

# The Language of Compression

Leif Walsh

Two Sigma Investments  
leif.walsh@gmail.com  
@leifwalsh

September 22, 2015

Today's talk is about **compression**

Today's talk is about **compression**:

- ▶ In data storage systems (databases, filesystems)

Today's talk is about **compression**:

- ▶ In data storage systems (databases, filesystems)
- ▶ Using general-purpose (lossless) algorithms

Today's talk is about **compression**:

- ▶ In data storage systems (databases, filesystems)
- ▶ Using general-purpose (lossless) algorithms
- ▶ On disk, not in memory or over the wire

We'll talk about systems like:

- ▶ MySQL (InnoDB, TokuDB)
- ▶ MongoDB (WiredTiger, TokuMX, RocksDB)
- ▶ Cassandra
- ▶ PostgreSQL
- ▶ Vertica
- ▶ zfs, btrfs

## Goal of the Talk

A **framework** for answering:



A **framework** for answering:

- ▶ How do compression algorithms even work?

A **framework** for answering:

- ▶ How do compression algorithms even work?
- ▶ How do storage systems use compression?

A **framework** for answering:

- ▶ How do compression algorithms even work?
- ▶ How do storage systems use compression?
- ▶ How should I evaluate the compression of a storage system?

A **framework** for answering:

- ▶ How do compression algorithms even work?
- ▶ How do storage systems use compression?
- ▶ How should I evaluate the compression of a storage system?
- ▶ How should I read articles about compression?

A **framework** for answering:

- ▶ How do compression algorithms even work?
- ▶ How do storage systems use compression?
- ▶ How should I evaluate the compression of a storage system?
- ▶ How should I read articles about compression?
- ▶ How should I write articles about compression?

## About Me

## Engineer at Two Sigma

- ▶ We have a lot of data
- ▶ We care a lot about compression

## Engineer at Two Sigma

- ▶ We have a lot of data
- ▶ We care a lot about compression

## Previously at Tokutek

- ▶ Worked on TokuMX, TokuFT
- ▶ We thought a lot about compression
- ▶ We evaluated a lot of compression algorithms
- ▶ We wrote a lot about compression



## How to Talk About Compression

## Scenario

I have a database which can store 1TB of “user data” in only 200GB of disk.

## Scenario

I have a database which can store 1TB of “user data” in only 200GB of disk.

How much is my database compressing?

## Scenario

I have a database which can store 1TB of “user data” in only 200GB of disk.

How much is my database compressing? 80%? 20%? 5x? 1/5? 5:1?

## Scenario

I have a database which can store 1TB of “user data” in only 200GB of disk.

How much is my database compressing? 80%? 20%? 5x? 1/5? 5:1?

Let's talk about what this number is going to mean to us...

## Why Compress?

# Why Compress?

Data storage is expensive (you've heard this)

# Why Compress?

Data storage is expensive (you've heard this)

- ▶ Replication magnifies your data costs



Data storage is expensive (you've heard this)

- ▶ Replication magnifies your data costs
- ▶ Maintenance/operations cost scales superlinearly with hardware

Data storage is expensive (you've heard this)

- ▶ Replication magnifies your data costs
- ▶ Maintenance/operations cost scales superlinearly with hardware
- ▶ SSD is expensive

# Why Compress?

Compression **magnifies** your capacity to store data at a fixed cost.

Compression **minimizes** your cost to provide a fixed capacity.

# Why Compress?

Compression **magnifies** your capacity to store data at a fixed cost.

We ask “by what **factor** does compression *multiply* my *capacity*?”

Compression **minimizes** your cost to provide a fixed capacity.

We ask “by what **factor** does compression *divide* my *cost*?”

Compression **magnifies** your capacity to store data at a fixed cost.

We ask “by what **factor** does compression *multiply* my *capacity*?”

Compression **minimizes** your cost to provide a fixed capacity.

We ask “by what **factor** does compression *divide* my *cost*?”

We should always talk about compression in terms of the **multiplicative factor** by which you increase your **cost-effectiveness**.

Say **“5x compression”**, not “80% compression”.

# Cost Model

Compression is **more expensive** than decompression.



Compression is **more expensive** than decompression.

- ▶ Compression is **searching** for repeated patterns in data. Searching is expensive.

Compression is **more expensive** than decompression.

