



Philosophy of Observability

Leveraging Cloud-Native Technologies



Richard "RichiH" Hartmann

My Background

- Senior Developer Programs Director @ Grafana Labs
- Prometheus team member (CNCF graduated project)
 - PromCon lead, Prometheus Dev Summit chair
- CNCF Governing Board
- CNCF Technical Oversight Committee
- CNCF Technical Advisory Group Observability chair
- OpenMetrics founder
- OpenTelemetry member

My offer

- Built Europe's most modern datacenter, monitored through SNMP and Modbus/TCP only
- Maintainer of
 - https://github.com/prometheus/snmp_exporter
 - https://github.com/RichiH/modbus_exporter
 - Others..

If you have a cool / large factory to visit and want to trade a talk

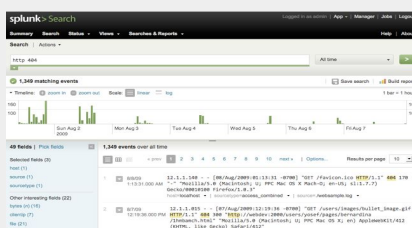
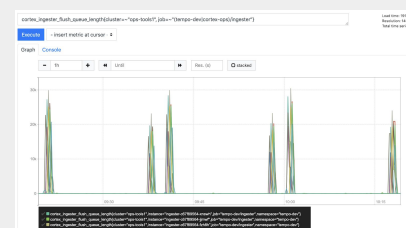
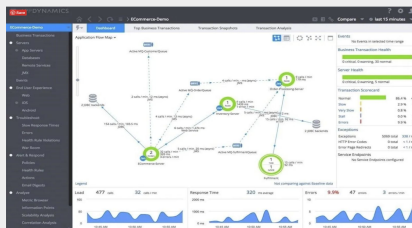
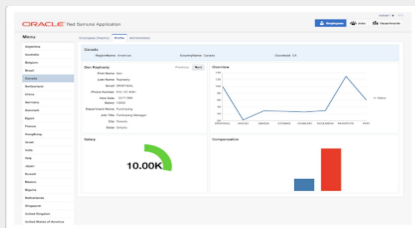
My Background

- Ran the backbone of an ISP for eleven years
 - Was the only person on call
 - My life & sanity depended on on-point monitoring & alerting
- Active in RIPE, IETF, DENOG, #networker etc
 - RFC to my name, changed RIPE NCC's IPv4 PI policies, etc.
- Prometheus transition for Germany's oldest ISP, ~5k devices
- Staffed world's largest IRC network for more than a decade
- Run conferences & monitoring from 100s to 18k attendees
 - DENOG, DebConf, FOSDEM, CCC, GrafanaCon, PromCon

Today's reality: Disparate systems. Disparate data.

