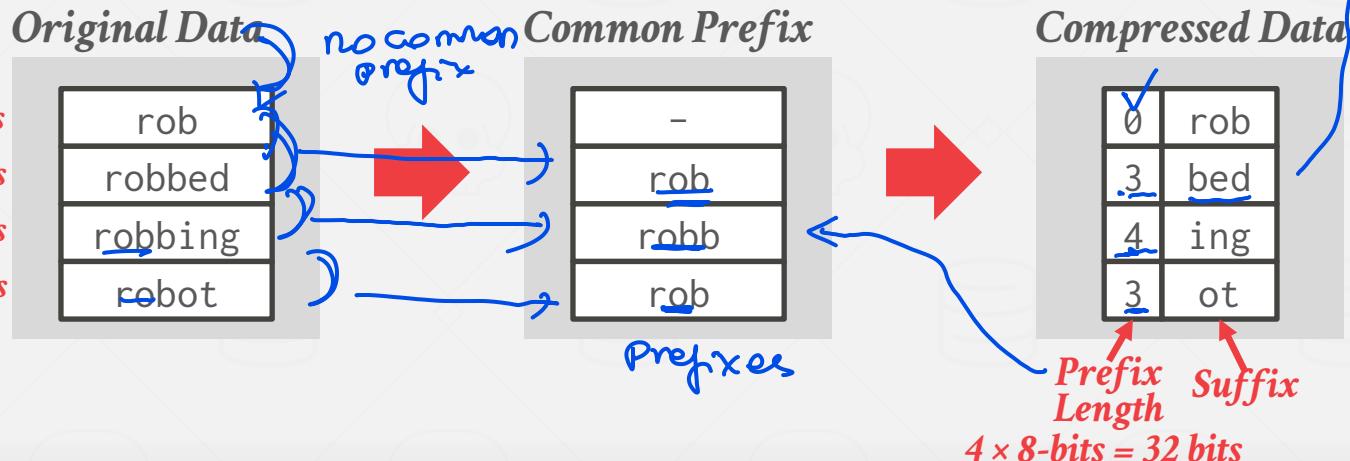
| time | temp |
|---------|------|
| 12:00 | 99.5 |
| (+1, 4) | -0.1 |
| +0.1 | +0.1 |
| +0.1 | -0.2 |

$32\text{-bits} + (2 \times 16\text{-bits})$
 $= 64 \text{ bits}$

INCREMENTAL ENCODING

Type of delta encoding that avoids duplicating common prefixes/suffixes between consecutive tuples. This works best with sorted data.

suffix or
left over data



DICTIONARY COMPRESSION

Build a data structure that maps variable-length values to a smaller integer identifier.

Replace those values with their corresponding identifier in the dictionary data structure.

- Need to support fast encoding and decoding.
- Need to also support range queries.

Most widely used compression scheme in DBMSs.



DICTIONARY COMPRESSION

We create a dictionary where each value (larger / variable) is

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



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

mapped to
a smaller

Code.

Original Data

| name |
|-----------|
| Andrea |
| Prashanth |
| Andy |
| Matt |
| Prashanth |



Compressed Data

| |
|----|
| 10 |
| 20 |
| 30 |
| 40 |
| 20 |

| value | code |
|-----------|------|
| Andrea | 10 |
| Prashanth | 20 |
| Andy | 30 |
| Matt | 40 |

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.

*cause hash fun might have to change
for diff data*

ORDER-PRESERVING ENCODING

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 |



Those
with And at start.

Stay at
first

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

Compressed Data

| name |
|------|
| 10 |
| 40 |
| 20 |
| 30 |
| 40 |

| value | code |
|-----------|------|
| Andrea | 10 |
| Andy | 20 |
| Matt | 30 |
| Prashanth | 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 |
|------|
| 10 |
| 40 |
| 20 |
| 30 |
| 40 |

| value | code |
|-----------|------|
| Andrea | 10 |
| Andy | 20 |
| Matt | 30 |
| Prashanth | 40 |



Sorted Dictionary

CONCLUSION

It is important to choose the right storage model for the target workload:

- OLTP = Row Store →
- OLAP = Column Store →

DBMS^s 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 move data back-and-forth from disk.

← Next