- ▶ Compression is **searching** for repeated patterns in data. Searching is expensive.
- ▶ Decompression is **copying** bytes out in the order described by encoding, which isn't very hard.

Compression is **more expensive** than decompression.

- ▶ Compression is **searching** for repeated patterns in data. Searching is expensive.
- ▶ Decompression is **copying** bytes out in the order described by encoding, which isn't very hard.

Bandwidth speeds for typical compression algorithms (cp is a no-op)  
(my laptop, Haswell CPU, Samsung SSD, 362MB tarball of /usr/include):

(MB/s)	zlib	bz2	lzma	lzo	lz4	zstd	cp
Compress	39	8	3	366	405	293	1466
Decompress	179	28	138	395	774	500	1466

(higher is better)

How does compression impact **perceived** performance?

How does compression impact **perceived** performance?

Compression:

- ▶ Usually infrequent and done in the background
- ▶ Can reduce overall **throughput**

How does compression impact **perceived** performance?

Compression:

- ▶ Usually infrequent and done in the background
- ▶ Can reduce overall **throughput**

Decompression:

- ▶ More frequent (“Write Once, Read Many”) and on the critical path
- ▶ High impact on user-visible **latency**



1. Do compression in the background and in large batches
  - ▶ Implement backpressure to avoid falling behind
  - ▶ If backpressure reaches users, try a faster compressor



1. Do compression in the background and in large batches
  - ▶ Implement backpressure to avoid falling behind
  - ▶ If backpressure reaches users, try a faster compressor
2. Be sensitive to decompression latency
  - ▶ Hit the highest nail: other latency sources may be more important
  - ▶ Experiment with block sizes and faster compression algorithms

## How Compression Even Works

- All\* compression algorithms, at their core, use a form of *dictionary encoding*:
- ▶ Write down a dictionary of “common phrases” with shorter names
  - ▶ Encode the input stream by referencing the short names in the dictionary

- All\* compression algorithms, at their core, use a form of *dictionary encoding*:
- ▶ Write down a dictionary of “common phrases” with shorter names
  - ▶ Encode the input stream by referencing the short names in the dictionary

\*A *Universal Algorithm for Sequential Data Compression*, J. Ziv, A. Lempel 1977

All\* compression algorithms, at their core, use a form of *dictionary encoding*:

- ▶ Write down a dictionary of “common phrases” with shorter names
- ▶ Encode the input stream by referencing the short names in the dictionary

abbabbabbcdcdcdcdabb => abb|abb|abb|cd|cd|cd|cd|abb

\*A *Universal Algorithm for Sequential Data Compression*, J. Ziv, A. Lempel 1977

All\* compression algorithms, at their core, use a form of *dictionary encoding*:

- ▶ Write down a dictionary of “common phrases” with shorter names
- ▶ Encode the input stream by referencing the short names in the dictionary

abbabbabbcdcdcdcdabb => abb|abb|abb|cd|cd|cd|cd|abb

Symbol	Phrase
x	abb
y	cd

, xxxyyyyyx

\*A Universal Algorithm for Sequential Data Compression, J. Ziv, A. Lempel 1977

All\* compression algorithms, at their core, use a form of *dictionary encoding*:

- ▶ Write down a dictionary of “common phrases” with shorter names
- ▶ Encode the input stream by referencing the short names in the dictionary

abbabbabbcdcdcdcdabb => abb|abb|abb|cd|cd|cd|cd|abb

Symbol	Phrase
x	abb
y	cd

, xxxyyyyyx

To decompress: read the dictionary, use it to interpret the compressed stream.

\*A *Universal Algorithm for Sequential Data Compression*, J. Ziv, A. Lempel 1977

Most compressors have a dynamic dictionary which is modified (optimized) as it compresses the input.



Most compressors have a dynamic dictionary which is modified (optimized) as it compresses the input.

The dictionary takes up some space in the file header, so to be worthwhile, we want to compress a lot of input with it at once.

We cannot seek directly to an offset in the decompressed output because:

We cannot seek directly to an offset in the decompressed output because:

- ▶ We need to read the compressed stream to modify the dictionary

We cannot seek directly to an offset in the decompressed output because:

- ▶ We need to read the compressed stream to modify the dictionary
- ▶ We don't know how much output any given chunk of input will produce

We cannot seek directly to an offset in the decompressed output because:

- ▶ We need to read the compressed stream to modify the dictionary
- ▶ We don't know how much output any given chunk of input will produce

We also can't update a compressed file without recompressing the whole thing.

