

Carnegie Mellon University DVANCE ATABAS Database Compression

@Andy_Pavlo // 15-721 // Spring 2020

UPCOMING DATABASE EVENTS

Oracle Tech Talk

 \rightarrow Wednesday Feb 12th @ 4:30pm

 \rightarrow NSH 4305





3

Integer Numbers

Data Type	Size	Size (Not Null)	Synonyms	Min Value	Max Value
BOOL* (see note below)	2 bytes	1 byte	BOOLEAN	-128	127
BIT	9 bytes	8 bytes]	
TINYINT	2 bytes	1 byte		-128	127
SMALLINT	4 bytes	2 bytes		-32768	32767
MEDIUMINT	4 bytes	3 bytes		-8388608	8388607
INT	8 bytes	4 bytes	INTEGER	-2147483648	2147483647
BIGINT	12 bytes	8 bytes		-2 ** 63	(2 ** 63) - 1

Note:

• BOOL and BOOLEAN are synonymous with TINYINT. A value of 0 is considered FALSE, non-zero values are considered TRUE.



TODAY'S AGENDA

Compression Background Naïve Compression OLAP Columnar Compression OLTP Index Compression



OBSERVATION

I/O is the main bottleneck if the DBMS has to fetch data from disk.

In-memory DBMSs are more complicated.

- Key trade-off is **speed** vs. **compression ratio**
- \rightarrow In-memory DBMSs (always?) choose speed.
- → Compressing the database reduces DRAM requirements and processing.

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



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 <u>late materialization</u>.

Goal #3: Must be a lossless scheme.

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...



DATA SKIPPING

Approach #1: Approximate Queries (Lossy)

- \rightarrow Execute queries on a sampled subset of the entire table to produce approximate results.
- → Examples: <u>BlinkDB</u>, <u>SnappyData</u>, <u>XDB</u>, <u>Oracle</u> (2017)

Approach #2: Zone Maps (Loseless)

- \rightarrow Pre-compute columnar aggregations per block that allow the DBMS to check whether queries need to access it.
- \rightarrow Examples: <u>Oracle</u>, Vertica, MemSQL, <u>Netezza</u>



ZONE MAPS

Pre-computed aggregates for blocks of data. DBMS can check the zone map first to decide whether it wants to access the block.



CMU·DB

OBSERVATION

If we want to add compression to our DBMS, the first question we have to 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

 \rightarrow Compress a block of tuples for the same 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

 \rightarrow 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.

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

Considerations

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

MYSQL INNODB COMPRESSION



NAÏVE COMPRESSION

The DBMS must decompress data first before it can be read and (potentially) modified. \rightarrow This limits the "scope" of the compression scheme.

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



OBSERVATION

We can perform exact-match comparisons and natural joins on compressed data if predicates and data are compressed the same way. \rightarrow Range predicates are trickier...



COLUMNAR COMPRESSION

Null Supression Run-length Encoding Bitmap Encoding Delta Encoding **Incremental Encoding** Mostly Encoding **Dictionary Encoding**

NULL SUPPRESSION

Consecutive zeros or blanks in the data are replaced with a description of how many there were and where they existed. \rightarrow Example: Oracle's Byte-Aligned Bitmap Codes (BBC)

Useful in wide tables with sparse data.



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.







id sex (M,0,3) 1 (F,3,1) 2 SELECT sex, COUNT(*) (M,4,1) 3 FROM users **GROUP BY** sex (F,5,1) 4 (M,6,2) 6 7 **RLE Triplet** - Value 8 - Offset 9

Compressed Data

- Length











Store a separate bitmap for each unique value for an attribute where an offset in the vector corresponds to a tuple.

- \rightarrow The ith position in the Bitmap corresponds to the ith tuple in the table.
- \rightarrow Typically segmented into chunks to avoid allocating large blocks of contiguous memory.

Only practical if the value cardinality is low.













BITMAP ENCODING: EXAMPLE

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

Assume we have 10 million tuples. 43,000 zip codes in the US. → 1000000 × 32-bits = 40 MB → 1000000 × 43000 = 53.75 GB

Every time a txn inserts a new tuple, the DBMS must extend 43,000 different bitmaps.

BITMAP ENCODING: COMPRESSION

Approach #1: General Purpose Compression

- \rightarrow Use standard compression algorithms (e.g., LZ4, Snappy).
- \rightarrow Have to decompress before you can use it to process a query. Not useful for in-memory DBMSs.

Approach #2: Byte-aligned Bitmap Codes

 \rightarrow Structured run-length encoding compression.



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 **1**s.

