# Clickhouse

as long term storage
for metrics, events, logs
from K8s

exness

# About us

- Platform architecture
- Operations and maintenance
- Troubleshooting
- Deployments

\* и Танцы с бубнами - это про нас)

**exness**

# About

www.exness.com

# Agenda

- Introduction
- Long time ago)
- 1nd implementation (Rancher)
- 2nd implementation (K8s)
- Questions

**exness**

# About

**what we have now**

- **in 2 datacenters**

- **500+ service**

- **2500+ containers**

- **10Krps metrics  up to 200+Krps**

- **2Krps Logs up to 100+Krps**

exness

# Introduction

**Clickhouse in production now**

- **3 clusters**

- **10+ servers**

- **200+ cores, 1Tb ram,   20+ Tb SSD**

exness

# Introduction

**Clickhouse in production now**

- **Easy to replicate**

- **Easy to shard**

- **Easy to use**

- **Easy to manage**

- **Nice support**

- **K8S Operator :)**

exness

# Questions for "full" Clickhouse

A lot of new versions ???
Access rights so simple
ZooKeeper only ?
Clouds ?
UI access ? ( tabix, superset )

1. Zookeeper replacement => Etcd
2. Cloud messages brokers Kinesis, Pub/Sub
3. Auhtorization LDAP, SSO
4. Prometheus metrics exporter
5. GraphQL interface out of the box
6. Clickhouse as Prometheus long-term storage
7. Auto retention policy
8. Detach/Drop/Freeze parts (not only partition as a whole)

exness

# Logs

# Long time ago)



Host

Service #1

Service #2

Service #2

logs

metrics

Graylog

Graphite

Elastic

Whisper

20+ services

# Long time ago)



Host

Service #1

Service #2

logs

Service #2

metrics

Graylog

**1 instance
max 20K rps**

Elastic

**14 days only**

Graphite

Whisper

20+ services

# Long time ago)

# Long time ago)

# Long time ago)

# Long time ago)

exness

Host

Service #1

Service #2

logs

Service #2

metrics

Graylog

**1 instance
max 200K rps ???**

Elastic

**14 days only**

Graphite

Whisper

**2000+** services ???

Throttling ?
Maintenance ?
Long term solution for 1 year

# Long time ago)

# Disadvantages

- **UDP: Throttling ?**
- **Rate: 100+ Krps how ?**
- **Elastic: huge load and a lot of resources**
- **Graylog: Java, fixed load on each node**
- **Retention policy: 1 year or more**

# 1st implementation (Rancher)
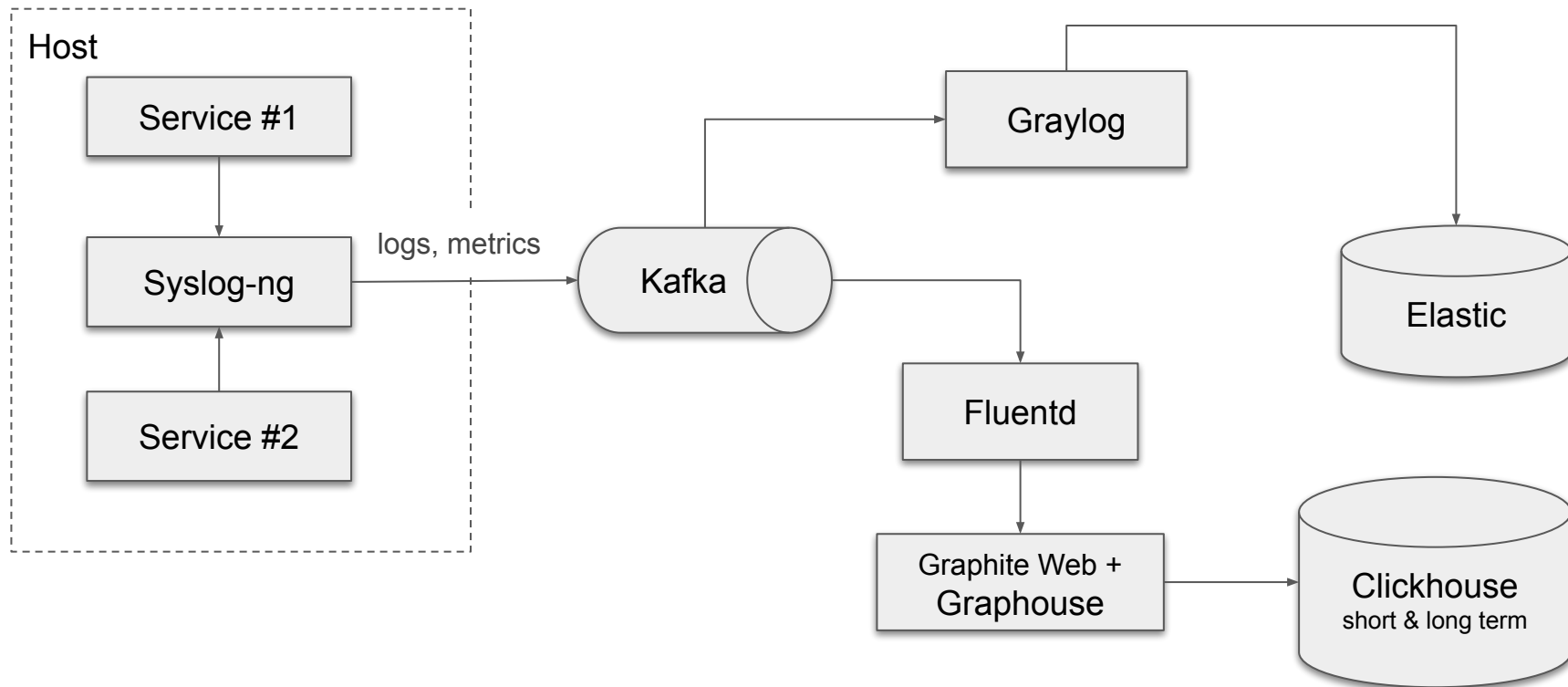
# 1st implementation (Rancher)

# 1nd implementation (disadvantages)

- Not supported for tags (Graphite, Grafana, Graphouse)

- Custom message format (hard to understand)

- Syslog-ng (one point of failure on host, module is written on Java)

- Fluentd pitfalls (offsets, not supported groups, configuration)

- Graphite Web is too slow (issue with long term queries)

- Graphouse is fast (but still written on Java, consumes memory)

trend to next K8s...

# 2nd implementation (K8s)

# Transformation by Consumer



* K8s namespace-related database with it's own metrics, logs, events (ReplicatedMergeTree) and tags (ReplicatedReplacingMergeTree)

# 2nd implementation (advantages)

- No one point of failure (collector on board)

- Fast and robust (Collector & Consumer have written on Golang)

- Throttling on service side

- Standard message format & version support (Telegraf)

- Tags support across logs, metrics and events

- Graylog + Elastic for logs, Prometheus for metrics (short term)

- Clickhouse for long term (metrics, logs, events)

# Thank you!
## Questions?