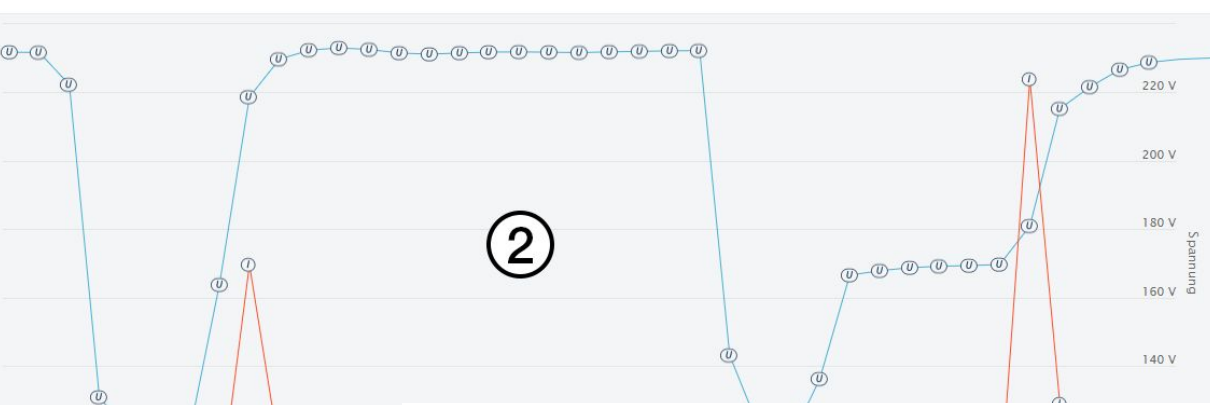# Cloud native observability with Prometheus and beyond