Systems that provide seeking in compressed data do so by dividing the input into **blocks**, and compressing them individually.

Systems that provide seeking in compressed data do so by dividing the input into **blocks**, and compressing them individually.

- ▶ When writing, recompress the whole block being written (but not the whole data set)

Systems that provide seeking in compressed data do so by dividing the input into **blocks**, and compressing them individually.

- ▶ When writing, recompress the whole block being written (but not the whole data set)
- ▶ When reading, decompress the whole block being read



Systems that provide seeking in compressed data do so by dividing the input into **blocks**, and compressing them individually.

- ▶ When writing, recompress the whole block being written (but not the whole data set)
- ▶ When reading, decompress the whole block being read
- ▶ Overall compression ratio depends on the size of the blocks

## Block Sizes

Compression algorithms are pattern finders. Give them more data to search in, and they find more patterns.

Compression algorithms are pattern finders. Give them more data to search in, and they find more patterns.

Compressors use **block sizes** to limit their runtime and memory usage.

Compression algorithms are pattern finders. Give them more data to search in, and they find more patterns.

Compressors use **block sizes** to limit their runtime and memory usage.

As block size increases:

Compression algorithms are pattern finders. Give them more data to search in, and they find more patterns.

Compressors use **block sizes** to limit their runtime and memory usage.

As block size increases:

- ▶ Compression throughput decreases

Compression algorithms are pattern finders. Give them more data to search in, and they find more patterns.

Compressors use **block sizes** to limit their runtime and memory usage.

As block size increases:

- ▶ Compression throughput decreases
- ▶ Compression and decompression memory usage increases

Compression algorithms are pattern finders. Give them more data to search in, and they find more patterns.

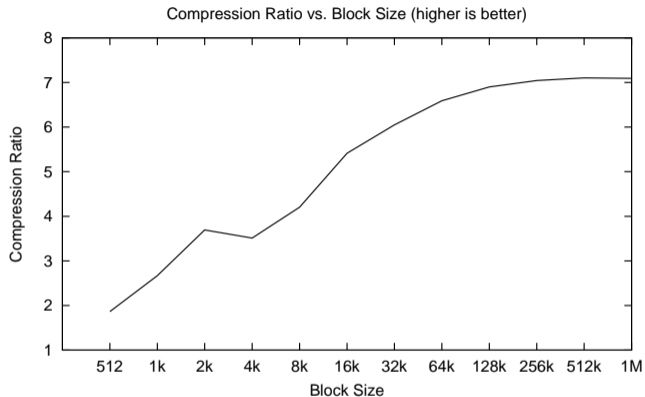
Compressors use **block sizes** to limit their runtime and memory usage.

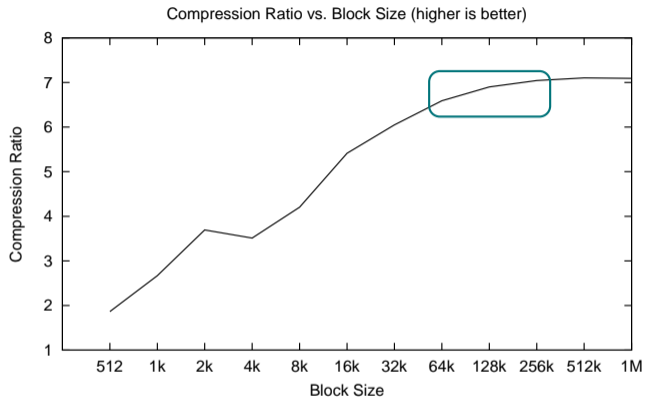
As block size increases:

- ▶ Compression throughput decreases
- ▶ Compression and decompression memory usage increases
- ▶ Decompression throughput may increase if disk throughput increases

