

# Whispering in the trees

Scaling go-carbon, the go-graphite storage stack at [Booking.com](https://www.booking.com)

Xiaofan Hu @ [Booking.com](https://www.booking.com)

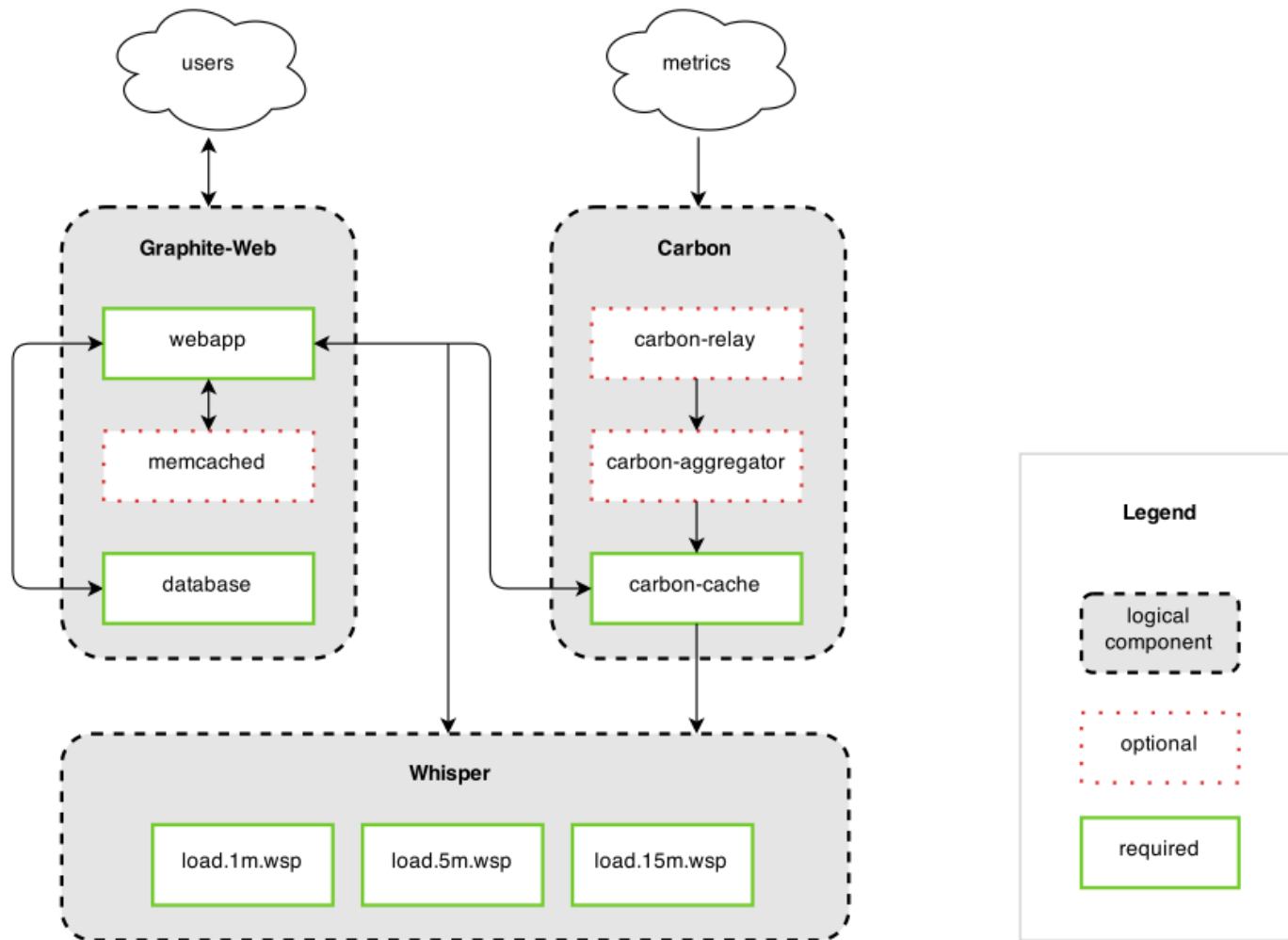
(note: puns intended in the title)