ORACLE

APPDYNAMICS



Grafana

splunk>

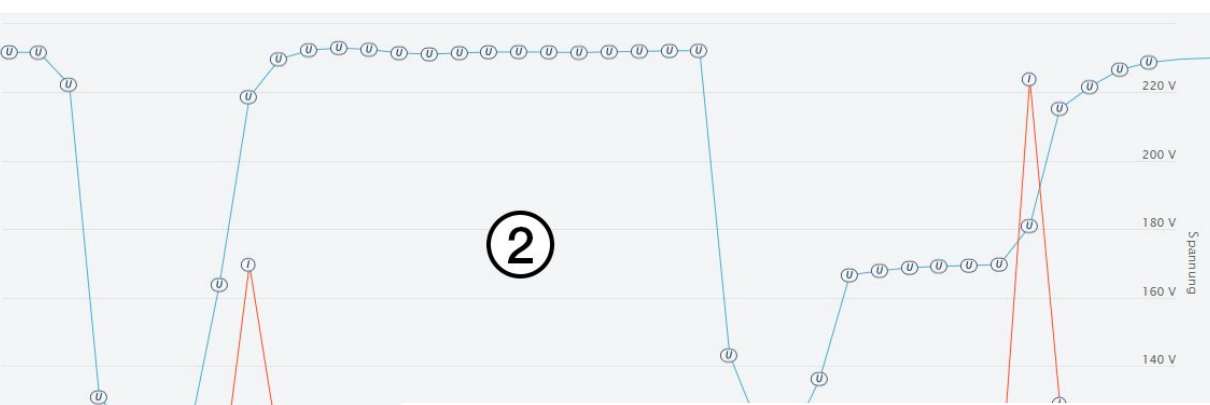
elasticsearch



Back to the basics

Let's rethink this

How humanity deals with data



Event Browser

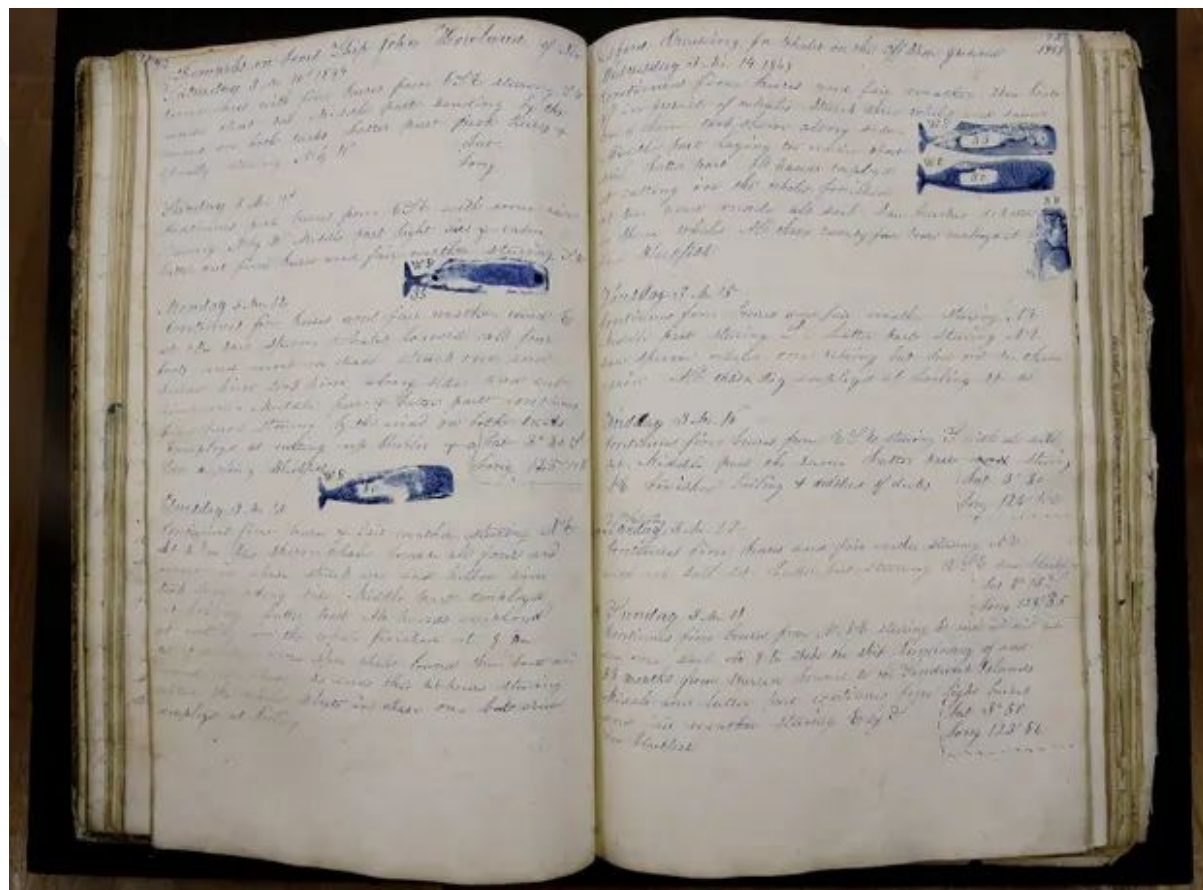
Geräte auswählen

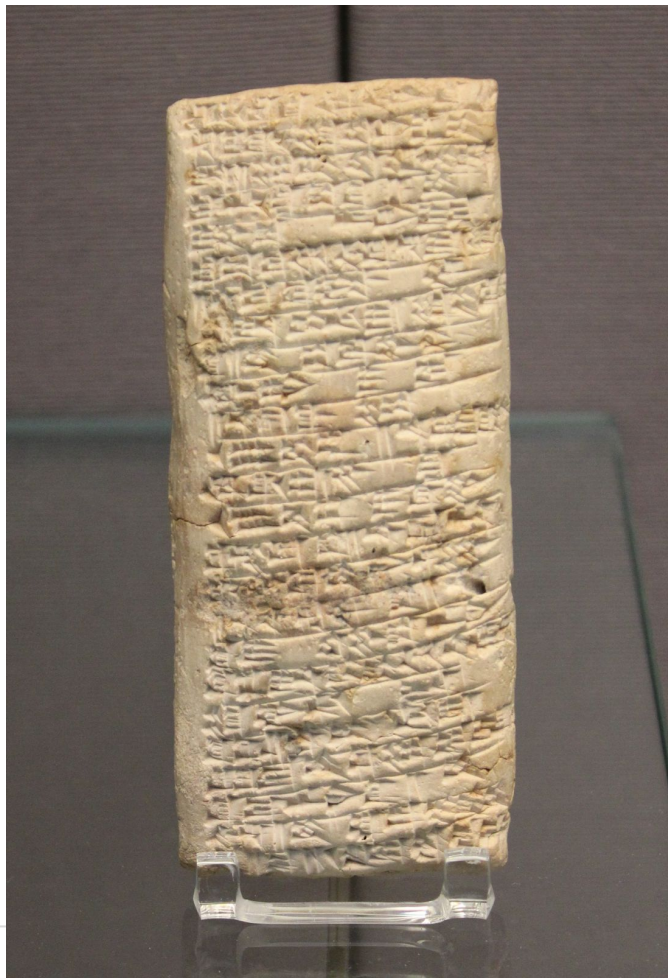
19.12.2019 - 19.12.2019
00:00 - 23:59

Statistische Auswertung	Gerätename	Typ	Phase	Start	Ende	Dauer	Wert	
Ereignisse 8	UMG 512 - TD	Einhüllende	[1] [2] [3] [4] [5] [6]	19.12.2019 01:15:40'663	---	---	---	Analysieren
Transienten 14	UMG 512 - TD	Einhüllende	[1] [2] [3] [4] [5] [6]	19.12.2019 01:15:40'663	---	---	---	Analysieren
Ereignis Übersicht	UMG 512 - TD	Einhüllende	[1] [2] [3] [4] [5] [6]	19.12.2019 01:15:40'663	---	---	---	Analysieren
Durchschnittliche Dauer	UMG 512 - TD	Einhüllende	[1] [2] [3] [4] [5] [6]	19.12.2019 01:15:39'345	---	---	---	Analysieren
Längste Dauer	UMG 512 - TD	Einhüllende	[1] [2] [3] [4] [5] [6]	19.12.2019 01:15:39'345	---	---	---	Analysieren
Unterspannung	UMG 512 - TD	Einhüllende	[1] [2] [3] [4] [5] [6]	19.12.2019 01:15:39'345	---	---	---	Analysieren
Transienten Übersicht	UMG 512 - TD	Einhüllende	[1] [2] [3] [4] [5] [6]	19.12.2019 01:15:39'345	---	---	---	Analysieren
Einhüllende	UMG 508 - TD	Unterspannung	[1] [2] [3] [4] [5] [6]	19.12.2019 00:30:59'199	19.12.2019 00:30:59'199	220 ms	117.836 V (MIN)	Analysieren
Betroffene Phasen	UMG 508 - TD	Unterspannung	[1] [2] [3] [4] [5] [6]	19.12.2019 00:30:59'979	19.12.2019 00:30:59'199	220 ms	117.806 V (MIN)	Analysieren
Phase L1	UMG 508 - TD	Unterspannung	[1] [2] [3] [4] [5] [6]	19.12.2019 00:30:58'979	19.12.2019 00:30:59'199	220 ms	117.825 V (MIN)	Analysieren
Phase L2	UMG 508 - TD	Unterspannung	[1] [2] [3] [4] [5] [6]	19.12.2019 00:30:58'979	19.12.2019 00:30:59'199	220 ms	117.830 V (MIN)	Analysieren
Phase L3	UMG 508 - TD	Unterspannung	[1] [2] [3] [4] [5] [6]	19.12.2019 00:30:58'559	19.12.2019 00:30:58'659	100 ms	118.640 V (MIN)	Analysieren