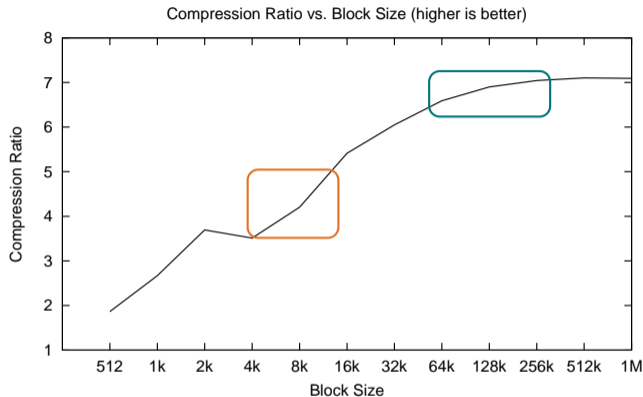


# Block Sizes



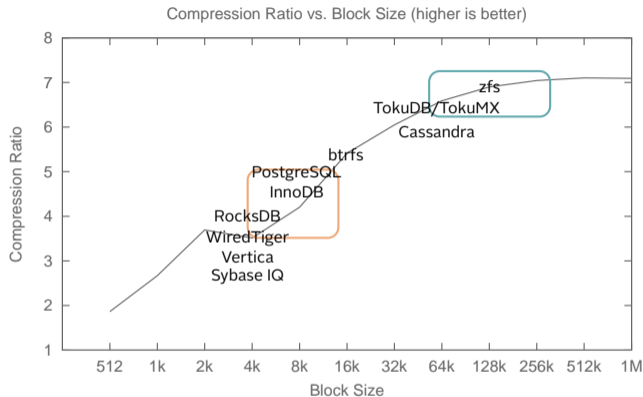


The compression ratio sweet spot is  $\sim$ **128k**, for gzip on this data set.



The compression ratio sweet spot is  $\sim 128k$ , for gzip on this data set.

Most systems use small blocks  $\sim 8k$ , to reduce decompression latency.



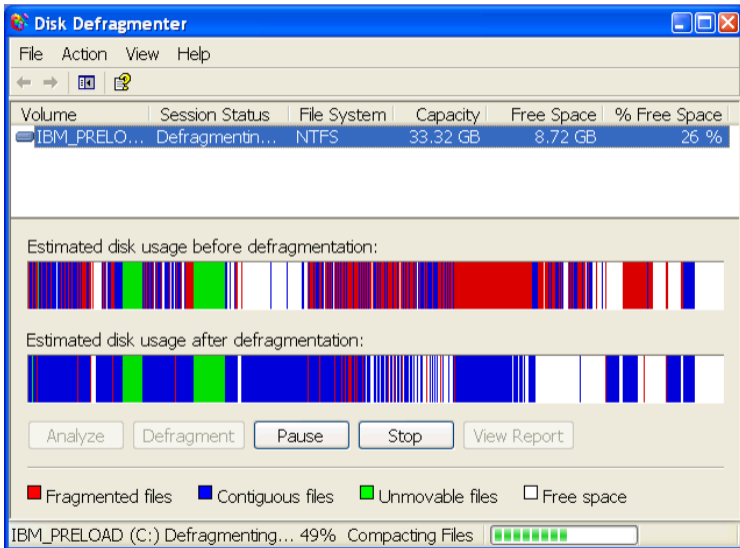
The compression ratio sweet spot is  $\sim 128k$ , for gzip on this data set.

Most systems use small blocks  $\sim 8k$ , to reduce decompression latency.

Another corollary of compressing in blocks is **fragmentation**.

Another corollary of compressing in blocks is **fragmentation**.

Blocks need to be allocated locations on disk. As the data grows, shrinks, and moves around, these locations (and for some systems, allocation sizes) change.



The screenshot shows the Disk Defragmenter utility window. The title bar reads "Disk Defragmenter". The menu bar includes "File", "Action", "View", and "Help". Below the menu is a toolbar with navigation and help icons. A table displays disk information:

Volume	Session Status	File System	Capacity	Free Space	% Free Space
IBM_PRELO...	Defragmentin...	NTFS	33.32 GB	8.72 GB	26 %

Below the table, two bar charts illustrate disk usage. The first chart, "Estimated disk usage before defragmentation:", shows a highly fragmented disk with many small, scattered blocks of data. The second chart, "Estimated disk usage after defragmentation:", shows the same disk with data blocks consolidated into larger, contiguous segments, significantly reducing fragmentation.