# Disclaimer

# What is Graphite

Graphite is a time-series database. It was originally written in python (mainly), the whole tool consists of multiple components like:

- frontend carbon api for returning timeseries data or graphite-web for rendering graph
- relay (for scaling and duplicating data)
- storage: carbon and whisper
- admin tooling: buckytools



credit: <https://github.com/graphite-project/whisper>

# Graphite at Booking

No longer a vanilla setup, various components are rewritten (some more than once), for example:

- [carbonapi](#), rewritten by [Damian Gryski](#), [Vladimir Smirnov](#) and many others.
- relay is now [nanotube](#) written by Roman Grytskiv, Gyanendra Singh and Andrei Vereha from our Graphite team, (it was preceded by [carbon-c-relay](#) written by [Fabian Groffen](#))
- [go-carbon](#) for storage, written by [Roman Lomonosov](#)

My story today is mainly about the storage program: go-carbon.

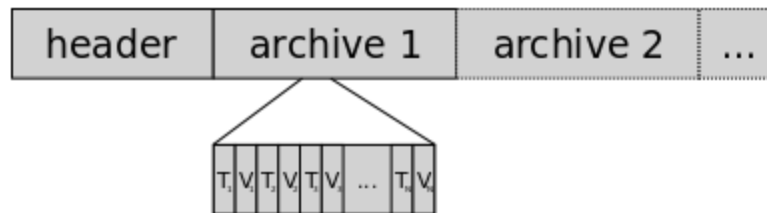
# Graphite Metric

- An example graphite metric: `sys.cpu.loadavg.app.host-0001`
- An example graphite retention policy: `1s:2d,1m:30d,1h:2y`
  - size of the retention example:  $(864002 + 144030 + 24 \times 730) \times 12 = 2,802,240$  bytes
  - 1s:2d is called an archive (same for 1m:30d and 1h:2y)

# What is Whisper

In graphite, each metric is saved in a file, using the a round-robin database format, named whisper. Important properties:

- Data point addressable: given a random timestamp and a target archive, its location could be inferred in the whisper file, which means that it is programmably trivial to support out-of-order data and rewrite
- Fixed size: each data point has a fixed size of 12 bytes (4 bytes for timestamp, 8 bytes for value and yes, one more thing to fix before 2038)



# What is Gorilla compression

- An compression algorithm published in VLDB '15: Gorilla: Facebook's Fast, Scalable, In-Memory Time Series Database
- It has great compression performace for time series data (even though it's payload dependent)
- It has seen wide adoption since then: [M3DB](#), [Prometheus](#), [Timescale](#), [VictoriaMetrics](#), etc.