Phase L4	UMG 508 - TD	Unterspannung	[1] [2] [3] [4] [5] [6]	19.12.2019 00:30:58'559	19.12.2019 00:30:58'659	100 ms	118.612 V (MIN)	Analysieren
	UMG 508 - TD	Unterspannung	[1] [2] [3] [4] [5] [6]	19.12.2019 00:30:58'559	19.12.2019 00:30:58'659	100 ms	118.622 V (MIN)	Analysieren
	UMG 508 - TD	Unterspannung	[1] [2] [3] [4] [5] [6]	19.12.2019 00:30:58'559	19.12.2019 00:30:58'659	100 ms	118.631 V (MIN)	Analysieren
	UMG 512 - TD	Einhüllende	[1] [2] [3] [4] [5] [6]	19.12.2019 00:28:40'565	---	---	---	Analysieren
	UMG 512 - TD	Einhüllende	[1] [2] [3] [4] [5] [6]	19.12.2019 00:28:40'565	---	---	---	Analysieren
	UMG 512 - TD	Einhüllende	[1] [2] [3] [4] [5] [6]	19.12.2019 00:28:40'565	---	---	---	Analysieren
	UMG 512 - TD	Einhüllende	[1] [2] [3] [4] [5] [6]	19.12.2019 00:28:40'130	---	---	---	Analysieren
	UMG 512 - TD	Einhüllende	[1] [2] [3] [4] [5] [6]	19.12.2019 00:28:40'130	---	---	---	Analysieren

20 < 1 >

6472







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

Again and again and again



Observability & SRE

Or: Buzzwords, and their useful parts

Observability, the buzzword

- 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 new questions on the fly
 - Ask **why** it's not working instead of just knowing that it's not
- Terms such as "Observability 2.0", "Observability 3.0", and "Observability 4.0", are other examples of buzzwordiness



[...] observations are [...] approximations to the truth [...] this can be accomplished in no other way than by a suitable combination of more observations than the number absolutely requisite for the determination of the unknown quantities

Carl Friedrich Gauß, 1809



Observability is a measure of how well internal states of a system can be inferred from knowledge of its external outputs.

Rudolf Emil Kálmán, 1960



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 (observability...)
 - Can **not** be reduced


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, hard drive, compute nodes, your lunch 🌻

Cloud-native vs client-server vs mainframe vs...

