

Lessons from Self-Hosting ClickHouse

Boris Tane

 boristane.com

 @boristane



About me

Boris Tane



About Baselime

application observability for the next
frontier of cloud computing

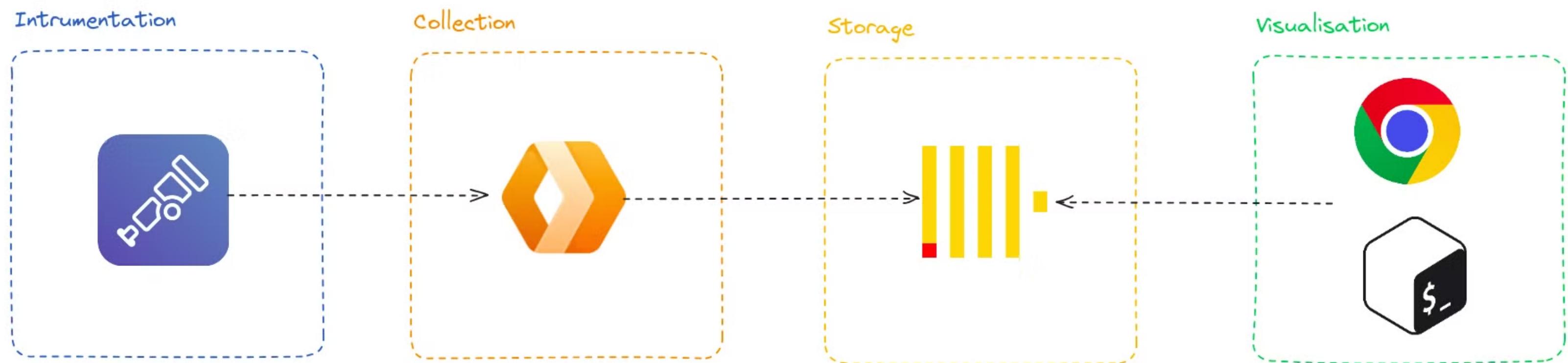
What is observability?

understanding how your software works
by looking at its external outputs

How much data?

~15,000+ events per second

How does it work?



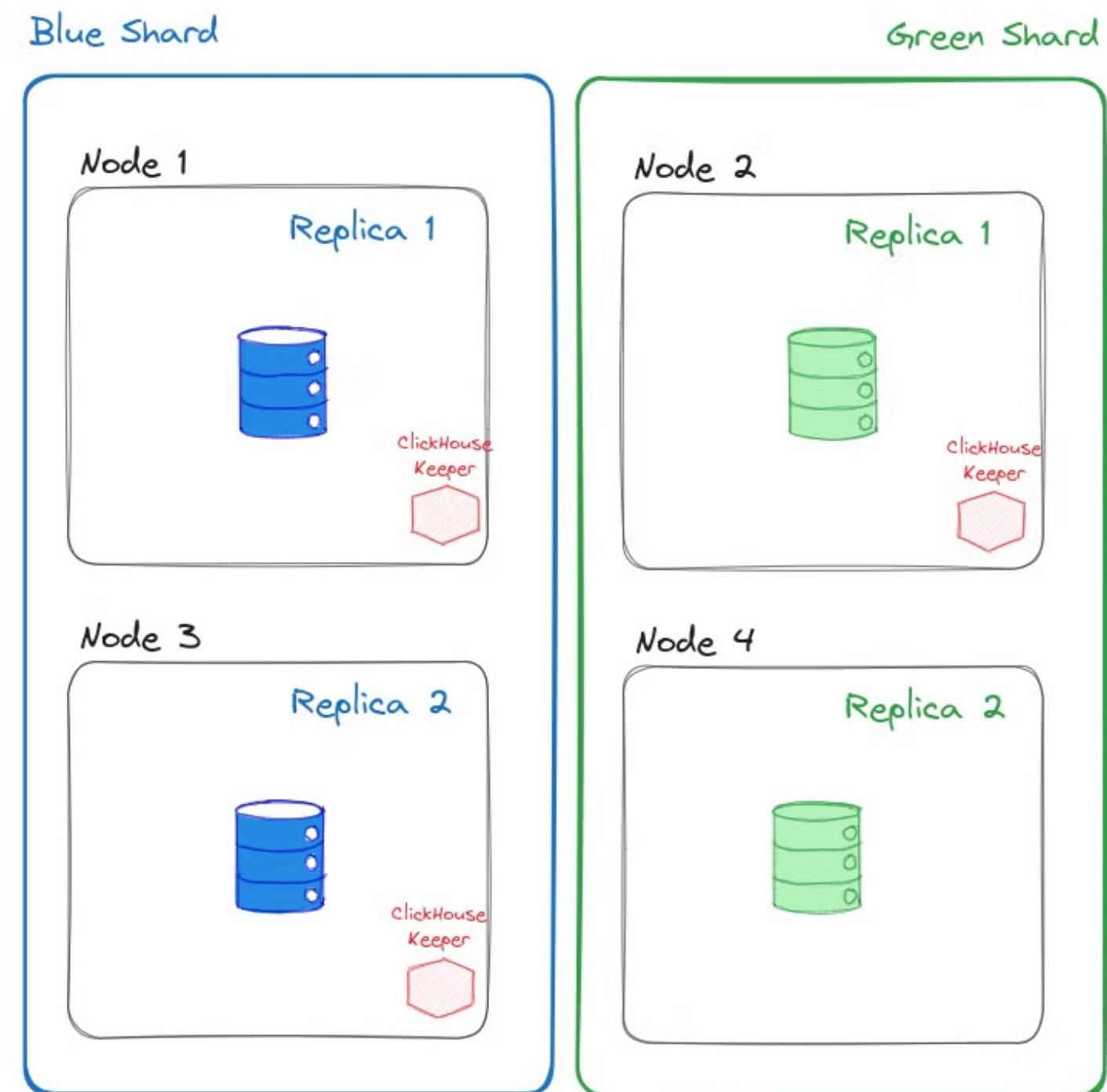
Why ClickHouse?

