

Frontiers of Reliability Engineering

SRECon EMEA 2024 - Heinrich Hartmann

👋 I'm Heinrich - Reliability Engineer

Personal Experience

- Led Zalando SRE for 2.5 years
- Now Senior Principal SRE
- 10 years of Reliability Engineering
- Chief Data Scientist @ Circonus
- Math PhD



- A Leading Fashion Platform in Europe
- 3K Software Engineers
- 3K+ Micro services
- 250 Kubernetes Clusters
- 50M+ customers
- 14.6 bn EUR Revenue

Find me on [LinkedIn](#)

Menu

1. What have we achieved?
2. Principles
3. Where are we going?
 - a. Managing for Reliability
 - b. Mobile Observability
 - c. Data Operations



**What have
we
Achieved?**



Hardware Provisioning & Capacity Planning

2014

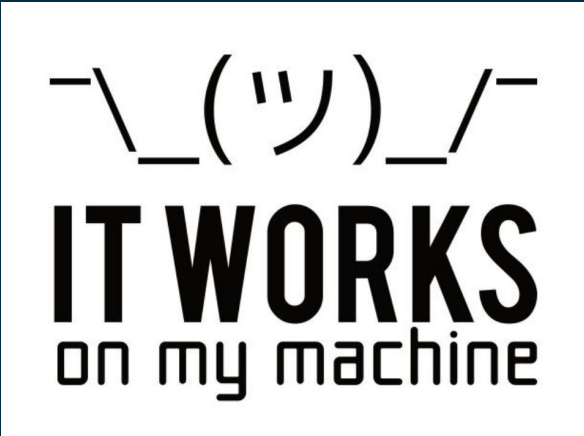


2024



Packaging and Deployment

2014



2024



Monitoring

2014



StatsD chat on gitter docker pulls 11M

Ganglia

collectd

The system statistics collection daemon



2024



Prometheus



Observability

2014

Dapper at Google

Google Technical Report dapper-2010-1, April 2010

Dapper, a Large-Scale Distributed Systems Tracing Infrastructure

Benjamin H. Sigelman, Luiz André Barroso, Mike Burrows, Pat Stephenson, Manoj Plakal, Donald Beaver, Saul Jaspán, Chandan Shanbhag

Abstract

Modern Internet services are often implemented as complex, large-scale distributed systems. These applications are constructed from collections of software modules that may be developed by different teams, perhaps in different programming languages, and could span many thousands of machines across multiple physical facilities. Tools that aid in understanding system behavior and reasoning about performance issues are invaluable in such an environment.

Here we introduce the design of Dapper, Google's production distributed systems tracing infrastructure, and describe how it is used to understand system behavior and reasoning about performance issues in such an environment.

because large collections of small servers are a particularly cost-efficient platform for Internet services workloads [4]. Understanding system behavior in this context requires observing related activities across many different programs and machines.

A web-search example will illustrate some of the challenges such a system needs to address. A front-end service may distribute a web query to many hundreds of query servers, each searching within its own piece of the index. The query may also be sent to a number of other sub-systems that may process advertisements, check spelling, or look for specialized results, including images, videos, news, and so on. Results from all

Scuba: Diving into Data at Facebook

Lior Abraham, John Allen, Oleksandr Barykin, Vinayak Borkar, Bhuwan Chopra, Ciprian Gereia, Daniel Merl, Josh Metzler, David Reiss, Subbu Subramanian, Janet L. Wiener, Okay Zed
Facebook, Inc. Menlo Park, CA

ABSTRACT

Facebook takes performance monitoring seriously. Performance issues can impact over one billion users so we track thousands of servers, hundreds of PB of daily network traffic, hundreds of daily code changes, and many other metrics. We require latency of under a minute from events occurring (a client request on a phone, a bug report filed, a code change checked in) to graphs showing those events on developers' monitors.

