

### Carnegie Mellon University ADVANCED DATABASE SYSTEMS

# Database Compression

Andy Pavlo // 15-721 // Spring 2023

### LAST CLASS

We discussed methods for the DBMS to skip data for sequential scans in OLAP queries via filters (zone maps) and indexes.

### REAL-WORLD DATA CHARACTERISTICS

Data sets tend to have highly **<u>skewed</u>** distributions for attribute values.

 $\rightarrow$  Example: Zipfian distribution of the <u>Brown Corpus</u>

Data sets tend to have high <u>correlation</u> between attributes of the same tuple.

 $\rightarrow$  Example: Zip Code to City, Order Date to Ship Date

# OBSERVATION

I/O is (traditionally) the main bottleneck during query execution. If the DBMS still needs to read data, we need to ensure that it maximizes the amount of useful information it can extract from it.

Key trade-off is <u>speed</u> vs. <u>compression ratio</u>
 → Compressing the data reduces DRAM requirements and processing.

### **TODAY'S AGENDA**

Background Naïve Page Compression Native Columnar Compression Intermediate Data

### DATABASE COMPRESSION

Reduce the size of the database physical representation to increase the # of values accessed and processed per unit of computation or I/O.

**Goal #1:** Must produce fixed-length values. **Goal #2:** Must be a <u>lossless</u> scheme.

**Goal #3:** Ideally postpone decompression for as long as possible during query execution.



### LOSSLESS VS. 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.

Reading less than the entire data set during query execution is sort of like of compression.  $\rightarrow$  Approximate Query Processing



### DATABASE COMPRESSION

If we want to add compression to our DBMS, the first question we must ask ourselves is what is what do want to compress.

This determines what compression schemes are available to us...



### **COMPRESSION GRANULARITY**

#### Choice #1: Block-level

→ Compress a block (e.g., database page, RowGroup) of tuples in a table.

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

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

#### Choice #3: Attribute-level

- $\rightarrow$  Compress a single attribute value within one tuple.
- $\rightarrow$  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.

 $\rightarrow \underline{\text{LZO}} (1996), \underline{\text{LZ4}} (2011), \underline{\text{Snappy}} (2011), \underline{\text{Brotli}} (2013), \\ \underline{\text{Oracle OZIP}} (2014), \underline{\text{Zstd}} (2015)$ 

#### Considerations

- $\rightarrow$  Computational overhead
- $\rightarrow$  Compress vs. decompress speed

# NAÏVE COMPRESSION

Compress data using a general-purpose algorithm. Scope of compression is only based on the data provided as input.

 $\rightarrow \underline{LZO} (1996), \underline{LZ4} (2011), \underline{Snappy} (2011), \underline{Brotli} (2013), \\ \underline{Oracle OZIP} (2014), \underline{Zstd} (2015)$ 

#### Considerations

- $\rightarrow$  Computational overhead
- $\rightarrow$  Compress vs. decompress speed

### MYSQL INNODB COMPRESSION



Source: MySQL 5.7 Documentation

# NAÏVE COMPRESSION

The DBMS must decompress data first before it can be read and (potentially) modified.

- $\rightarrow$  Even if the algo uses dictionary compression, the DBMS cannot access the dictionary's contents.
- $\rightarrow$  This limits the practical scope of the compression scheme.

These schemes also do not consider the high-level meaning or semantics of the data.

### COLUMNAR COMPRESSION

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

Compress runs of the same value in a single column into triplets:

- $\rightarrow$  The value of the attribute.
- $\rightarrow$  The start position in the column segment.
- $\rightarrow$  The # of elements in the run.

Requires the columns to be sorted intelligently to maximize compression opportunities.

Sometimes also called *null suppression* if the DBMS only tracks empty space.