## Or: philosophy of Observability

Richard "RichiH" Hartmann

# How humans deal with data

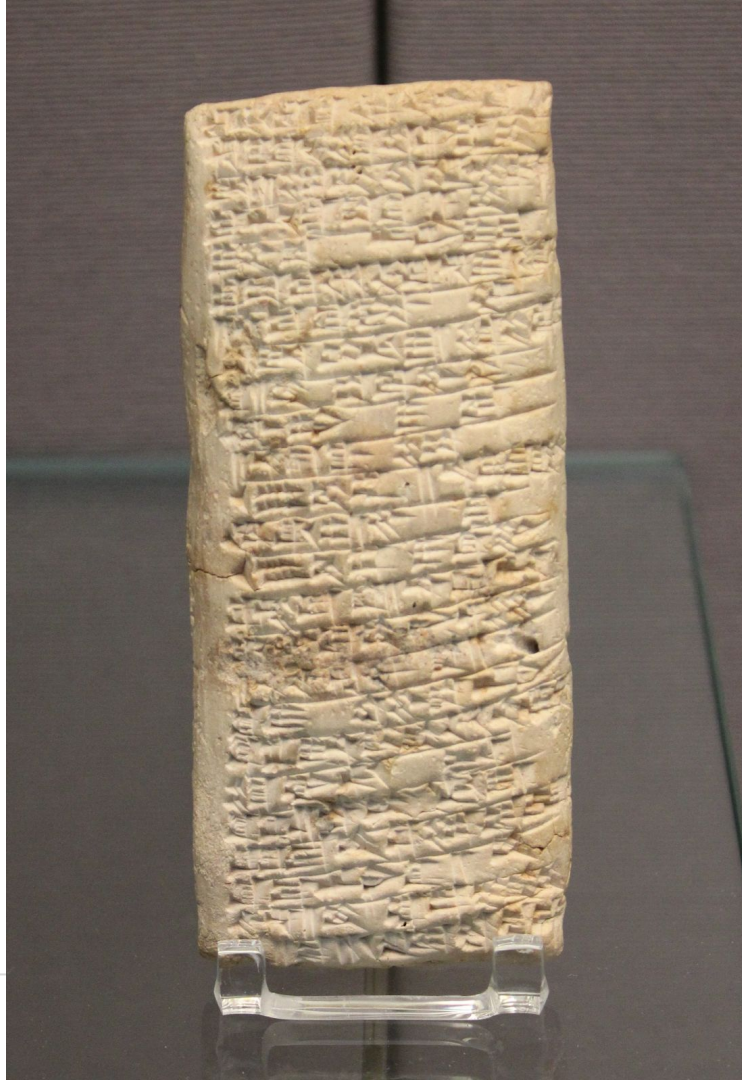**Event Browser**

Geräte auswählen | Geräte suchen

19.12.2019 - 19.12.2019
00:00 - 23:59

**Statistische Auswertung**

| | |
|---|---|
| Ereignisse | 8 |
| Transienten | 14 |

**E Ereignis Übersicht**

| | |
|---|---|
| Durchschnittliche Dauer | 160 ms |
| Längste Dauer | 220 ms |
| Unterspannung | 8 |

**T Transienten Übersicht**

| | |
|---|---|
| Einhüllende | 14 |

**Betroffene Phasen**

| | T | E |
|---|---|---|
| Phase L1 | 2 | 2 |
| Phase L2 | 4 | 2 |
| Phase L3 | 4 | 2 |
| Phase L4 | 4 | 2 |

| | Gerätename | Typ | Phase | Start | Ende | Dauer | Wert | |
|---|---|---|---|---|---|---|---|---|
| T | UMG 512 - TD | Einhüllende | L2 | 19.12.2019 01:15:40'663 | --- | --- | --- | Analysieren |
| T | UMG 512 - TD | Einhüllende | L3 | 19.12.2019 01:15:40'663 | --- | --- | --- | Analysieren |
| T | UMG 512 - TD | Einhüllende | L4 | 19.12.2019 01:15:40'663 | --- | --- | --- | Analysieren |
| T | UMG 512 - TD | Einhüllende | L1 | 19.12.2019 01:15:39'345 | --- | --- | --- | Analysieren |
| T | UMG 512 - TD | Einhüllende | L2 | 19.12.2019 01:15:39'345 | --- | --- | --- | Analysieren |
| T | UMG 512 - TD | Einhüllende | L3 | 19.12.2019 01:15:39'345 | --- | --- | --- | Analysieren |
| T | UMG 512 - TD | Einhüllende | L4 | 19.12.2019 01:15:39'345 | --- | --- | --- | Analysieren |
| E | UMG 508 - TD | Unterspannung | L1 | 19.12.2019 00:30:59'199 | 19.12.2019 00:30:59'199 | 220 ms | 117.836 V (MIN) | Analysieren |
| E | UMG 508 - TD | Unterspannung | L2 | 19.12.2019 00:30:59'199 | 19.12.2019 00:30:59'199 | 220 ms | 117.806 V (MIN) | Analysieren |
| E | UMG 508 - TD | Unterspannung | L3 | 19.12.2019 00:30:59'199 | 19.12.2019 00:30:59'199 | 220 ms | 117.825 V (MIN) | Analysieren |
| E | UMG 508 - TD | Unterspannung | L4 | 19.12.2019 00:30:58'979 | 19.12.2019 00:30:59'199 | 220 ms | 117.830 V (MIN) | Analysieren |
| E | UMG 508 - TD | Unterspannung | L1 | 19.12.2019 00:30:58'559 | 19.12.2019 00:30:58'659 | 100 ms | 118.640 V (MIN) | Analysieren |
| E | UMG 508 - TD | Unterspannung | L2 | 19.12.2019 00:30:58'559 | 19.12.2019 00:30:58'659 | 100 ms | 118.612 V (MIN) | Analysieren |
| E | UMG 508 - TD | Unterspannung | L3 | 19.12.2019 00:30:58'559 | 19.12.2019 00:30:58'659 | 100 ms | 118.622 V (MIN) | Analysieren |
| E | UMG 508 - TD | Unterspannung | L4 | 19.12.2019 00:30:58'659 | 19.12.2019 00:30:58'659 | 100 ms | 118.631 V (MIN) | Analysieren |
| T | UMG 512 - TD | Einhüllende | L2 | 19.12.2019 00:28:40'565 | --- | --- | --- | Analysieren |
| T | UMG 512 - TD | Einhüllende | L3 | 19.12.2019 00:28:40'565 | --- | --- | --- | Analysieren |
| T | UMG 512 - TD | Einhüllende | L4 | 19.12.2019 00:28:40'565 | --- | --- | --- | Analysieren |
| T | UMG 512 - TD | Einhüllende | L1 | 19.12.2019 00:28:40'130 | --- | --- | --- | Analysieren |
| T | UMG 512 - TD | Einhüllende | L2 | 19.12.2019 00:28:40'130 | --- | --- | --- | Analysieren |