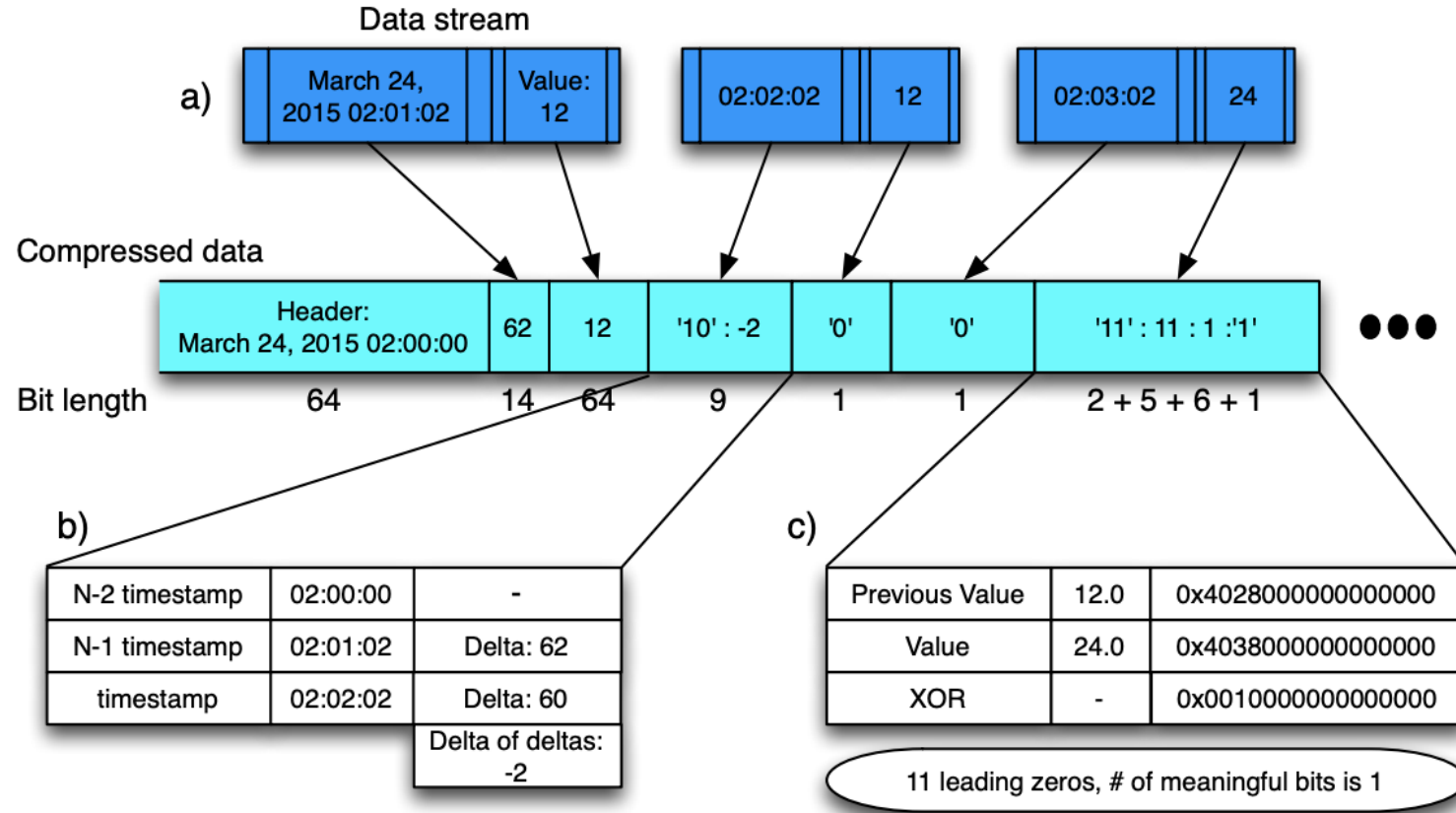
# The core of the Gorilla algorithm

- Delta encoding for timestamps

To be precise, it's actually the delta of delta

- XOR for values

Built on the assumption that timeseries data tend to have constant/repetitive values, or values fluctuating within a certain range, this means that XOR with the previous value often has leading and trailing zeros, and we can only save mostly just the meaningful bits



**Figure 2: Visualizing the entire compression algorithm. For this example, 48 bytes of values and time stamps are compressed to just under 21 bytes/167 bits.**