#### **Original Data**

id	lit
1	Y
2	Y
3	Y
4	Ν
6	Y
7	Ν
8	Y
9	Y

#### **Compressed Data**



#### **Compressed Data**

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

SELECT lit, COUNT(\*)
 FROM users
 GROUP BY lit





#### **Original Data**

id	lit
1	Y
2	Y
3	Y
4	Ν
6	Y
7	Ν
8	Y
9	Y

#### **Compressed Data**



**ECMU-DB** 15-721 (Spring 2023)

#### Sorted Data

id	lit
1	Y
2	Y
3	Y
6	Y
8	Y
9	Y
4	Ν
7	N

#### **Compressed Data**





### DICTIONARY COMPRESSION

Replace frequent values with smaller fixed-length codes and then maintain a mapping (dictionary) from the codes to the original values  $\rightarrow$  Typically, one code per attribute value.

 $\rightarrow$  Most widely used native compression scheme in DBMSs.

The ideal dictionary scheme supports fast encoding and decoding for both point and range queries.



### DICTIONARY COMPRESSION

When to construct the dictionary? What is the scope of the dictionary? What data structure do we use for the dictionary? What encoding scheme to use for the dictionary?

### DICTIONARY CONSTRUCTION

#### Choice #1: All-At-Once

- $\rightarrow$  Compute the dictionary for all the tuples at a given point of time.
- $\rightarrow$  New tuples must use a separate dictionary, or the all tuples must be recomputed.
- $\rightarrow$  This is easy to do if the file is immutable.

#### Choice #2: Incremental

- $\rightarrow$  Merge new tuples in with an existing dictionary.
- $\rightarrow$  Likely requires re-encoding to existing tuples.

### DICTIONARY SCOPE

#### Choice #1: Block-level

- $\rightarrow$  Only include a subset of tuples within a single table.
- → DBMS must decompress data when combining tuples from different blocks (e.g., hash table for joins).

#### Choice #2: Table-level

- $\rightarrow$  Construct a dictionary for the entire table.
- $\rightarrow$  Better compression ratio, but expensive to update.

#### Choice #3: Multi-Table

- $\rightarrow$  Can be either subset or entire tables.
- $\rightarrow$  Sometimes helps with joins and set operations.

### ENCODING / DECODING

**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.



### ORDER-PRESERVING ENCODING

The encoded values need to support sorting in the same order as original values.



#### **Original Data**







#### **Compressed Data**

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



### ORDER-PRESERVING ENCODING

SELECTnameFROMusersWHEREnameLIKE'And%'



Still must perform seq scan

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



Only need to access dictionary

#### **Original Data**



#### **Compressed Data**

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



### DICTIONARY DATA STRUCTURES

#### Choice #1: Array

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

#### Choice #2: Hash Table

- $\rightarrow$  Fast and compact.
- $\rightarrow$  Unable to support range and prefix queries.

#### Choice #3: B+Tree

- $\rightarrow$  Slower than a hash table and takes more memory.
- $\rightarrow$  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.





15-721 (Spring 2023)

### DICTIONARY: SHARED-LEAVES B+TREE



SECMU.DB 15-721 (Spring 2023)

### EXPOSING DICTIONARY TO DBMS

Parquet / ORC do not provide an API to directly access a file's compression dictionary.

This means the DBMS cannot perform predicate pushdown and operate directly on compressed data before decompressing it.

Google's Artus proprietary format for <u>Procella</u> supports this.



### **COLUMNAR COMPRESSION**

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

### BITMAP ENCODING

Using bitmaps to represent a column can reduce its storage size if the column's cardinality is low.

#### **Original Data**



#### **Compressed Data**

	<b>1</b> i	it
id	Y	Ν
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

Using bitmaps to represent a column can reduce its storage size if the column's cardinality is low.



#### **Compressed Data**



### BITMAP INDEX: COMPRESSION

#### **Approach #1: General Purpose Compression**

- $\rightarrow$  Use standard compression algorithms (e.g., Snappy, zstd).
- $\rightarrow$  Must decompress entire data chunk before DBMS can use it to process a query.

#### Approach #2: Byte-aligned Bitmap Codes

 $\rightarrow$  Structured run-length encoding compression.

#### Approach #3: Roaring Bitmaps

 $\rightarrow$  Modern hybrid of run-length encoding and value lists.

- Divide bitmap into chunks that contain different categories of bytes:
- $\rightarrow$  **Gap Byte:** All the bits are **0**s.
- $\rightarrow$  Tail Byte: Some bits are 1s.