20 | 1 2

## Unlegierte Werkzeugstähle / Carbon tool steels

| Werkstoff-Nr. / Standard No. | Gustav Grimm Stahl- und Hammerwerk Düsseldorf | Haeckerstahl Willy Haecker & Co. Remscheid-Hasten | Hagenbütter & Co. Hückeswagen-Hohenhagen Solingen | Henricot Stahlwerke Court-Saint-Etienne Belgien | Stachelhauser Stahl- u. Walzwerke Hessenbruch & Co. G.m.b.H. Remscheid | Chr. Höver & Sohn Burghausen Bez. Aachen |
|---|---|---|---|---|---|---|

*(Carbon tool steels and high-speed steels comparison table — data cells not legibly resolvable at this resolution.)*

## Schnellarbeitsstähle / High speed steels

DEUTSCHLAND — GERMANY-WEST

## Unlegierte Werkzeugstähle / Carbon tool steels

| Werkstoff-Nr. / Standard No. | Gebrüder Höver Edelstahlwerk | Hofmann & Co. Edelstahle Köln-Rhein | Horbach & Schmitz GmbH | Georg Hufnagel Edelstähle Nürnberg | IBACH Stahlwerke Remscheid-Vgbn. | Kind & Co. Edelstahlwerk | Fritz Kölsche K.G. Edelstähle Herdecke 2 |
|---|---|---|---|---|---|---|---|

*(Data cells not legibly resolvable at this resolution.)*

## Schnellarbeitsstähle / High speed steels

"

**Humanity has optimized detailed accounts into key events into numbers for millenia, again and again**

Grafana Labs

# Observability & SRE

Or: Buzzwords, and their useful parts

# 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

"

**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**

# Services

- What's a 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

# 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

Grafana Labs

# Shared understanding

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

# 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

# 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

# Concepts & guarantuees

- 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

- New release candidate every six weeks

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
```

Grafana Labs

# PromQL

All partitions in my entire infrastructure with more than 100GB capacity that are not mounted on root?

```
node_filesystem_bytes_total{mountpoint!="/"} / 1e9 > 100

 {device="sda1", mountpoint="/home", instance="10.0.0.1"} 118.8
 {device="sda1", mountpoint="/home", instance="10.0.0.2"} 118.8
 {device="sdb1", mountpoint="/data", instance="10.0.0.2"} 451.2
 {device="xdvc", mountpoint="/mnt", instance="10.0.0.3"} 320.0
```

# PromQL

What's the ratio of request errors across all service instances?