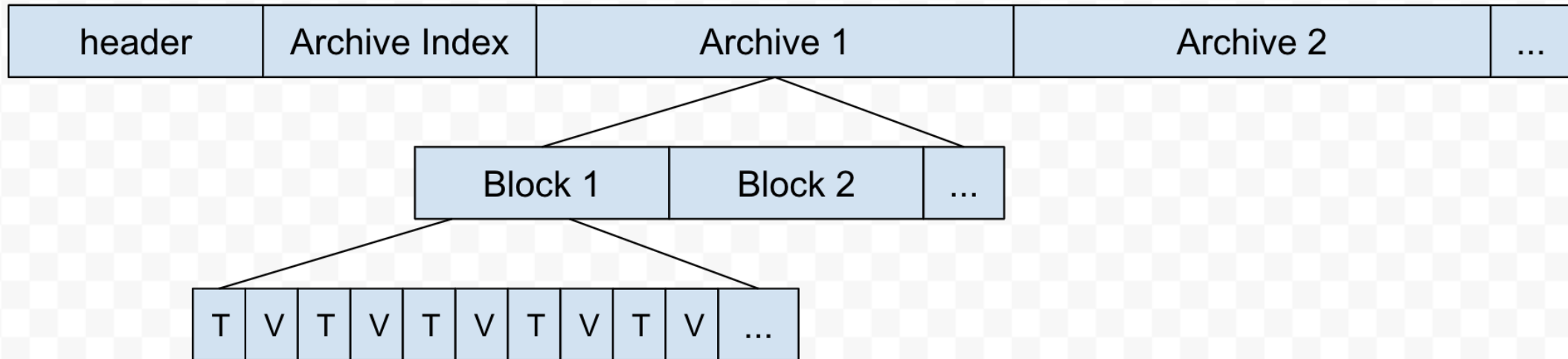
## Best case example

#	Timestamp	Value
#1	1598475390	0
#2	1598475391	0
#3	1598475392	0
...	...	0
#100	1598475493	0

With the compression algorithms introduced in the gorilla paper, orther than the first two data points, the rest of them could be compressed with 2 bits.

# How to combine Gorilla and Whisper

A new file format needs to be designed from scratch in order to compress data points using the gorilla algorithm.



# CWhisper (Compressed Whisper)

- Still a round robin database
- File size isn't fixed (would grow/extend over time)
- Archives are split into many blocks (ideally consist of 7200 data points per archive)
- No longer data point addressable (means hard to support rewrite and limited out-of-order range)

# Result

Metrics	Whisper (standard)	CWhisper (compressed)
Total Metrics	50.6 Millions	53.1 Millions
Num of Servers	32	9
Disk Usage ( <b>45.75% less</b> )	32.28 TB	14.77 TB
Total Disk Space (2.9TB Per Server)	92.8 TB	26.1 TB
Theoretical Capacity Per Server (Metrics)	~4.5 Millions	~10.43 Millions

# Globbing graphite metrics

A most simple graphite query: `sys.cpu.loadavg.app.host-0*`

It's basically the same as globbing in shell: `ls /sys/cpu/loadavg/app/host-0*`

## **filepath.Match/Glob (Go stdlib)**

Pro: simple to implement

Con: high performance cost in a large file tree (millions of files)

filepath.Glob in Go is an userspace implementation, so it first needs to ask the kernel for all the files and then globs over it. Therefore the overhead is a high when serving millions of files.



# Trigram (part 1)

There is alternative implementation in go-carbon, which is using trigram, originally implemented by Damian Gryski.

TLDR: it breaks down all the metrics as trigrams, and maps the trigram to the metrics (an inverted index). A glob query is also converted as a trigrams, then intersects the metric trigrams and query trigrams, then it would use the glob to make sure the files match the query.

# Trigram (part 2)

Pro:

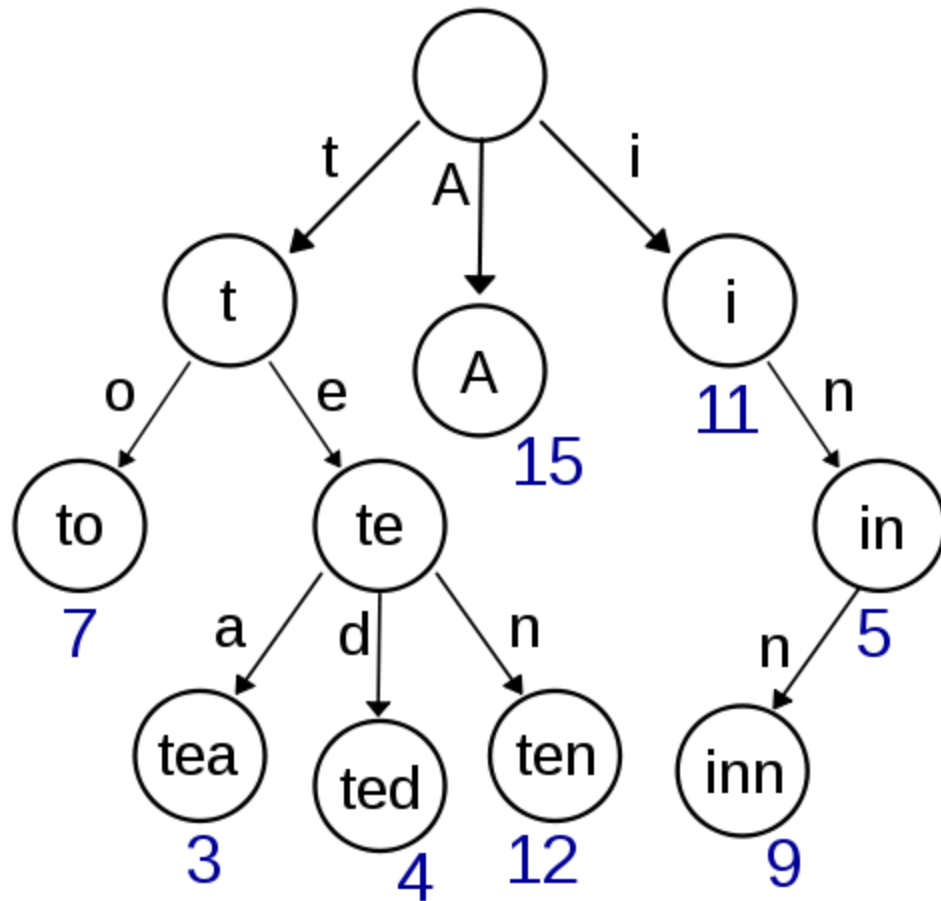
- faster than standard library (no syscalls after index, and file list are cached in memory)

Con:

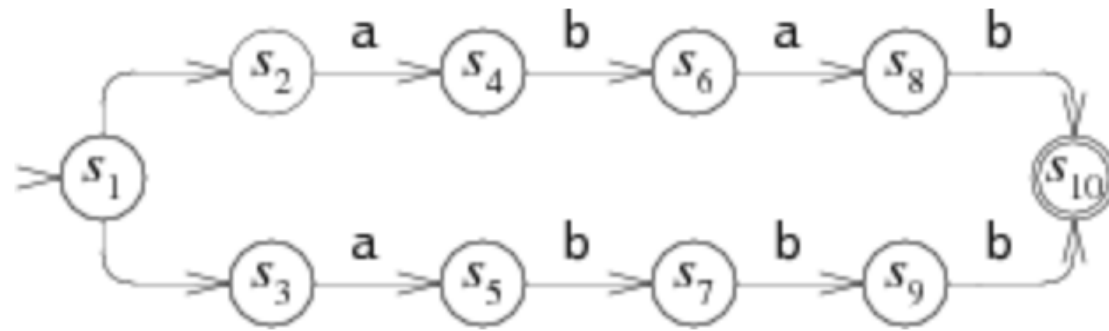
- index is expensive to build when dealing higher number of metrics (above 5 millions or more)
- result returned by trigram index aren't always matching the query, so it still falls back to `filepath.Match` to double check

(trigram itself is a pretty big topic, so sorry that I can't explain all its glory too much)