Encode each <u>chunk</u> that consists of some Gap Bytes followed by some Tail Bytes.

- $\rightarrow$  Gap Bytes are compressed with RLE.
- $\rightarrow$  Tail Bytes are stored uncompressed unless it consists of only 1-byte or has only one non-zero bit.



#### Bitmap

00000000	00000000	00010000
00000000	00000000	00000000
00000000	00000000	00000000
00000000	00000000	00000000
00000000	00000000	00000000
00000000	01000000	00100010

#### **Compressed Bitmap**

Bitmap	Gap Bytes	Tail Bytes	
00000000	00000000	00010000	#
00000000	00000000	00000000	
00000000	00000000	00000000	
00000000	00000000	00000000	
00000000	00000000	00000000	
00000000	01000000	00100010	

#### **Compressed Bitmap**

Ì	Bitmap	Gap Bytes	Tail Bytes	
	00000000	00000000	00010000	#1
	00000000	00000000	00000000	
	00000000	00000000	00000000	
	00000000	00000000	00000000	#2
	00000000	00000000	00000000	
	00000000	01000000	00100010	

#### **Compressed Bitmap**

#### Bitmap

00000000	00000000	00010000	#1
00000000	00000000	00000000	
00000000	00000000	000000000	
00000000	00000000	000000000	
00000000	00000000	000000000	
00000000	01000000	00100010	
00000000	01000000	00100010	

#### **Compressed Bitmap**

**#1** (010)(1)(0100)1-3 4 5-7

#### Chunk #1 (Bytes 1-3)

Header Byte:

- $\rightarrow$  Number of Gap Bytes (Bits 1-3)
- $\rightarrow$  Is the tail special? (Bit 4)
- $\rightarrow$  Number of verbatim bytes (if Bit 4=0)
- $\rightarrow$  Index of one-bit in tail byte (if Bit 4=1)

No gap length bytes since gap length < 7 No verbatim bytes since tail is special

#### Bitmap

00000000	00000000	00010000	
00000000	000000000	00000000	
00000000	00000000	00000000	
00000000	00000000	00000000	#]
00000000	00000000	00000000	
00000000	01000000	00100010	

#### **Compressed Bitmap**

- **#1** (010)(1)(0100)
- **#2** (111)(0)(0010)00001101 01000000 00100010

### Chunk #2 (Bytes 4-18)

- Header Byte:
- $\rightarrow$  13 gap bytes, two tail bytes
- $\rightarrow$  # of gaps is > 7, so must use extra byte

One gap length byte gives gap length = 13 Two verbatim bytes for tail.

#### Bitmap

00000000	00000000	000 <b>1</b> 0000
00000000	00000000	00000000
00000000	00000000	00000000
00000000	00000000	00000000
00000000	00000000	00000000
00000000	01000000	00100010

#### **Compressed Bitmap**

**#1** (010)(1)(0100)

**#2(111)**(0)(0010)00001101 0100000 00100010

#### Chunk #2 (Bytes 4-18)

Header Byte:

- $\rightarrow$  13 gap bytes, two tail bytes
- $\rightarrow$  # of gaps is > 7, so must use extra byte

One gap length byte gives gap length = 13 Two verbatim bytes for tail.

#### Bitmap

00000000	00000000	000 <b>1</b> 0000
00000000	00000000	00000000
00000000	00000000	00000000
00000000	00000000	00000000
00000000	00000000	00000000
00000000	01000000	00100010

#### **Compressed Bitmap**

**#1** (010)(1)(0100) **#2** (111)(0)(0010)00001101 01000000 00100010

### Source: Brian Babcock

#### Chunk #2 (Bytes 4-18)

- Header Byte:
- $\rightarrow$  13 gap bytes, two tail bytes
- $\rightarrow$  # of gaps is > 7, so must use extra byte

One gap length byte gives gap length = 13 Two verbatim bytes for tail.