Scuba is the data management system Facebook uses for most real-time analysis. Scuba is a fast, scalable, distributed, in-memory database built at Facebook. It currently ingests millions of rows (events) per second and expires data at the same rate. Scuba stores data completely in memory on hundreds of servers each with 144 GB RAM. To process each query, Scuba aggregates data from all servers. Scuba processes almost a million queries per day. Scuba is used extensively for interactive, ad-hoc, analysis queries that run in under a second over live data. In addition, Scuba is the workhorse behind Facebook's code regression analysis, bug report monitoring, ads revenue monitoring, and performance debugging.

Originally, we relied on pre-aggregated graphs and a carefully managed, hand-coded, set of scripts over a MySQL database of performance data. By 2011, that solution became too rigid and slow. It could not keep up with the growing data ingestion and query rates. Other query systems within Facebook, such as Hive [2] and Pigrite [3], query data that is written to HDFS with a long (typically one day) latency before data is made available to queries and queries themselves take minutes to run.

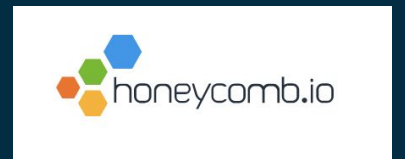
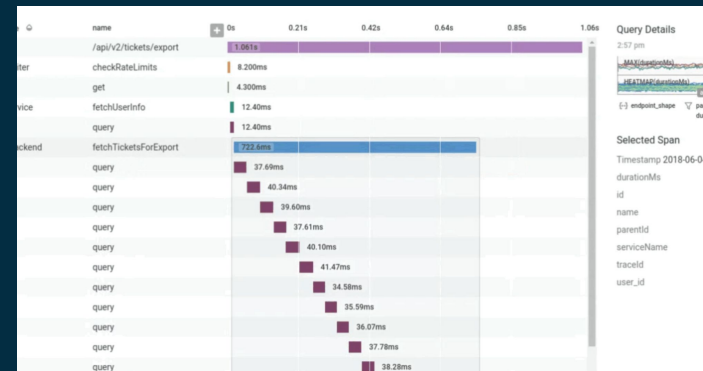
Therefore, we built Scuba, a fast, scalable, in-memory database. Scuba is a significant evolution in the way we collect and analyze data from the variety of systems that keep the site running every day. We now use Scuba for most real-time, ad-hoc analysis of arbitrary data. We compare Scuba to other data management systems later in the paper, but we know of no other system that both ingests data as fast and runs complex queries as fast as Scuba.

Today, Scuba runs on hundreds of servers each with 144 GB RAM in a shared-nothing cluster. It stores around 70 TB of compressed data for over 1000 tables in memory, distributed by partitioning each table randomly across all of the servers. Scuba ingests millions of rows per second. Since Scuba is memory-bound,