- A mainframe application and a microservices fleet are fundamentally the same
 - You can *move* system-inherent complexity, but...
- Microservices broke up old service and system boundaries
 - Enabling horizontal scalability, arguably at the cost of vertical scalability
- Previous-generation tooling is designed to understand system complexity along existing service boundaries
 - Cloud native tooling is able to deal with this increased complexity
 - NB: This means previous-generation complexity is even easier to observe

SRE

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

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

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
- Rich ecosystem, 1,000s of instrumentations & exporters
- Cloud-native default

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)

Cloud-native default

- Kubernetes =~ Borg
- Prometheus =~ Borgmon
- Google couldn't have run Borg without Borgmon
- Kubernetes & Prometheus are designed and written with each other in mind

Prometheus 101

- 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

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

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
```

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
```

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



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

Mimir @ Grafana

- 1,000,000,000 Active Series - in one cluster
- 1,500 machines
- 7,000 CPU cores
- 30 TiB RAM



Loki

Loki 101

- For Logs
- 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

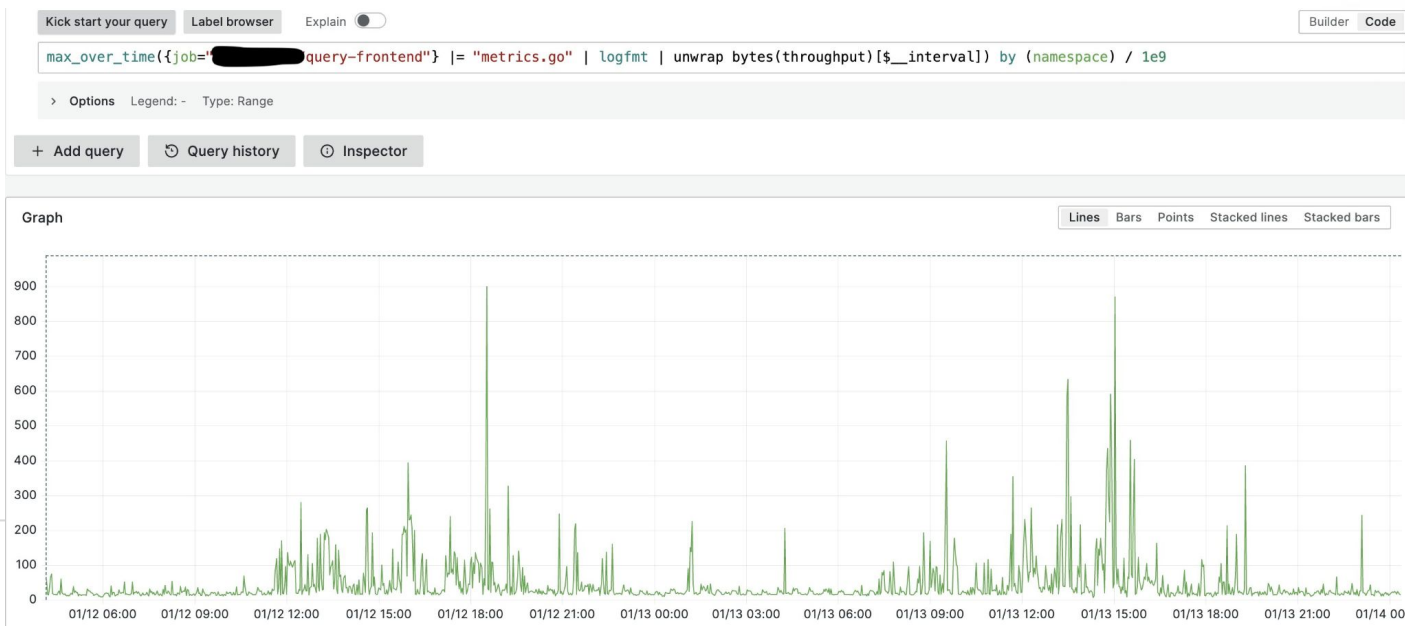
log line

indexed

unindexed

Loki @ Grafana Labs

- Largest user cluster: 180 TiB per day
- Queries regularly peak at **900GB/s**
- Query 10TB in 12 seconds, including complex processing of result sets





Tempo

Tempo

- For **Traces**
- Historic problem: Traces require extremely rich metadata for analysis
 - Expensive, slow, and mandates sampling
- Exemplars: Leverage the extracted logs & metrics
 - Exemplars work at Google scale, with the ease of Grafana
 - Native to Prometheus, Cortex, Thanos, and Loki
- Index and search by labelsets available for those who need it
- 100% compatible with OpenTelemetry Tracing, Zipkin, Jaeger

Tempo @ Grafana Labs

- 1,500,000 samples per second @ 450 MiB/s
 - 560 MiB/s peak
- 14-day retention @ 3 copies stored
- Latencies:
 - p99 - 2.5s
 - p90 - 2.3s
 - p50 - 1.6s



Grafana

Pyroscope

Pyroscope

- For **Profiling**
- 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>

Data (and cost) savings

Logs to metrics

- Full text indexing: 10 TiB logs → ~20 TiB index
- Loki: 10 TiB logs → ~200 MiB index
- Logs @ Grafana ~600 Byte 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



HACKATHON

Bringing it together

From logs to traces

The screenshot displays the Grafana Labs interface. On the left, a sidebar contains navigation icons. The main panel shows a log entry with various labels and their values. A 'Derived fields' configuration panel is overlaid on the right, allowing users to create new fields from log messages.