#### Bitmap

00000000	00000000	000 <b>1</b> 0000
00000000	00000000	00000000
00000000	00000000	00000000
00000000	00000000	00000000
00000000	00000000	00000000
00000000	01000000	00100010

#### **Compressed Bitmap**

**#1** (010)(1)(0100)

#2 (111)(0)(0010)00001101 01000000 00100010

Verbatim Tail Bytes

Source: Brian Babcock

#### Chunk #2 (Bytes 4-18)

- Header Byte:
- $\rightarrow$  13 gap bytes, two tail bytes
- $\rightarrow$  # of gaps is > 7, so must use extra byte

One gap length byte gives gap length = 13 Two verbatim bytes for tail.

Original: 18 bytes BBC Compressed: 5 bytes

### OBSERVATION

#### Oracle's BBC is an obsolete format.

- $\rightarrow$  Although it provides good compression, it is slower than recent alternatives due to excessive branching.
- $\rightarrow$  <u>Word-Aligned Hybrid</u> (WAH) encoding is a patented variation on BBC that provides better performance.

#### None of these support random access.

 $\rightarrow$  If you want to check whether a given value is present, you must start from the beginning and decompress the whole thing.

### ROARING BITMAPS

Store 32-bit integers in a compact two-level indexing data structure.

- $\rightarrow$  Dense chunks are stored using bitmaps
- $\rightarrow$  Sparse chunks use packed arrays of 16-bit integers.

Used in Lucene, Hive, Spark, Pinot.



### ROARING BITMAPS



15-721 (Spring 2023)

For each value  $\mathbf{k}$ , assign it to a chunk based on  $\mathbf{k}/2^{16}$ .  $\rightarrow$  Store  $\mathbf{k}$  in the chunk's container.

If # of values in container is less than 4096, store as array. Otherwise, store as Bitmap.

k=1000 1000/2<sup>16</sup>=0 1000%2<sup>16</sup>=1000

### ROARING BITMAPS



15-721 (Spring 2023)

For each value **k**, assign it to a chunk based on  $k/2^{16}$ .  $\rightarrow$  Store k in the chunk's container. If # of values in container is less than 4096, store as array. Otherwise, store as Bitmap. k=1000 k=199658 199658/216=3  $1000/2^{16}=0$  $1000\%2^{16} = 1000$   $199658\%2^{16} = 50$ 

#### 40

### DELTA ENCODING

Recording the difference between values that follow each other in the same column.

- $\rightarrow$  Store base value <u>in-line</u> or in a separate <u>look-up table</u>.
- $\rightarrow$  Combine with RLE to get even better compression ratios.



# BIT PACKING

If the values for an integer attribute is <u>smaller</u> 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.  $\rightarrow$  Like in BitWeaving/Vertical

#### **Original Data**





# BIT PACKING

If the values for an integer attribute is <u>smaller</u> 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.  $\rightarrow$  Like in BitWeaving/Vertical

8 × 32-bits = Original Data 256 bits



8 × 8-bits = 64 bits

#### 43

### MOSTLY ENCODING

A variation of bit packing for when an attribute's values are "mostly" less than the largest size, store them with smaller data type.

 $\rightarrow$  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: <u>Redshift Documentation</u>

15-721 (Spring 2023)

# MOSTLY ENCODING

A variation of bit packing for when an attribute's values are "mostly" less than the largest size, store them with smaller data type.

 $\rightarrow$  The remaining values that cannot be compressed are stored in their raw form.



SECMU.DB 15-721 (Spring 2023)





mostly8	offset	value		
13	3	999999999		
181				
XXX				
92				
81				
120				
231				
172				

 $(8 \times 8$ -bits) + 16-bits + 32-bits = 112 bits

### INTERMEDIATE RESULTS

After the evaluating a predicate on compressed data, the DBMS will decompress it as it moves from the scan operator to the next operator.

→ Example: Execute a hash join on two tables that use different compression schemes.

The DBMS (typically) does not recompress data during query execution. Otherwise, the system needs to embed decompression logic throughout the entire execution engine.

### PARTING THOUGHTS

Dictionary encoding is not always the most effective compression scheme, but it is the most used.

The DBMS can combine different approaches for even better compression.

### **NEXT CLASS**

#### Query Execution!