At the bottom of the window, there are five buttons: "Analyze", "Defragment", "Pause", "Stop", and "View Report". A legend below the buttons identifies the colors used in the charts: red for "Fragmented files", blue for "Contiguous files", green for "Unmovable files", and white for "Free space".

The status bar at the very bottom indicates the current operation: "IBM\_PRELOAD (C:) Defragmenting... 49% Compacting Files" followed by a progress bar with 10 green segments.

Fragmentation hurts you in two ways:



Fragmentation hurts you in two ways:

1. Fragmented files occupy more **effective space** than defragmented ones

Fragmentation hurts you in two ways:

1. Fragmented files occupy more **effective space** than defragmented ones
2. Fragmentation degrades range query throughput by reducing data locality

Fragmentation hurts you in two ways:

1. Fragmented files occupy more **effective space** than defragmented ones
2. Fragmentation degrades range query throughput by reducing data locality

For some systems, the overall compression ratio will be reduced once fragmentation develops.

# Entropy

Not all data compresses equally!

Not all data compresses equally!

*Information Theory*\* can tell us how much **real information** is present in a set of data (“bits of entropy”).

Not all data compresses equally!

*Information Theory*\* can tell us how much **real information** is present in a set of data (“bits of entropy”).

\**A Mathematical Theory of Communication*, C. E. Shannon, 1948

Not all data compresses equally!

*Information Theory*\* can tell us how much **real information** is present in a set of data (“bits of entropy”).

A general-purpose, lossless compression algorithm can't hope to compress data smaller than that.

\**A Mathematical Theory of Communication*, C. E. Shannon, 1948



Not all data compresses equally!

*Information Theory*\* can tell us how much **real information** is present in a set of data (“bits of entropy”).

A general-purpose, lossless compression algorithm can't hope to compress data smaller than that.

If it could, it would have to produce the same compressed output for multiple inputs, which would mean it isn't *lossless*.

\**A Mathematical Theory of Communication*, C. E. Shannon, 1948

Not all data compresses equally!

*Information Theory*\* can tell us how much **real information** is present in a set of data (“bits of entropy”).

A general-purpose, lossless compression algorithm can't hope to compress data smaller than that.

If it could, it would have to produce the same compressed output for multiple inputs, which would mean it isn't *lossless*.

*High entropy data is highly uncompressible.*

*Low entropy data is easily compressed.*

\**A Mathematical Theory of Communication*, C. E. Shannon, 1948

# Entropy: Experiment

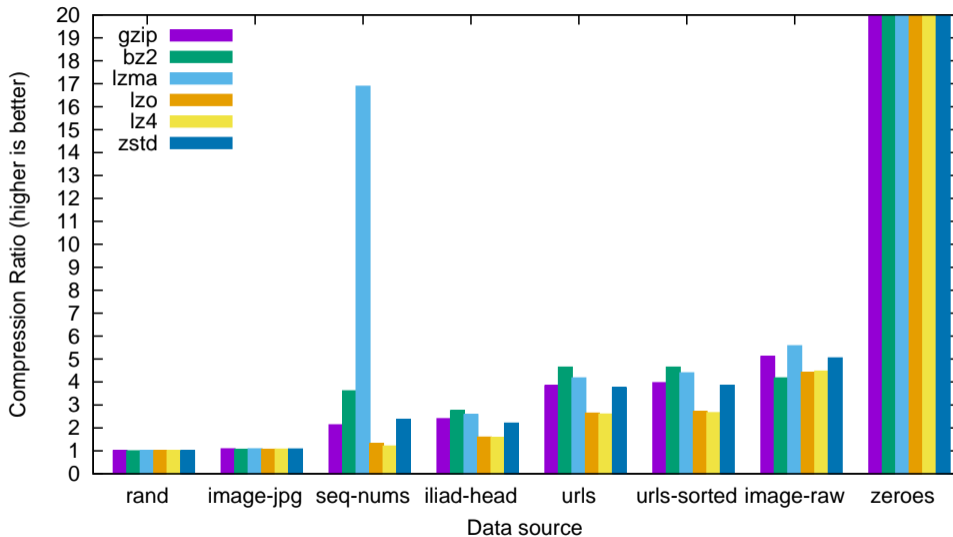
# Entropy: Experiment

Built 8 data sources ( $\sim 50k$  each):