Log Labels:

Label	Value
pod	app-76bb7d944c-r7gb7
stdout	
traceID	1e38524b7f6e13f
_2020_11_24T15_20_29_996153289Z	
cluster	tns
container	app
level	info
namespace	tns
filename	/var/log/pods/tns_app-76bb7d944c-r7gb7_68f6413f-be42-486c-88f0-ed86badd1766/app/3.log
job	tns/app
level_extracted	info
msg	HTTP client success
duration	53.027253ms
name	app
status	500
F	
pod_template_hash	76bb7d944c
url	http://db

Derived fields

Derived fields can be used to extract new fields from the log message and create link from it's value.

Name	Regex
TraceID	(?<traceID>trace_id)=(\w+)

Query: `${_value.raw}`

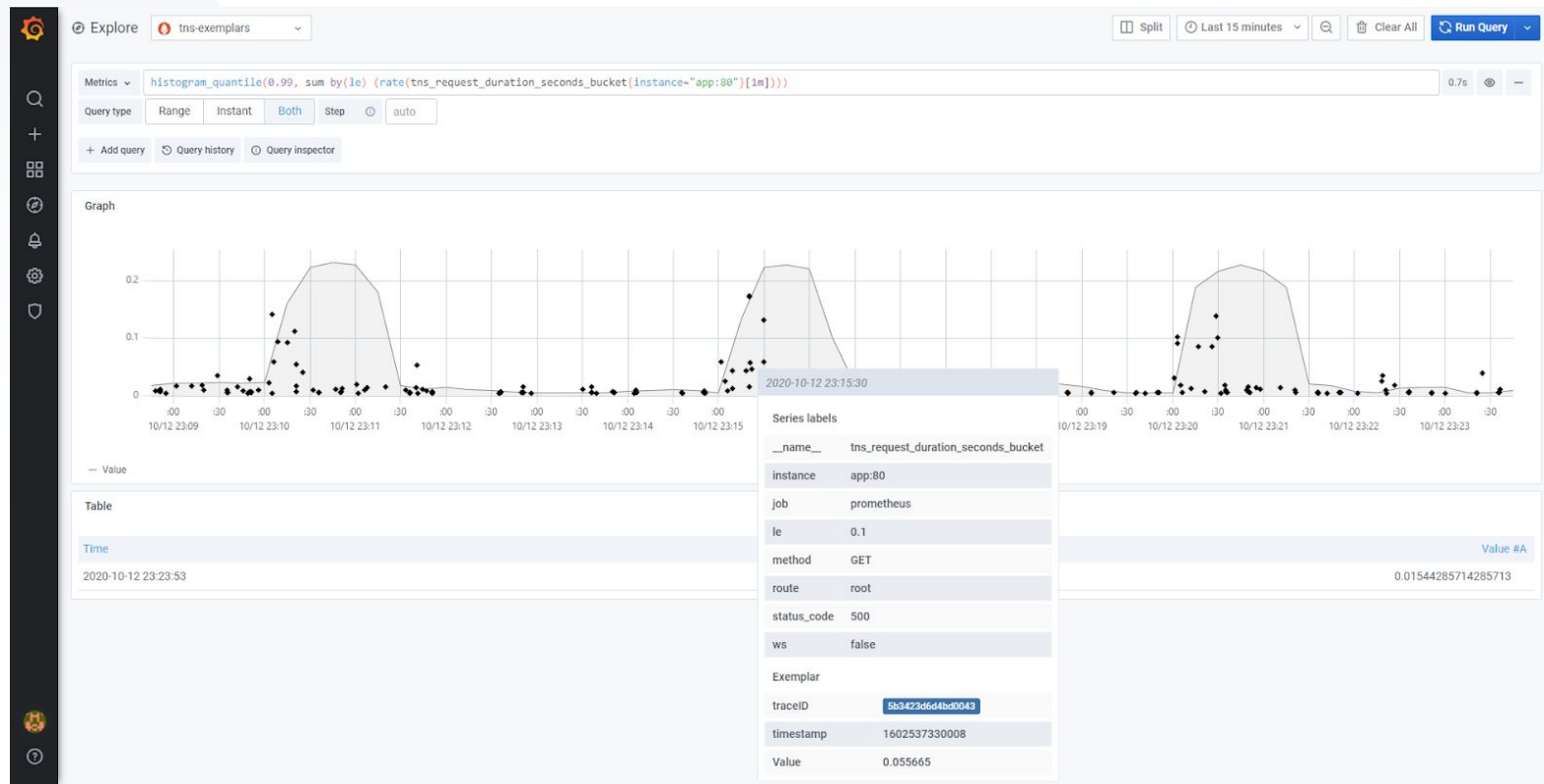
Internal link: ☒ Tempo

[+ Add](#) [Show example log message](#)

