# Intro to open source observability with Grafana, Prometheus, Loki, and Tempo

Richard "RichiH" Hartmann
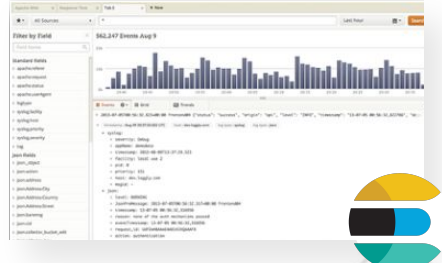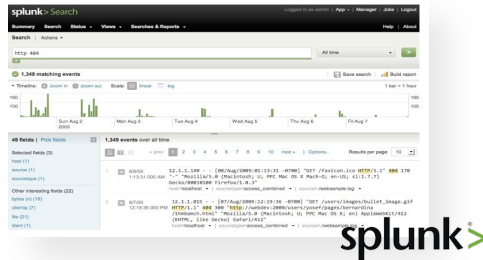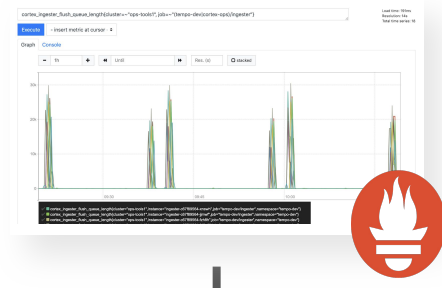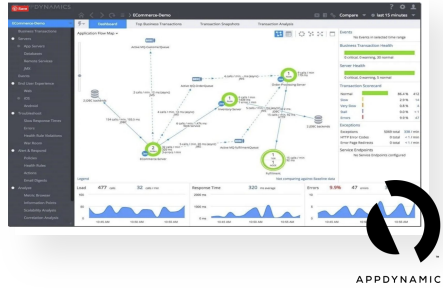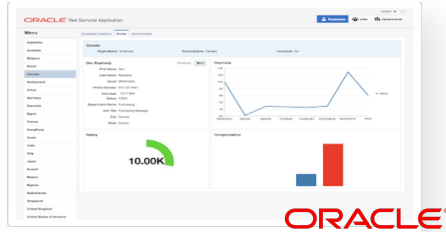
"

I come from the trenches of tech

Bad tools are worse than no tools

# Today's reality: Disparate systems. Disparate data.

"

# Back to the basics

# Let's rethink this

# Observability & SRE

Or: Buzzwords, and their useful parts

# Dirty secret

- What is "cloud native scale" is "Internet scale" of networks two decades ago
- The combination of metrics and events has been standard in power measurement for half a century
- Modern tools with good engineering practices map very well onto brownfield technology

# Buzzword alert!

- Cool new term, almost meaningless by now, what does it mean?
  - Pitfall alert: Cargo culting
  - It's about changing the behaviour, not about changing the name
- "Monitoring" has taken on a meaning of collecting, not using data
  - One extreme: Full text indexing
  - Other extreme: Data lake
- "Observability" is about enabling humans to understand complex systems
  - Ask **why** it's not working instead of just knowing that it's not

Grafana Labs

"

If you can't ask new questions on the fly, it's not observability

# Complexity

- Fake complexity, a.k.a. bad design
  - Can be **reduced**
- Real, system-inherent complexity
  - Can be **moved** (monolith vs client-server vs microservices)
  - Must be **compartmentalized** (service boundaries)
  - Should be **distilled meaningfully**

# SRE, an instantiation of DevOps

- At its core: Align incentives across the org
  - Error budgets allow devs, ops, PMs, etc. to optimize for shared benefits
- Measure it!
  - SLI: Service Level Indicator: What you measure
  - SLO: Service Level Objective: What you need to hit
  - SLA: Service Level Agreement: When you need to pay

# Shared understanding

- Everyone uses the same tools & dashboards
  - Shared incentive to invest into tooling
  - Pooling of institutional system knowledge
  - Shared language & understanding of services

# Services

- Service?
  - Compartmentalized complexity, with an interface
  - Different owners/teams
  - Contracts define interfaces
- Why "contract": Shared agreement which MUST NOT be broken
  - Internal and external customers rely on what you build and maintain
- Other common term: layer
  - The Internet would not exist without network layering
  - Enables innovation, parallelizes human engineering
- Other examples: CPUs, harddisk, compute nodes, your lunch