- Schemaless
- High Cardinality and Dimensionality
- Fast

Why Self-Host?

- More fun :)
- ClickHouse Cloud wasn't around yet
- Unit Economics

What's the architecture?



How's it going?



How's it going?

10 Lessons from operating ClickHouse

Lesson 1: Sending data in small batches is a very bad idea

Lesson 1: Sending data in small batches is a very bad idea

TOO_MANY_PARTS

Lesson 1: Sending data in small batches is a very bad idea

TOO_MANY_PARTS

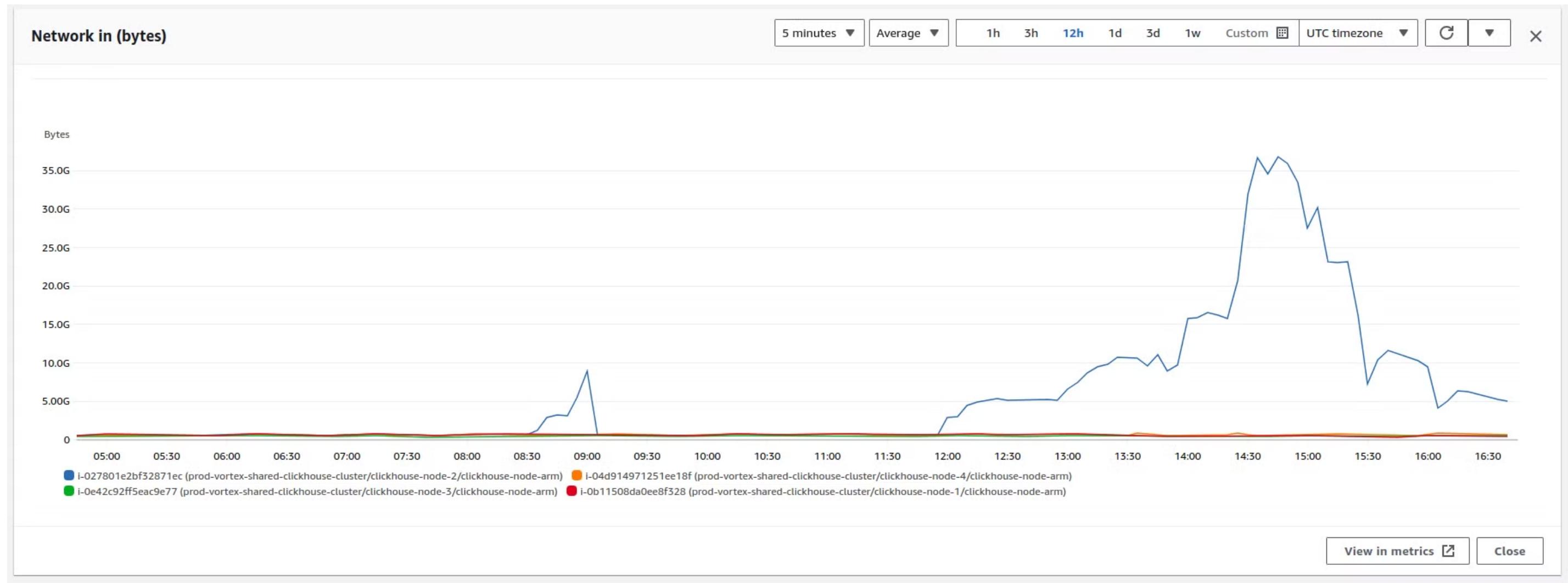
ClickHouse is struggling to merge all the parts as data is coming through