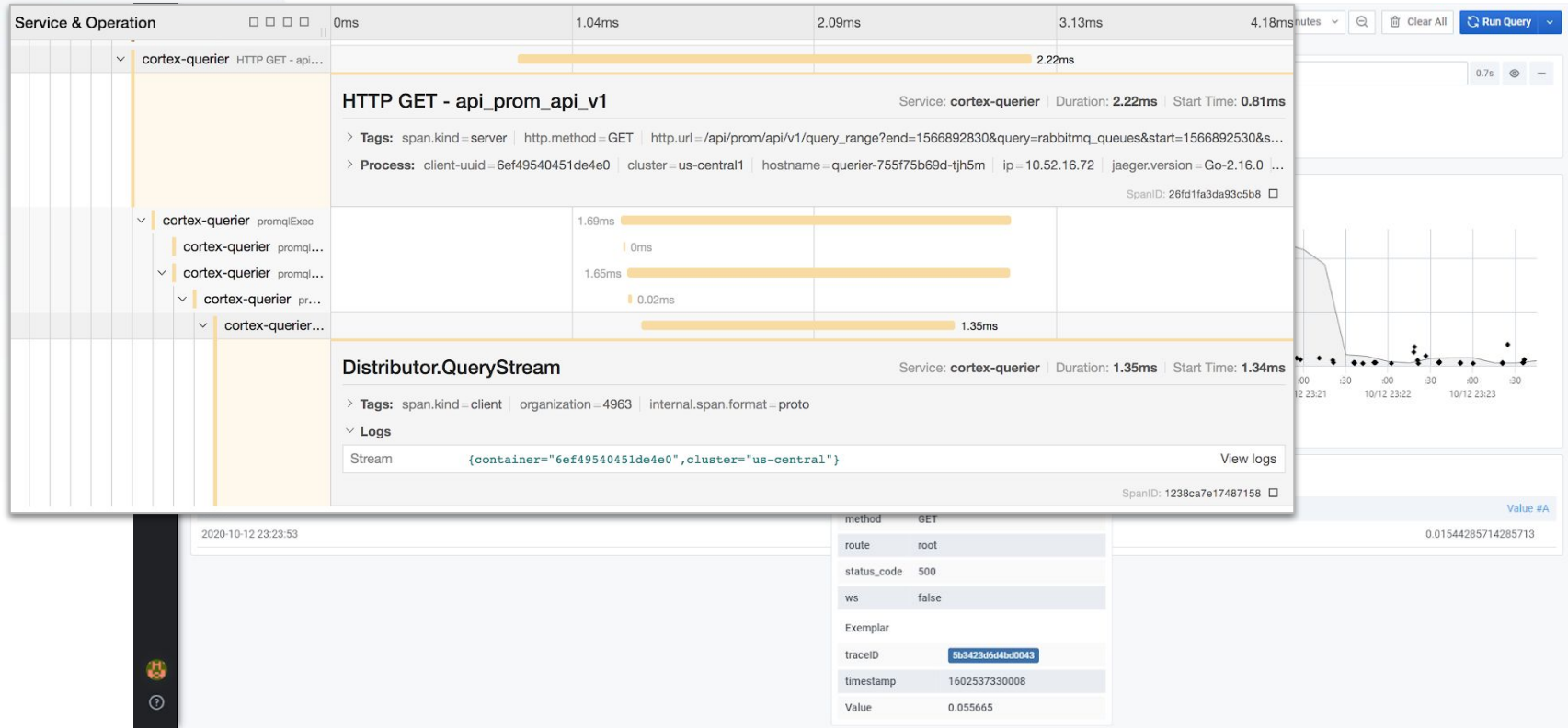
Parsed Fields:

Field	Value
TraceID	1e38524b7f6e13f Tempo
duration	53.027253ms
level	info

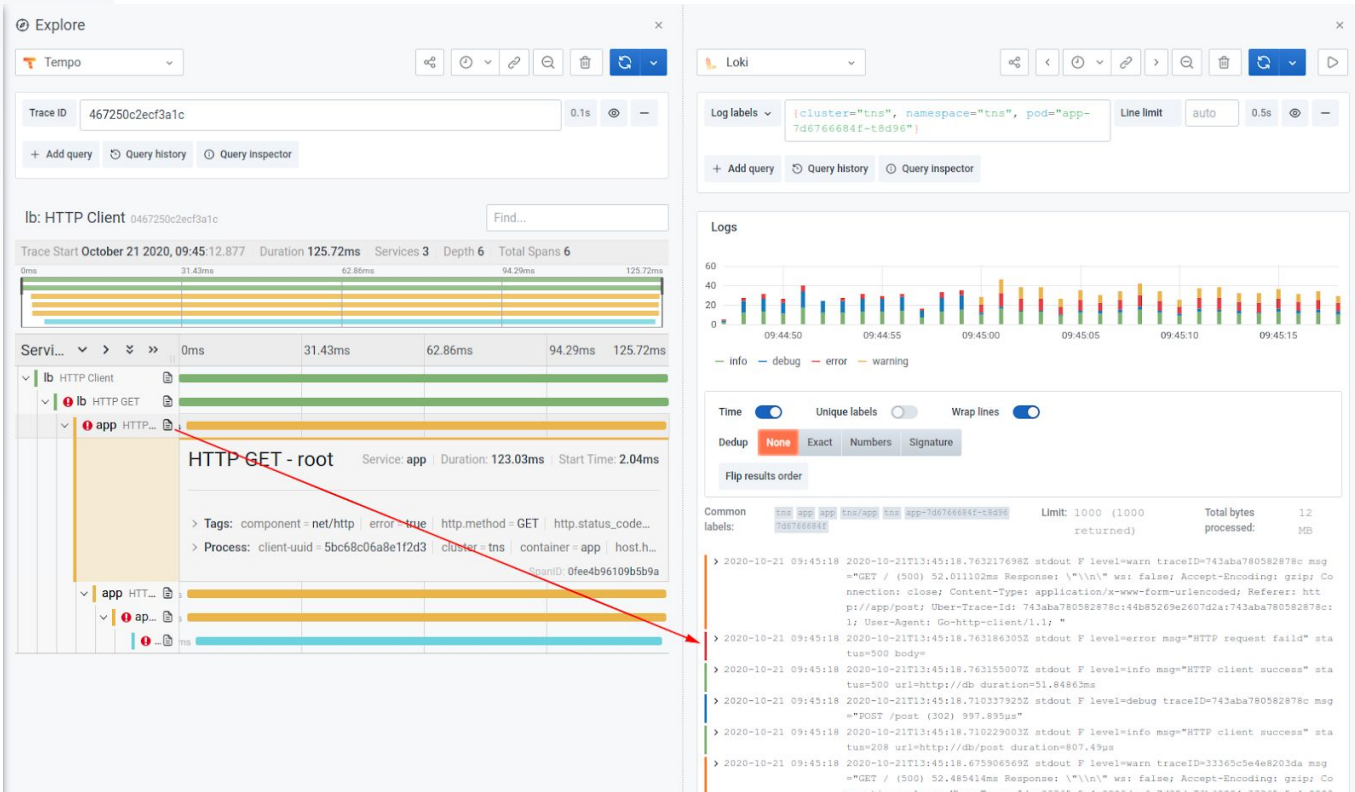
From metrics to traces

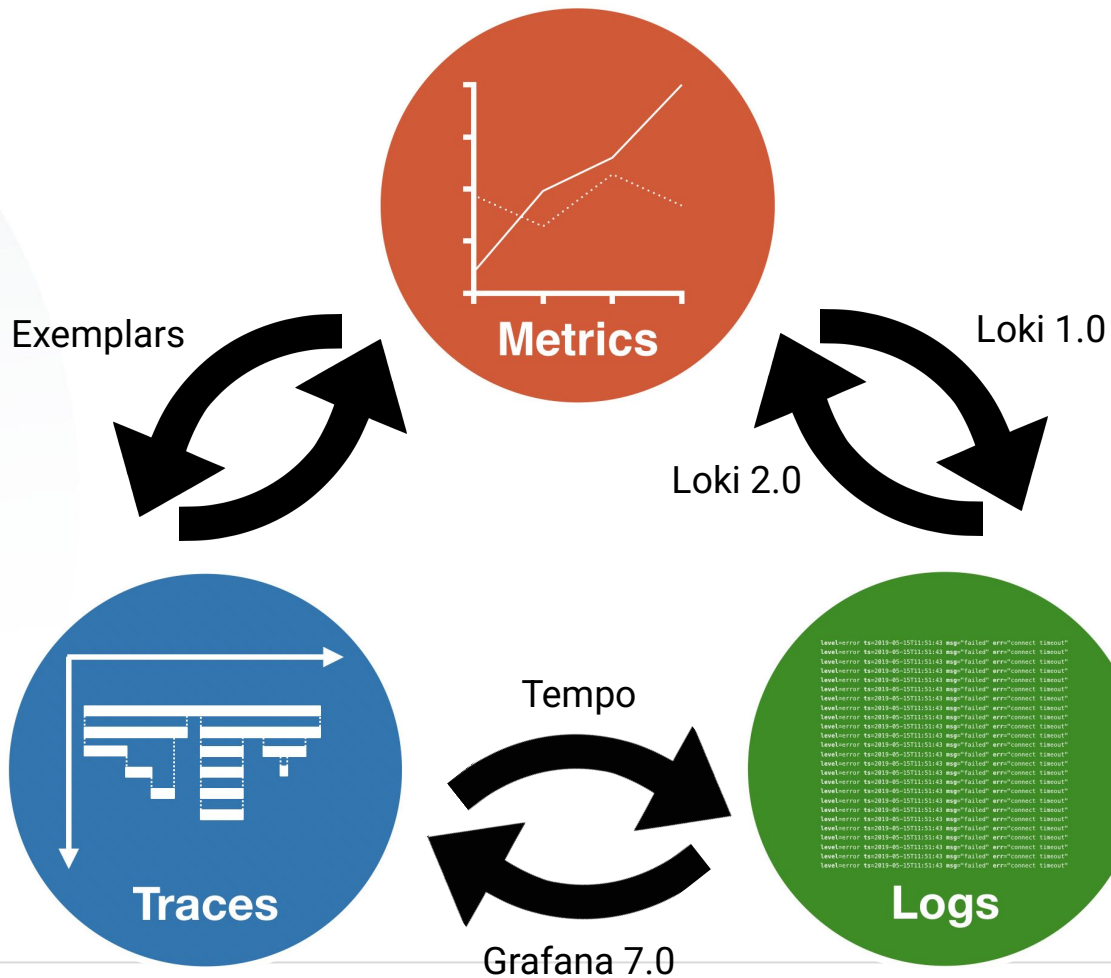


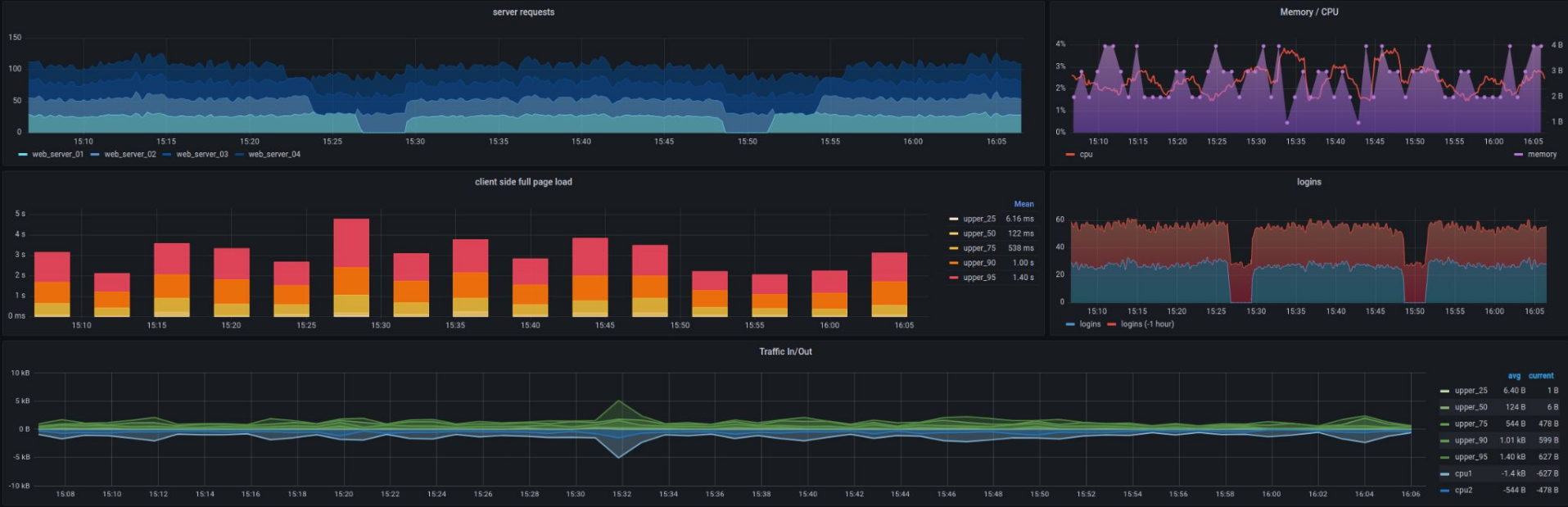
From metrics to traces