# Alerting

- Customers care about services being up, not about individual components

- Discern between different SLIs
  - Primary: service-relevant, for alerting
  - Secondary: informational, debugging, might be underlying's primary

**Anything currently or imminently impacting customer service must be alerted upon**
**But nothing(!) else**

# Prometheus

Grafana Labs

# Prometheus 101

- Inspired by Google's Borgmon

- Time series database

- unit64 millisecond timestamp, float64 value

- Instrumentation & exporters

- Not for event logging

- Dashboarding via Grafana

# Main selling points

- Highly dynamic, built-in service discovery

- No hierarchical model, n-dimensional label set

- PromQL: for processing, graphing, alerting, and export

- Simple operation

- Highly efficient

# Main selling points

- Prometheus is a pull-based system

- Black-box monitoring: Looking at a service from the outside (Does the server answer to HTTP requests?)

- White-box monitoring: Instrumenting code from the inside (How much time does this subroutine take?)

- Every service should have its own metrics endpoint

- Hard API commitments within major versions

Grafana Labs

# Time series

- Time series are recorded values which change over time
- Individual events are usually merged into counters and/or histograms
- Changing values are recorded as gauges
- Typical examples
  - Requests to a webserver (counter)
  - Temperatures in a datacenter (gauge)
  - Service latency (histograms)

# Super easy to emit, parse & read

```
http_requests_total{env="prod",method="post",code="200"} 1027
http_requests_total{env="prod",method="post",code="400"} 3
http_requests_total{env="prod",method="post",code="500"} 12
http_requests_total{env="prod",method="get",code="200"} 20
http_requests_total{env="test",method="post",code="200"} 372
http_requests_total{env="test",method="post",code="400"} 75
```

# Scale

- Kubernetes is Borg

- Prometheus is Borgmon

- Google couldn't have run Borg without Borgmon (plus Omega and Monarch)

- Kubernetes & Prometheus are designed and written with each other in mind

# Prometheus scale

- 1,000,000+ samples/second no problem on current hardware

- ~200,000 samples/second/core

- 16 bytes/sample compressed to 1.36 bytes/sample

**The highest we saw in production on a single Prometheus instance were 15 million active times series at once!**

# Long term storage

- Two long term storage solutions have Prometheus-team members working on them
  - Thanos
    - Historically easier to run, but slower
    - Scales storage horizontally
  - Cortex
    - Easy to run these days
    - Scales both storage, ingester, and querier horizontally

# Cortex @ Grafana (largest cluster, 2021-09)

- ~65 million active series (just the one cluster)

- 668 CPU cores

- 3,349 GiB RAM

One customer running at 3 billion active series

# Loki

Grafana Labs

# Loki 101

- Following the same label-based system like Prometheus

- No full text index needed, incredible speed

- Work with logs at scale, without the massive cost

- Access logs with the same label sets as metrics

- Turn logs into metrics, to make it easier to work with them

- Make direct use of syslog data, via promtail

2019-12-11T10:01:02.123456789Z `{env="prod",instance="1.1.1.1"}` GET /about

**Timestamp**

with nanosecond precision

**Prometheus-style Labels**

key-value pairs

**Content**

log line

indexed                    unindexed

Grafana Labs

# Loki @ Grafana Labs

- Queries regularly see 40 GiB/s

- Query terabytes of data in under a minute
  - Including complex processing of result sets

Grafana Labs

# Tempo

Grafana Labs

# Tempo

- Exemplars: Jump from relevant logs & metrics
  - Native to Prometheus, Cortex, Thanos, and Loki
  - Exemplars work at Google scale, with the ease of Grafana
- Index and search by labelsets available for those who need it
- Object store only: No Cassandra, Elastic, etc.
- 100% compatible with OpenTelemetry Tracing, Zipkin, Jaeger
- 100% of your traces, no sampling

# Tempo @ Grafana Labs (2021-07)

- 2,200,000 samples per second @ 350 MiB/s

- 14-day retention @ 3 copies stored

- ~240 CPU cores (includes compression cost)

- ~450 GiB RAM

- 132 TiB object storage

- Latencies:

  - p99 - 2.5s

  - p90 - 2.3s

  - p50 - 1.6s

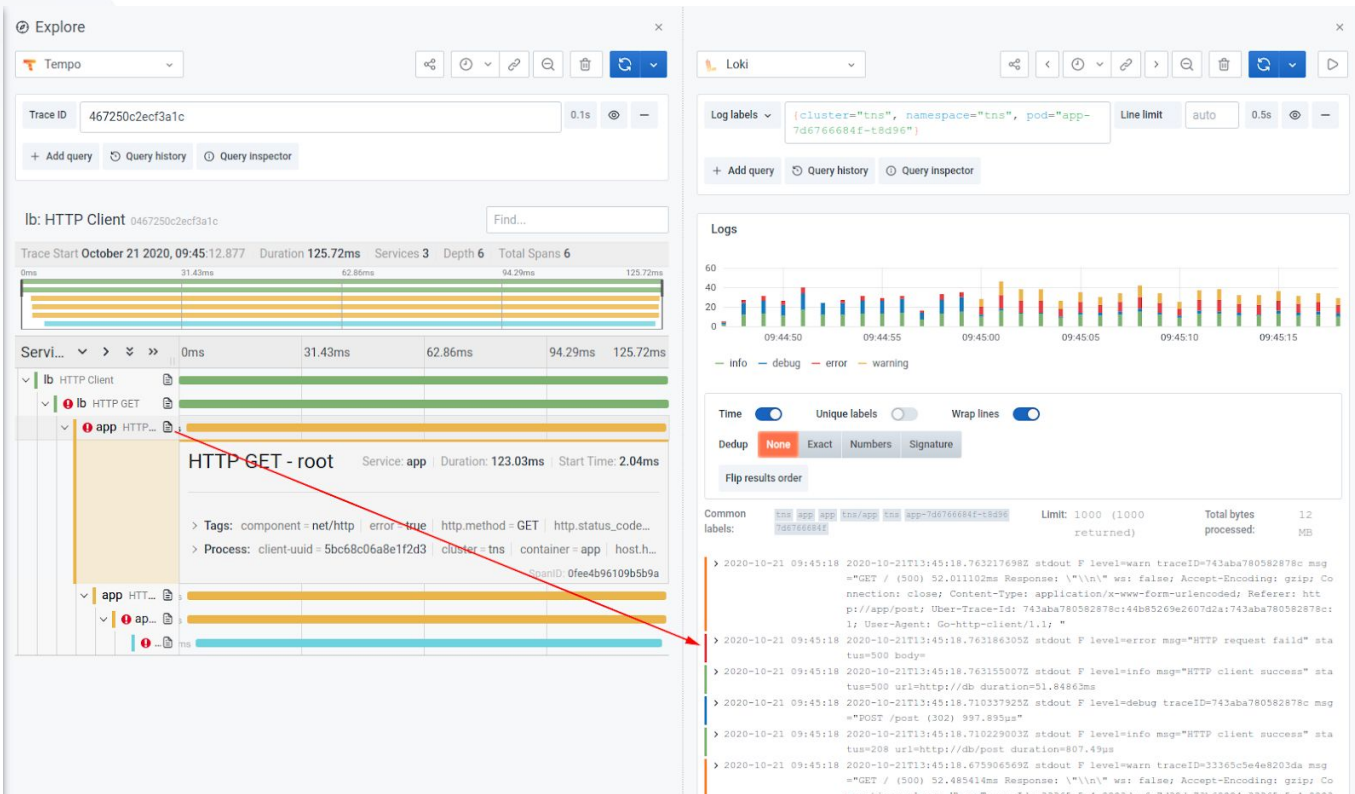Grafana Labs

# Bringing it together

# From logs to traces

# From metrics to traces

# From metrics to traces

# ...and from traces to logs

> "
**All of this is Open Source and you can run it yourself**

## server requests

- web_server_01
- web_server_02
- web_server_03
- web_server_04

## Memory / CPU

- cpu
- memory

## client side full page load

| | Mean |
|---|---|
| upper_25 | 6.16 ms |
| upper_50 | 122 ms |
| upper_75 | 538 ms |
| upper_90 | 1.00 s |
| upper_95 | 1.40 s |

## logins

- logins
- logins (-1 hour)

## Traffic In/Out

| | avg | current |
|---|---|---|
| upper_25 | 6.40 B | 1 B |
| upper_50 | 124 B | 6 B |
| upper_75 | 544 B | 478 B |
| upper_90 | 1.01 kB | 599 B |
| upper_95 | 1.40 kB | 627 B |
| cpu1 | -1.4 kB | -627 B |
| cpu2 | -544 B | -478 B |

## Hide Grid Lines

- A-series
- B-series
- C-series

For a deeper dive into
Open Source Observability, go to:

# grafana.com/go/obscon2021

# Thank you!

richih@richih.org
twitter.com/TwitchiH
github.com/RichiH/talks