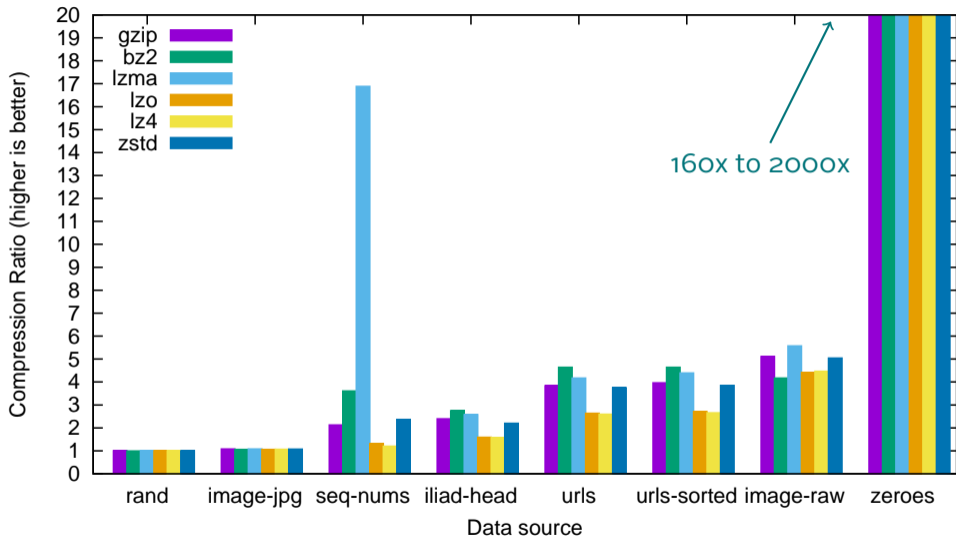
Built 8 data sources (~50k each):

1. Random bytes
2. Sequential numbers, encoded as ASCII decimals
3. All zeroes
4. The beginning of *The Iliad*
5. 1000 random Wikipedia URLs
6. 1000 random Wikipedia URLs, sorted
7. RAW image (CR2)
8. JPEG-compressed image

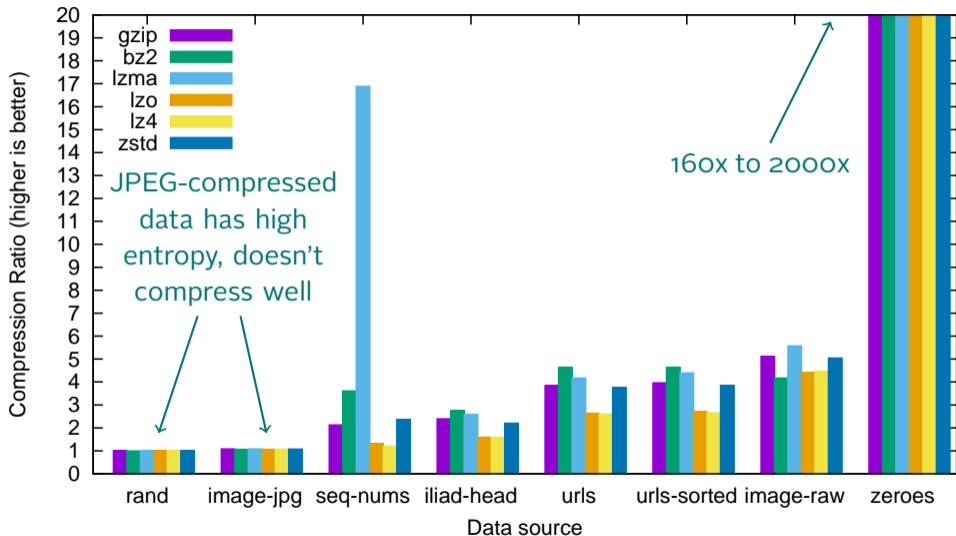
# Entropy: Experiment



# Entropy: Experiment



# Entropy: Experiment





*Homogeneous* data has lower entropy than *heterogeneous* data.

*Homogeneous* data has lower entropy than *heterogeneous* data.

- ▶ Integers compress better than documents with complex internal structure

*Homogeneous* data has lower entropy than *heterogeneous* data.

- ▶ Integers compress better than documents with complex internal structure

**Column stores** have a compression advantage over **row stores**.

Know your data!  
Don't waste your time compressing JPEG blobs

Know your data!