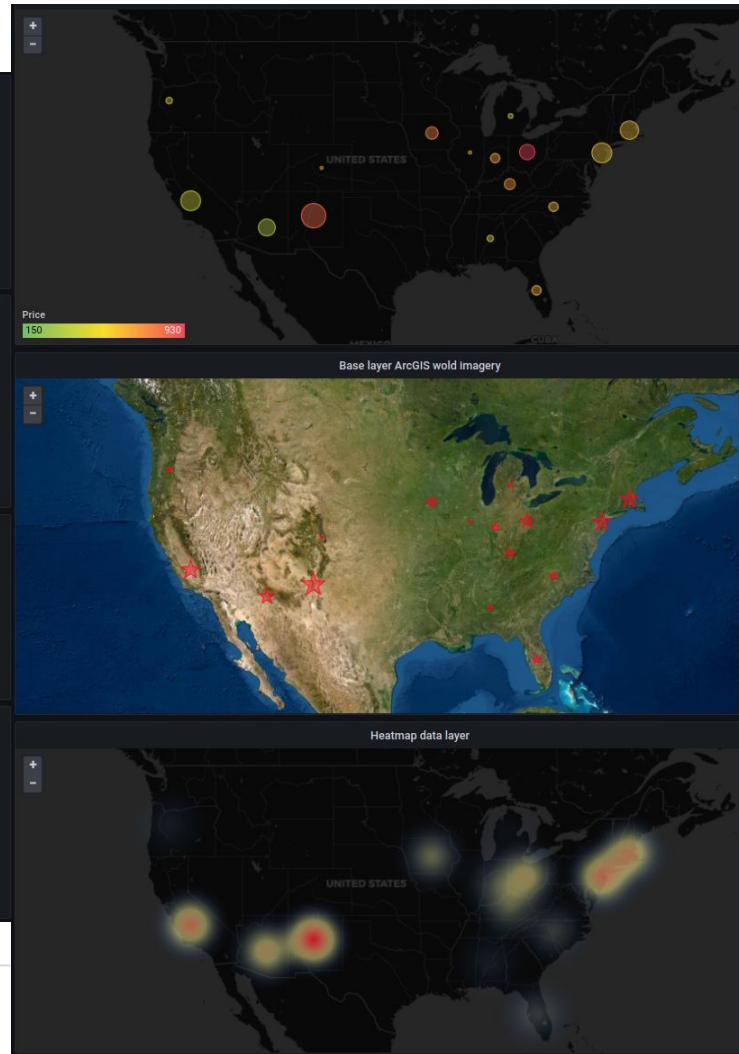


...and from traces to logs









“

All of this is Open Source and you can run it yourself
(But we will also sell it to you happily)





Thank you!

richih@richih.org

<https://chaos.social/@RichiH>

<https://github.com/RichiH/talks>