Lesson 2: If in doubt, check that you have enough disk space

Lesson 2: If in doubt, check that you have enough disk space

```
$ df -h
Filesystem      Size  Used Avail Use% Mounted on
/dev/root        88G  7.3G  80G  9% /
/devtmpfs       7.7G    0  7.7G  0% /dev
tmpfs           7.7G    0  7.7G  0% /dev/shm
tmpfs           1.6G  972K  1.6G  1% /run
tmpfs           5.0M    0  5.0M  0% /run/lock
tmpfs           7.7G    0  7.7G  0% /sys/fs/cgroup
/dev/loop0        22M   22M    0 100% /snap/amazon-ssm-agent/7629
/dev/loop1        50M   50M    0 100% /snap/core18/2794
/dev/nvme0n1p15  98M  6.3M  92M  7% /boot/efi
/dev/loop2        36M   36M    0 100% /snap/snapd/20298
/dev/loop3        92M   92M    0 100% /snap/lxd/24065
/dev/nvme1n1      1.4T  1.4T  11G 100% /clickhouse/data
/dev/loop4        60M   60M    0 100% /snap/core20/2019
$ █
```

Lesson 2: If in doubt, check that you have enough disk space



Lesson 3: Beware of the disks you select

Lesson 3: Beware of the disks you select

16TB Limit for General Purpose SSDs
On AWS

Lesson 4: Mutating data is a really bad idea

Lesson 4: Mutating data is a really bad idea

```
ALTER TABLE <table> ON CLUSTER <cluster>
    MODIFY COLUMN <column> Float64;
```

Lesson 4: Mutating data is a really bad idea

ClickHouse creates a new version
of the data and then merges it
with the existing data

Lesson 4: Mutating data is a really bad idea

- Resource-intensive
- Significantly performance impact
- Potentially returns stale data to users

Lesson 5: ClickHouse TTLs mutate data

Lesson 5: ClickHouse TTLs mutate data

- Resource-intensive
- Performance impact

Lesson 5: ClickHouse TTLs mutate data

Instead, delete partitions individually

PARTITION BY toDate(timestamp)

Lesson 5: ClickHouse TTLs mutate data

Instead, delete partitions individually

```
ALTER TABLE <table> ON CLUSTER <cluster> DROP  
PARTITION "2024-05-21"
```

Lesson 6: ClickHouse default values are pretty conservative, don't hesitate to increase them

Lesson 6: ClickHouse default values are pretty conservative, don't hesitate to increase them

max_partition_size_to_drop: 50GB

max_table_size_to_drop: 50GB

keep_alive_timeout: 10s

max_concurrent_queries: now unlimited (was 100)

