15-721 DATABASE SYSTEMS

Lecture #16 – Database Compression

Andy Pavlo // Carnegie Mellon University // Spring 2016

TODAY'S AGENDA

Background Naïve Compression OLAP Columnar Compression OLTP Index Compression

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

In-memory DBMSs are more complicated
→ Compressing the database reduces DRAM requirements and processing.

Key trade-off is <u>speed</u> vs. <u>compression ratio</u> \rightarrow In-memory DBMSs (always?) choose speed.



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.

Goal #2: Allow the DBMS to postpone decompression 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 is has to be performed at the application level.

Some new DBMSs support approximate queries \rightarrow Example: <u>BlinkDB</u>.



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.

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. Requires DSM.

NAÏVE COMPRESSION

Compress data using a general purpose algorithm. Scope of compression is only based on the data provided as input. \rightarrow Examples: LZO (1996), LZ4 (2011), Snappy (2011).

Considerations

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

NAÏVE COMPRESSION

Choice #1: Entropy Encoding

 \rightarrow More common sequences use less bits to encode, less common sequences use more bits to encode.

Choice #2: Dictionary Encoding

→ Build a data structure that maps data segments to an identifier. Replace those segments in the original data with a reference to the segments position in the dictionary data structure.





Source: <u>MySQL 5.7 Documentation</u> CMU 15-721 (Spring 2016)

10





Source: <u>MySQL 5.7 Documentation</u> CMU 15-721 (Spring 2016)

10



10





Source: <u>MySQL 5.7 Documentation</u> CMU 15-721 (Spring 2016)

NAÏVE COMPRESSION

The data has to be decompressed 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 highlevel meaning or semantics of the data.

SELECT	* FROM	users
WHERE	name =	'Trump'

NAME	SALARY
Trump	Huge
Јоу	Small

SELECT	* FROM	users
WHERE	name =	'Trump'

NAME	SALARY	NAME	SALARY
Trump	Huge	XX	AA
Јоу	Small	YY	BB



COLUMNAR COMPRESSION

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

COMPRESSION VS. MSSQL INDEXES

The MSSQL columnar indexes were a second copy of the data (aka fractured mirrors). \rightarrow The original data was still stored as in NSM format.

We are now talking about compressing the primary copy of the data. Many of the same techniques are applicable.



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.



RUN-LENGTH ENCODING

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.



BITMAP ENCODING

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

 \rightarrow Can use the same compression schemes that we talked about for Bitmap indexes.

Only practical if the value cardinality is low.



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

- \rightarrow The base value can be stored in-line or in a separate look-up table.
- \rightarrow Can be combined 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

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

- \rightarrow The base value can be stored in-line or in a separate look-up table.
- \rightarrow Can be combined 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

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

- \rightarrow The base value can be stored in-line or in a separate look-up table.
- \rightarrow Can be combined 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



time	temp
12:00	99.5
+1	-1
+1	+1
+1	+1
+1	-2

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

- \rightarrow The base value can be stored in-line or in a separate look-up table.
- \rightarrow Can be combined 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

Compressed Data

time	temp
12:00	99.5
+1	-1
+1	+1
+1	+1
+1	-2

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

- \rightarrow The base value can be stored in-line or in a separate look-up table.
- \rightarrow Can be combined 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



time	temp
12:00	99.5
+1	-1
+1	+1
+1	+1
+1	-2

Compressed Data



Type of delta encoding whereby common prefixes or suffixes and their lengths are recorded so that they need not be duplicated. This works best with sorted data.




















MOSTLY ENCODING

When the values for an attribute are "mostly" less than the largest size, you can store them as a smaller data type.

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







Source: <u>Redshift Documentation</u> CMU 15-721 (Spring 2016)

20



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 should the scope be of the dictionary? How do we allow for range queries? How do we enable fast encoding/decoding?

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 Only works for a DSM database \rightarrow I'm not sure any DBMS actually implements this.



ENCODING / DECODING

- A dictionary needs to support two operations:
- → Encode: For a given uncompressed value, convert it into its compressed form.
- → **Decode:** For a given compressed value, convert it back into its original form.

No magic hash function will do this for us. We need two data structures to support operations in both directions.

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

Original Data



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



Compressed Data



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

SELECT * FROM users
WHERE name LIKE 'Tru%'

Original Data



Compressed Data

name	value	code
40	Andy	10
20	Јоу	20
10	Truman	30
30	Trump	40



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

SELECT * FROM users
WHERE name LIKE 'Tru%'



Original Data

name	
Trump	
Јоу	
Andy	
Truman	



value	code
Andy	10
Јоу	20
Truman	30
Trump	40

WHERE name BETWEEN 30 AND 40

SELECT * **FROM** users

Compressed Data

name

40

20

10

30

Original Data





Compressed Data

name	value	code
40	Andy	10
20	Јоу	20
10	Truman	30
30	Trump	40

SELECT name FROM users WHERE name LIKE 'Tru%'



Original Data





name	value	code
40	Andy	10
20	Јоу	20
10	Truman	30
30	Trump	40



SELECT name FROM users WHERE name LIKE 'Tru%'



Still have to perform seq scan

Original Data





name	value	code
40	Andy	10
20	Јоу	20
10	Truman	30
30	Trump	40



SELECT name FROM users WHERE name LIKE 'Tru%'



Still have to perform seq scan

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

???

Original Data

name	
Trump	
Јоу	
Andy	
Truman	

Compressed Data

0 0 00 0	ualua	eede
name	value	coae
40	Andy	10
20	Јоу	20
10	Truman	30
30	Trump	40

SELECT name FROM users WHERE name LIKE 'Tru%'



Still have to perform seq scan

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



Only need to access dictionary

Original Data

name	
Trump	
Јоу	
Andy	
Truman	



name	value	code
40	Andy	10
20	Јоу	20
10	Truman	30
30	Trump	40



DICTIONARY IMPLEMENTATIONS

Hash Table:

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

B+Tree:

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































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 Primary Secondary Indexes Indexes Tuples TPC-C 42.5% 24.0% 33.5% Articles 64.8% 22.6% 12.6% Voter 45.1% 54.9% 0%

OLTP INDEX OVERHEAD Secondary Primary Indexes Indexes Tuples **TPC-C** 42.5% 24.0% 57.5% 33.5% Articles 64.8% 22.6% 12.6% 35.2% 54.9% Voter 45.1% 0% 54.9%



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.
- \rightarrow Static Stage: Old data, compressed + read-only.

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

Bloom Filter







Bloom Filter


























CMU 15-721 (Spring 2016)

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

Query Planning & Optimization Working in a large code base

38

CMU 15-721 (Spring 2016)