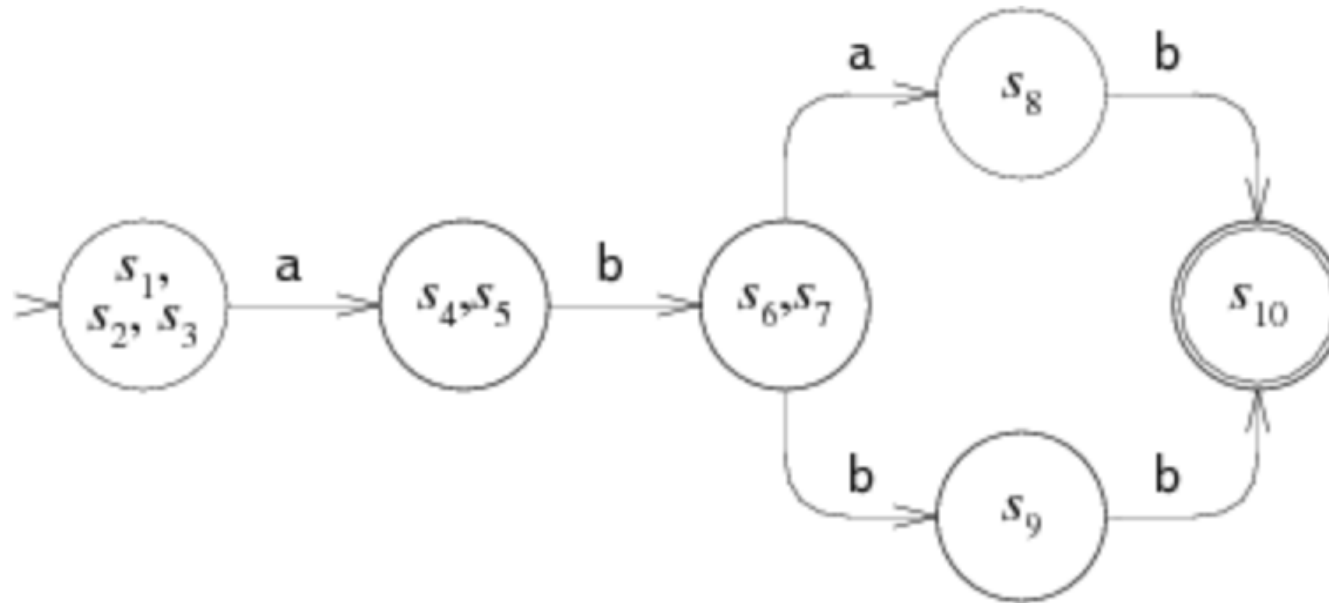
# Trie + NFA/DFA (part 1)



# Trie + NFA/DFA (part 2)



abab | abbb



abab | abbb

## Trie + NFA/DFA (part 3)

TLDR: index all the metrics in go-carbon instance with trie, compile the glob queries first as nfa (then dfa during walking). And walking over the trie and nfa/dfa at the same time.

More details about nfa and dfa could be found in

<https://swtch.com/~rsc/regexp/regexp1.html>

# Trie + NFA/DFA (part 4)

Pro:

- faster index time
- less memory usage
- no standard library fallback
- better/predictable performance

Con:

- Certain types of queries are faster using trigram (like `foo.*bar.zoo`, because of the leading star, the new index algorithm needs to travel the whole namespace, however, arguably, you can design your metric namespace properly to avoid this issue)

# Result

Time Range	Trigram	Trie+DFA
1μs-10μs	1621	0
10μs-100μs	104911	85662
100μs-1ms	20617	74514
1ms-10ms	18214	19454
10ms-100ms	34601	4164
1m40s-16m40s	11	418
100ms-1s	3851	6
1s-10s	219	0
10s+	21	0
Total	184066	184218
Queries Finished in 10ms	78.97%	97.51%

# Tips

More common names should come before less common and unique names in the metric: less memory usage and faster query.

`sys.cpu.loadavg.app.host-0001` performs better than `sys.app.host-0001.cpu.loadavg` using trie index + nfa/dfa.

Because in the first naming pattern, `sys.cpu.loadavg` is just one copy of string in the trie index and comparison is done only once.



# Production and Community usage status

- Challenges on rolling out compressed whisper
  - Out of order
  - Rewrite
- Trie+NFA/DFA index solution made it to our production!

# Retro

- Special thanks to Alexey Zhiltsov (best sysadmin) and our Graphite team!
- It was a great learning journey developing the two features!
- Challenging/Improving existing stack is hard
- Testing, debugging and tooling is important!