Don't waste your time compressing JPEG blobs

Some compressors are fantastic in specific data domains  
(VLQ, delta coding, JPEG, MP3, ...)

Know your data!

Don't waste your time compressing JPEG blobs

Some compressors are fantastic in specific data domains  
(VLQ, delta coding, JPEG, MP3, ...)

(But 95% of the time, gzip is fine)

Before we use compression, we need to understand the **costs** and **benefits** to our application.

# Benchmarking



When designing a compression benchmark, you should consider:

When designing a compression benchmark, you should consider:

- ▶ Execution

When designing a compression benchmark, you should consider:

- ▶ Execution
- ▶ Measurement

When designing a compression benchmark, you should consider:

- ▶ Execution
- ▶ Measurement
- ▶ Presentation

## Main Question

Is the workload representative of a real-world use-case?

1. Sample real data if you can get it.

1. Sample real data if you can get it.

If not, generate *plausibly realistic* data:

1. Sample real data if you can get it.

If not, generate *plausibly realistic* data:

- ▶ Zeroes: bad
- ▶ Random: bad
- ▶ 25% random and 75% zeroes: meh
- ▶ JSON blobs: good



2. Use a realistic read/insert/update mixture:

2. Use a realistic read/insert/update mixture:
  - ▶ Most applications are read-heavy
  - ▶ Favors fast decompressors

3. Use a realistic insert/update distribution:

3. Use a realistic insert/update distribution:
  - ▶ Most applications don't write uniformly over the keyspace

3. Use a realistic insert/update distribution:
  - ▶ Most applications don't write uniformly over the keyspace
  - ▶ Zipfian or Pareto (or sometimes sequential, or nearly) distributions are more realistic, and cache-friendlier

### 3. Use a realistic insert/update distribution:

- ▶ Most applications don't write uniformly over the keyspace
- ▶ Zipfian or Pareto (or sometimes sequential, or nearly) distributions are more realistic, and cache-friendlier
- ▶ Vadim wrote a `sysbench` workload generator that uses a Zipfian distribution:  
<http://j.mp/sysbench-zipf>

4. To measure **latency**, throttle your workload.

4. To measure **latency**, throttle your workload.
  - ▶ Full-throughput workloads will induce artificial latency spikes (fsyncs, GC)



4. To measure **latency**, throttle your workload.
  - ▶ Full-throughput workloads will induce artificial latency spikes (fsyncs, GC)

To measure **max throughput**, run at full speed.

4. To measure **latency**, throttle your workload.
  - ▶ Full-throughput workloads will induce artificial latency spikes (fsyncs, GC)

To measure **max throughput**, run at full speed.

You should do both.

5. Run for a *long time*. Lots of important properties don't become visible immediately (e.g. fragmentation), and you need to understand them.

5. Run for a *long time*. Lots of important properties don't become visible immediately (e.g. fragmentation), and you need to understand them.

Your application is hopefully going to run for months or years. You don't want to be surprised by degradation after you think everything's stable.

6. Parameterize your workload:

6. Parameterize your workload:
  - ▶ Read/insert/update mixture

6. Parameterize your workload:
  - ▶ Read/insert/update mixture
  - ▶ Write distribution

6. Parameterize your workload:
  - ▶ Read/insert/update mixture
  - ▶ Write distribution
  - ▶ Number of threads



6. Parameterize your workload:
  - ▶ Read/insert/update mixture
  - ▶ Write distribution
  - ▶ Number of threads
  - ▶ Data size

6. Parameterize your workload:
  - ▶ Read/insert/update mixture
  - ▶ Write distribution
  - ▶ Number of threads
  - ▶ Data size
  - ▶ Throttling

6. Parameterize your workload:
  - ▶ Read/insert/update mixture
  - ▶ Write distribution
  - ▶ Number of threads
  - ▶ Data size
  - ▶ Throttling
  - ▶ Duration

## 6. Parameterize your workload:

- ▶ Read/insert/update mixture
- ▶ Write distribution
- ▶ Number of threads
- ▶ Data size
- ▶ Throttling
- ▶ Duration
- ▶ System configuration (cache size, isolation levels, log commit)

## 6. Parameterize your workload:

- ▶ Read/insert/update mixture
- ▶ Write distribution
- ▶ Number of threads
- ▶ Data size
- ▶ Throttling
- ▶ Duration
- ▶ System configuration (cache size, isolation levels, log commit)

You are going to want to explore these parameter spaces. Save yourself the pain later and think about parameterization up front.