```
sum by(path) (rate(http_requests_total{status="500"}[5m])) /
 sum by(path) (rate(http_requests_total[5m]))

 {path="/status"} 0.0039
 {path="/"} 0.0011
 {path="/api/v1/topics/:topic"} 0.087
 {path="/api/v1/topics} 0.0342
```

Grafana Labs

# New features

- Remote Write Receiver (v2.25 (feature flag) v2.33 (stable))

- Trigonometric functions (v2.31)

- Agent mode (v2.32)

- Long term support versions (v2.27)

- Out of order ingestion (v2.39)

Next highlight feature: Native histograms

# Cloud native defaults

- 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

- Reliable into the tens of millions of active series

Grafana Labs

# Mimir

# Mimir

- For **M**etrics

- Prometheus -> Cortex -> Grafana Enterprise Metrics -> Mimir

- Scales to more than 1,000,000,000 Active Series

- Blazingly fast query performance

- Hard multi-tenancy, access control, and three-way replication

- Can ingest native OpenTelemetry, DataDog, Graphite, and Influx

Grafana Labs

# Mimir @ Grafana

- 1,000,000,000 Active Series - in one cluster

- 1,500 machines

- 7,000 CPU cores

- 30 TiB RAM

Grafana Labs

# Loki

# Loki 101

- For **L**ogs
- Following the same label-based system as Prometheus
  - Only index what you need often, query the rest
  - "Index the labels, query the data"
- Work with logs at scale, without the massive cost
  - Scalable low latency write path
  - Flexible schema on read
- Access logs with the same label sets as metrics