Lesson 7: Ship your ClickHouse Logs somewhere – don't rely on tail

Filter events - press enter to search

Clear 1m 30m 1h 12h Custom UTC timezone Display ▾

▶	Timestamp	Message
There are older events to load. Load more.		
▼	2024-06-19T11:46:26.254Z	2024.06.19 11:46:21.489890 [12814] {f35092e6-2fad-4e80-acbd-bee49371fbc1} <Warning> HedgedConnectionsFactory: Connection failed at try №1, ... 2024.06.19 11:46:21.489890 [12814] {f35092e6-2fad-4e80-acbd-bee49371fbc1} <Warning> HedgedConnectionsFactory: Connection failed at try №1, reason: Code: 209. DB::NetException: Timeout exceeded while writing to socket (10.0.175.24:9000, 5000 ms). (SOCKET_TIMEOUT) (version 23.9.1.1854 (official build))
▼	2024-06-19T12:59:52.255Z	2024.06.19 12:59:47.984852 [62389] {d9650879-f528-4a57-8aa0-eabe931dab5d} <Warning> HedgedConnectionsFactory: Connection failed at try №1, ... 2024.06.19 12:59:47.984852 [62389] {d9650879-f528-4a57-8aa0-eabe931dab5d} <Warning> HedgedConnectionsFactory: Connection failed at try №1, reason: Code: 209. DB::NetException: Timeout exceeded while writing to socket (10.0.128.18:9000, 5000 ms). (SOCKET_TIMEOUT) (version 23.9.1.1854 (official build))

No newer events at this moment. Auto retry paused. [Resume](#)

Lesson 8: Ship your ClickHouse Metrics somewhere

```
1 "metrics_collected": {
2     "disk": {
3         "measurement": ["used_percent"],
4         "metrics_collection_interval": 60,
5         "resources": ["*"]
6     },
7     "mem": {
8         "measurement": ["mem_used_percent"],
9         "metrics_collection_interval": 60
10    },
11    "diskio": {
12        "measurement": ["read_bytes", "write_bytes", "io_time"],
13        "metrics_collection_interval": 60,
14        "resources": ["/dev/root", "/dev/nvme1n1"]
15    },
16    "cpu": {
17        "measurement": ["usage_idle", "usage_system", "usage_user"],
18        "metrics_collection_interval": 60,
19        "totalcpu": true
20    },
21    "net": {
22        "measurement": ["bytes_sent", "bytes_recv", "packets_sent", "packets_recv", "drop_in", "drop_out"],
23        "metrics_collection_interval": 60,
24        "interface": "eth0"
25    }
26 }
```

Lesson 8: Ship your ClickHouse Metrics somewhere

```
1 <clickhouse>
2   <prometheus>
3     <endpoint>/metrics</endpoint>
4     <port>9363</port>
5     <metrics>true</metrics>
6     <events>true</events>
7     <asynchronous_metrics>true</asynchronous_metrics>
8     <errors>true</errors>
9   </prometheus>
10 </clickhouse>
11
```

Lesson 8: Ship your ClickHouse Metrics somewhere

818 metrics

Lesson 8: Ship your ClickHouse Metrics somewhere

Total time spent merging parts

MergesTimeMilliseconds

Lesson 8: Ship your ClickHouse Metrics somewhere

Bytes sent over the network, for reads,
writes and replications

NetworkSendBytes

Lesson 8: Ship your ClickHouse Metrics somewhere

Bytes received from the network, for reads, writes and replications

NetworkReceiveBytes

Lesson 8: Ship your ClickHouse Metrics somewhere

Slow writes

InsertQueryTimeMicroseconds

Lesson 8: Ship your ClickHouse Metrics somewhere

Slow reads

SelectQueryTimeMicroseconds

Lesson 9: Maintain a directory of useful ClickHouse queries

Lesson 9: Maintain a directory of useful ClickHouse queries

You'll need them when crap hits the fan
at 2am

Lesson 9: Maintain a directory of useful ClickHouse queries

Recent ClickHouse Errors



```
select * from system.errors order by last_error_time desc
```

Lesson 9: Maintain a directory of useful ClickHouse queries

Replication queue - view queues with errors



```
select * from system.replication_queue limit 1 format Vertical;
```

Lesson 9: Maintain a directory of useful ClickHouse queries

Partitions still in sync - potential issues with merges



```
select * from system.merges where elapsed > 10 and progress <  
0.99 limit 20 format Vertical;
```

Lesson 9: Maintain a directory of useful ClickHouse queries

List partitions - view potentially stale partitions



```
select partition, count(), sum(bytes_on_disk) as size from  
system.parts group by partition order by size desc
```

Lesson 10: If something looks odd, even just slightly, something is wrong

- check it before it gets worse

Thanks for having me here!

Boris Tane

 boristane.com

 @boristane



Want to make a presentation like this one?

Start with a fully customizable template, create a beautiful deck in minutes, then easily share it with anyone.

[Create a presentation \(It's free\)](#)