Great example: <https://github.com/ParsePlatform/flashback>

Captures a MongoDB workload with profiling, then replays operations either at their original timestamps, or at full speed.







Application metrics:

## Application metrics:

- ▶ Throughput
- ▶ Latency
- ▶ Aborted/retried transactions

Application metrics:

- ▶ Throughput
- ▶ Latency
- ▶ Aborted/retried transactions

Instrument your application so you know *which* operations are expensive.

System metrics:

## System metrics:

- ▶ CPU
- ▶ Memory (RSS)
- ▶ I/O
- ▶ Network
- ▶ Actual storage usage (du)

## System metrics:

- ▶ CPU
- ▶ Memory (RSS)
- ▶ I/O
- ▶ Network
- ▶ Actual storage usage (du)

`perf(1)`, `iostat(1)`, `dstat(1)`, `oprofile(1)`, `collectd(1)`, Datadog, Librato, ...

Database/filesystem metrics (product-specific):

Database/filesystem metrics (product-specific):

- ▶ Cache hits/misses
- ▶ Replication lag
- ▶ Checkpoint lag



Database/filesystem metrics (product-specific):

- ▶ Cache hits/misses
- ▶ Replication lag
- ▶ Checkpoint lag

Talk to your storage vendor about what's important.



1. Describe the workload, and make a case for why it's **realistic**

1. Describe the workload, and make a case for why it's **realistic**
2. Choose **key metrics** that reflect the *benefits* of compression (e.g. users stored per TB) as well as the *costs* (e.g. operation latency)

1. Describe the workload, and make a case for why it's **realistic**
2. Choose **key metrics** that reflect the *benefits* of compression (e.g. users stored per TB) as well as the *costs* (e.g. operation latency)
3. Demonstrate which **parameter choices** influence the costs and benefits you think are important

1. Describe the workload, and make a case for why it's **realistic**
2. Choose **key metrics** that reflect the *benefits* of compression (e.g. users stored per TB) as well as the *costs* (e.g. operation latency)
3. Demonstrate which **parameter choices** influence the costs and benefits you think are important
4. **Explain** which parameters have little or no effect on your metrics
5. Explain how much of your measurement is **overhead**.
6. *If* you show charts, **normalize** your data. Only present important differences.

# Review

Say **“5x compression”**



Compression is **slower** than decompression, but  
decompression is **more frequent**

**Large blocks** compress better

Fragmentation degrades  
**effective compression** over time

High entropy data is **less compressible**

Benchmark **realistic workloads** over a long period

Present **responsibly**  
(and distrust benchmarkers who don't)

- ▶ Tim and Mark Callaghan for being exemplar benchmarkers  
(<http://acmebenchmarking.com> and <http://smalldatum.blogspot.com>)
- ▶ Bohu Tang for introducing me to zstd
- ▶ Andrew Bolin, Corey Milloy, Effie Baram, Li Jin, Wil Yegelwel for making this talk better
- ▶ Tokutek engineering
- ▶ Percona (they're also good benchmarkers)

Questions?

Leif Walsh  
leif.walsh@gmail.com  
@leifwalsh



This document is being distributed for informational and educational purposes only and is not an offer to sell or the solicitation of an offer to buy any securities or other instruments. The information contained herein is not intended to provide, and should not be relied upon for investment advice. The views expressed herein are not necessarily the views of Two Sigma Investments, LP or any of its affiliates (collectively, "Two Sigma"). Such views reflect significant assumptions and subjective of the author(s) of the document and are subject to change without notice. The document may employ data derived from third-party sources. No representation is made as to the accuracy of such information and the use of such information in no way implies an endorsement of the source of such information or its validity.

The copyrights and/or trademarks in some of the images, logos or other material used herein may be owned by entities other than Two Sigma. If so, such copyrights and/or trademarks are most likely owned by the entity that created the material and are used purely for identification and comment as fair use under international copyright and/or trademark laws. Use of such image, copyright or trademark does not imply any association with such organization (or endorsement of such organization) by Two Sigma, nor vice versa.