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

Source: Brian Babcock



Bitmap	Gap Bytes	Tail Bytes
0000000	00000000	00010000 #1
0000000	00000000	00000000
0000000	00000000	00000000
0000000	00000000	00000000
0000000	00000000	00000000
0000000	01000000	00100010
0	1	

Compressed Bitmap

Source: Brian Babcock



]	Bitmap	Gap Bytes	Tail Bytes
	0000000	0 0000000	000 1 0000 #1
ſ	0000000	0 0000000	00000000
	0000000	0 0000000	00000000
	0000000	0 0000000) 00000000 #2
	0000000	0 0000000	00000000
	0000000	0 0100000	00100010

Compressed Bitmap

Source: Brian Babcock



Bitmap

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

Compressed Bitmap

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 1 bit in tail byte (if Bit 4=1)

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

Source: Brian Babcock

Bitmap

 00000000
 00000000
 00010000

 00000000
 00000000
 00000000
 00000000

 00000000
 00000000
 00000000
 #2

 00000000
 00000000
 00000000
 #2

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 have to use extra byte

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

Bitmap

Compressed Bitmap

#1 (010)(1)(0100)

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

Source: Brian Babcock

Chunk #2 (Bytes 4-18)

Header Byte:

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

One gap length byte gives gap length = 13 Two verbatim bytes for tail. 27
Bitmap

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

Chunk #2 (Bytes 4-18)

Header Byte:

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

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

Compressed Bitmap

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

Source: Brian Babcock

Bitmap

0000000000000000010010000000100010

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 have to use extra byte

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

Bitmap

0000000000000000010010000000100010

Compressed Bitmap

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

Source: Brian Babcock

CMU·DB

Chunk #2 (Bytes 4-18)

Header Byte:

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

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

Bitmap

000000000000000001000100000000100010

Compressed Bitmap

#1 (010)(1)(0100)

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

Source: Brian Babcock

Verbatim Tail Bytes

Chunk #2 (Bytes 4-18)

Header Byte:

- \rightarrow 13 gap bytes, two tail bytes
- \rightarrow # of gaps is > 7, so have to 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.

→ If you want to check whether a given value is present, you have to start from the beginning and decompress the whole thing.



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.



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 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.



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 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.



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

5 × 32-bits = 160 bits

Compressed Data



32-bits + (4 × 16-bits) = 96 bits

Compressed Data



32-bits + (2 × 16-bits) = 64 bits



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





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



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



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





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





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





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



MOSTLY ENCODING

When values for an attribute are "mostly" less than the largest size, store them as smaller data type.
→ The remaining values that cannot be compressed are stored in their raw form.



Source: Redshift Documentation

DICTIONARY COMPRESSION

Replace frequent patterns with smaller codes. Most pervasive compression scheme in DBMSs.

Need to support fast encoding and decoding. Need to also support 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.

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.
- \rightarrow Potentially lower compression ratio, but can add new tuples more easily.

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.



MULTI-ATTRIBUTE ENCODING

Instead of storing a single value per dictionary entry, store entries that span attributes. \rightarrow I'm not sure any DBMS implements this.



Compressed Data

val1+val2	V
XX	
YY	
XX	
ZZ	
YY	
XX	
ZZ	
YY	

val1	val2	code
А	202	XX
В	101	YY
С	101	ZZ

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.



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



SELECT name FROM users WHERE name LIKE 'And%'



Original Data





Compressed Data

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



SELECT name FROM users
WHERE name LIKE 'And%'



Still must perform seq scan

Original Data





Compressed Data

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

SELECTnameFROMusersWHEREnameLIKE'And%'



Still must perform seq scan

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

???

Original Data





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





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	Prashanth	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.

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.



SHARED-LEAVES B+TREE



OBSERVATION

An OLTP DBMS cannot use the OLAP compression techniques because we need to support fast random tuple access.
→ Compressing & decompressing "hot" tuples on-the-fly would be too slow to do during a txn.

Indexes consume a large portion of the memory for an OLTP database...

OLTP INDEX OVERHEAD





HYBRID INDEXES

Split a single logical index into two physical indexes. Data is migrated from one stage to the next over time.

- → **Dynamic Stage:** New data, fast to update.
- → **Static Stage:** Old data, compressed + read-only.

All updates go to dynamic stage. Reads may need to check both stages.



HYBRID INDEXES



COMPACT B+TREE





COMPACT B+TREE





COMPACT B+TREE




COMPACT B+TREE







PARTING THOUGHTS

Dictionary encoding is probably the most useful compression scheme because it does not require pre-sorting.

The DBMS can combine different approaches for even better compression.

It is important to wait as long as possible during query execution to decompress data.



NEXT CLASS

Logging + Checkpoints!