  - Turn logs into metrics, to make it easier & cheaper to work with them 📊

`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

- Largest user cluster (as of 2022-09): 180 TiB per day

- Queries regularly see 80 GiB/s

- Query terabytes of data in under a minute

  - Including complex processing of result sets

# Tempo

Grafana Labs

# Tempo

- For **T**races

- 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 (2022-09)

- 2,200,000 samples per second @ 350 MiB/s
  - 5,000,000 samples second peak
- 14-day retention @ 3 copies stored
- Latencies:
  - p99 - 2.5s
  - p90 - 2.3s
  - p50 - 1.6s

# Phlare

# Phlare

- For **P**rofiling

- Pronounced "Flare"

- Profiles

  - "How much CPU & RAM am I spending in what areas of the code?"

  - "...and how does this change over time?"

- Go: pprof

- Java: https://github.com/grafana/JPProf

# It's a numbers game

# Logs to metrics, the savings

- Full text indexing: 10 TiB logs -> ~20 TiB index

- Loki: 10 TiB logs -> ~200 MiB index

- Logs @ Grafana ~600 B average per line

- Metrics ~1.36 byte per metric sample

-> 99.8% reduction in storage size for first log line

~100% for every follow-up log line

# Bringing it together

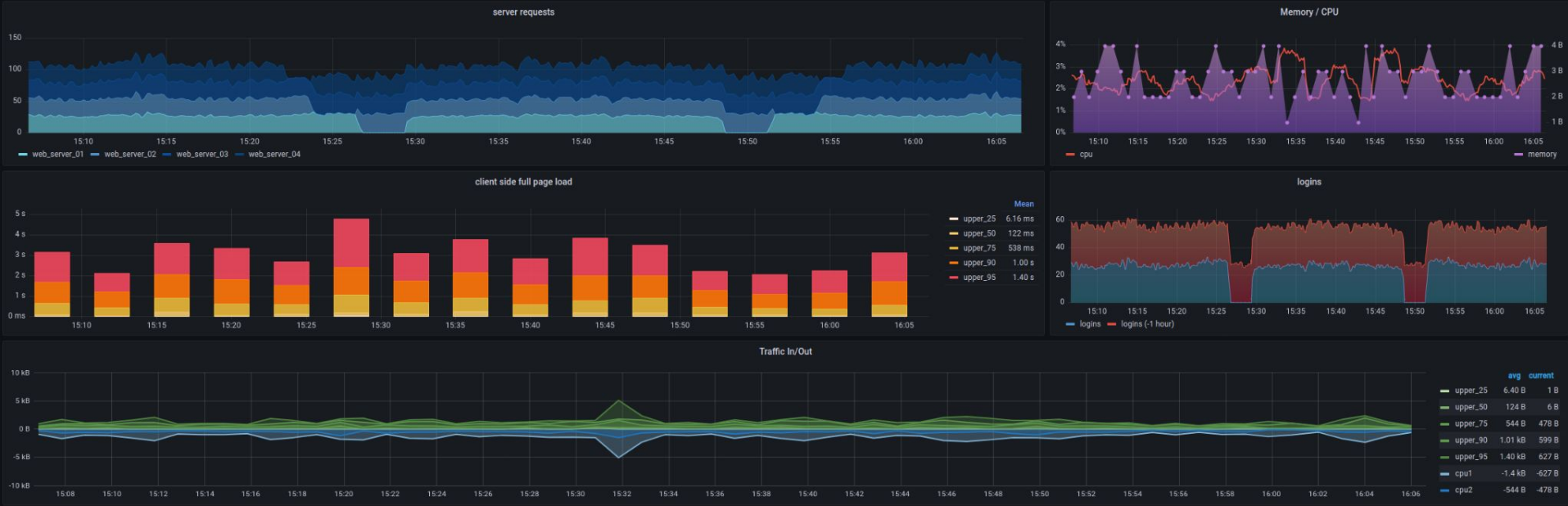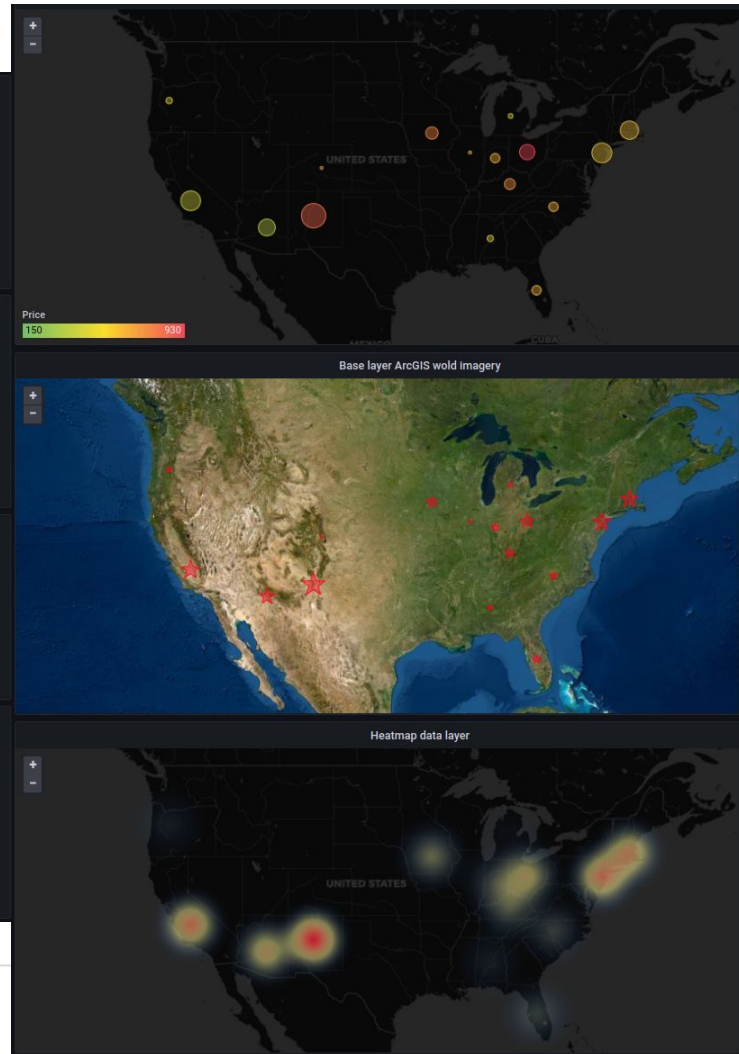# From logs to traces

# 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**

**(But we will also sell it to you happily)**

Grafana Labs

# Thank you!

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