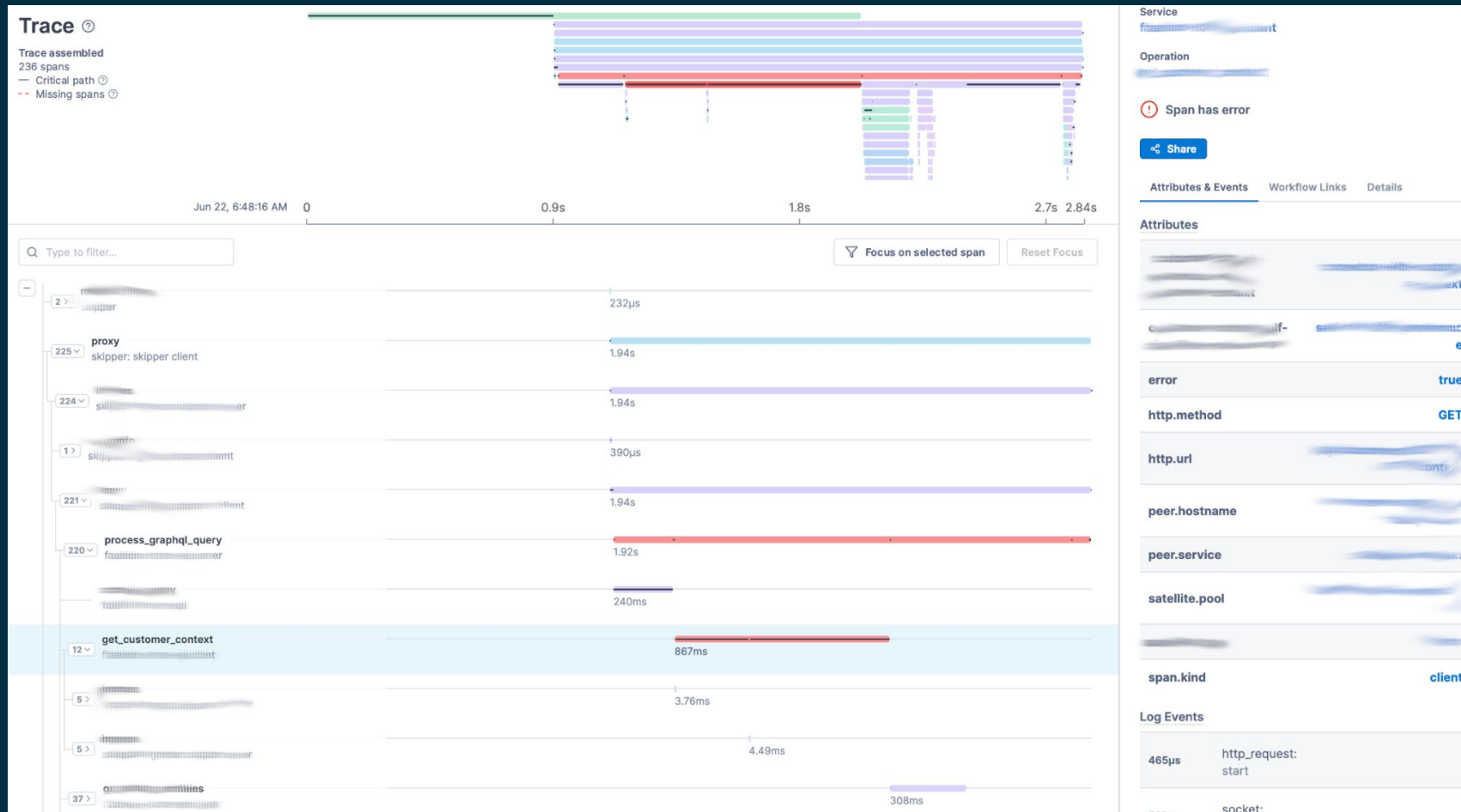
2024



Scuba at Facebook



Microservice Observability w/ Tracing



Metrics & SLOs can be derived from Tracing Data

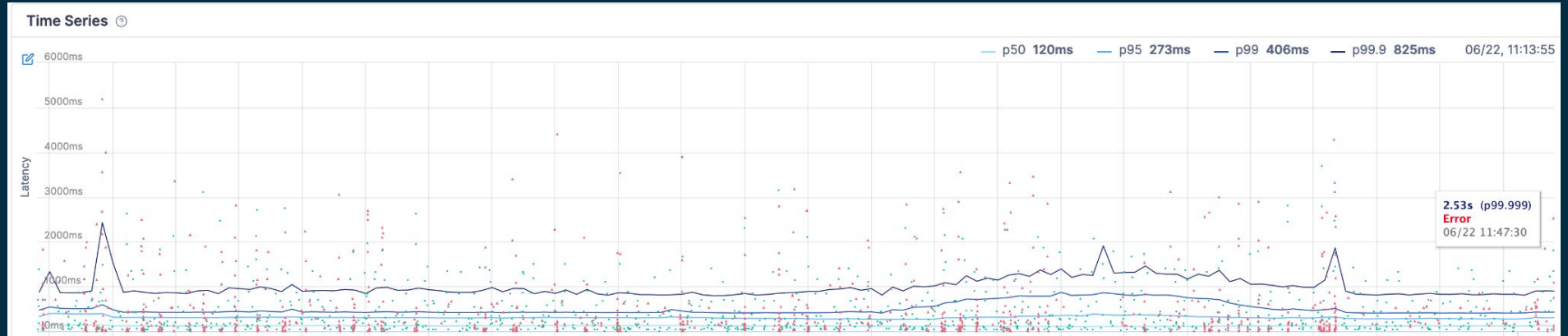
Requests



Errors



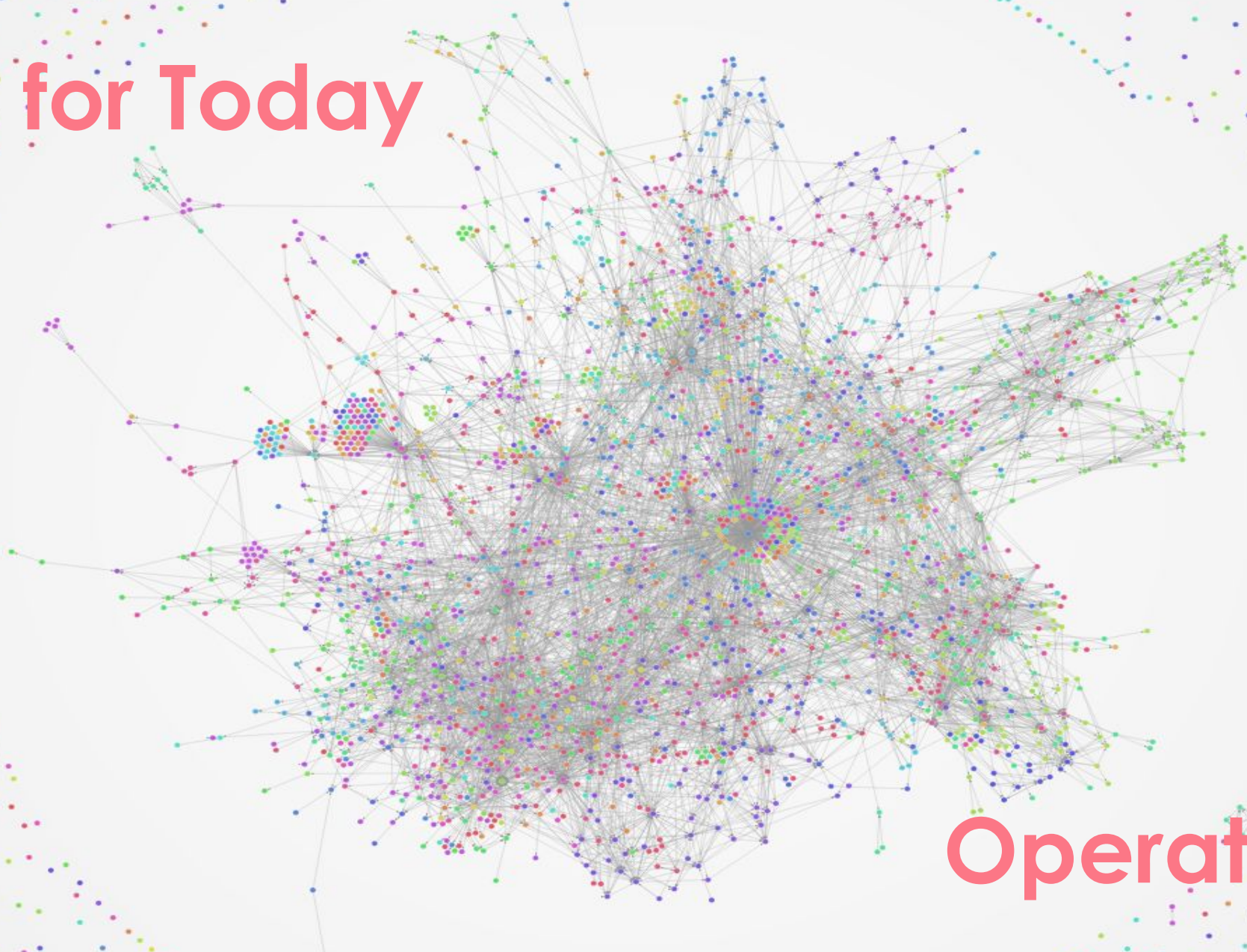
Duration



Principles



Quest for Today



Operate This!

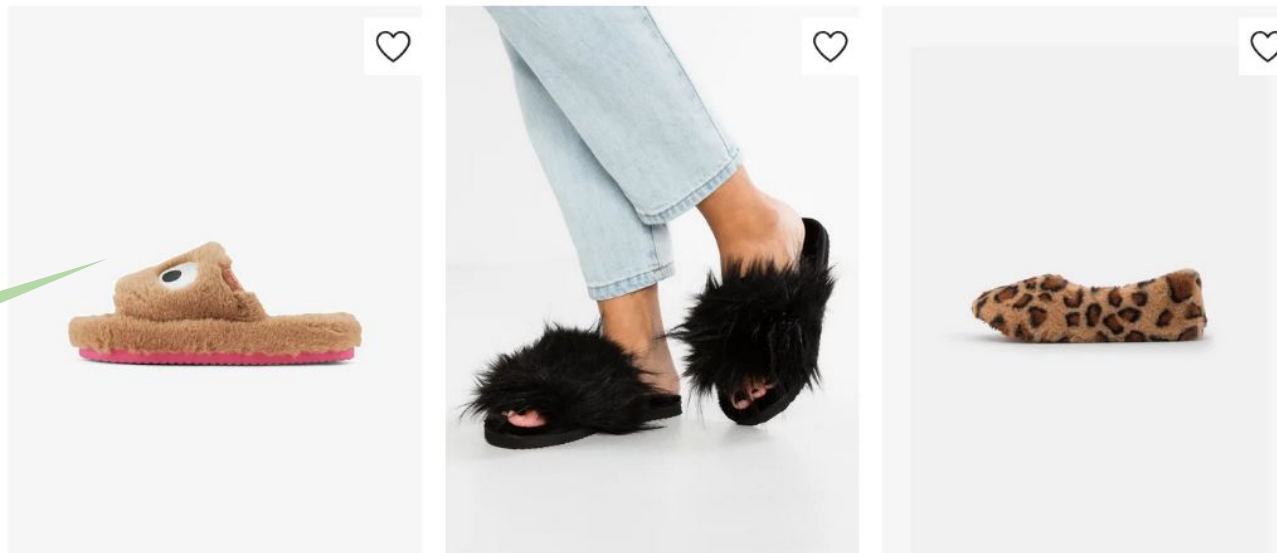
Protect the User Experience!

Shoes

- Sneakers
- Ankle boots
- Boots
- Slippers
- Flat shoes
- High heels
- Sandals
- Mules
- House shoes
- Sports shoes
- Ballerinas
- Bridal shoes
- Beach shoes

Sort by | Size | Brand | Colour | Price | Material | Pattern | Type of heel | Toe | Fastener | Lining | Show all filters


61 items ⓘ



Browse Catalog

View Product Detail

View Cart



I don't care
if your
Data Center
is on fire.

This is
fine!

“

**Reliability Engineering involves
People as much as Technology.**

Engineering Reliability at Scale

Small Company (~10 FTE)

- Alerts & Dashboards
- Logging
- On-call rotations

Medium Company (~100 FTE)

- Playbooks
- Incident Management
- Observability

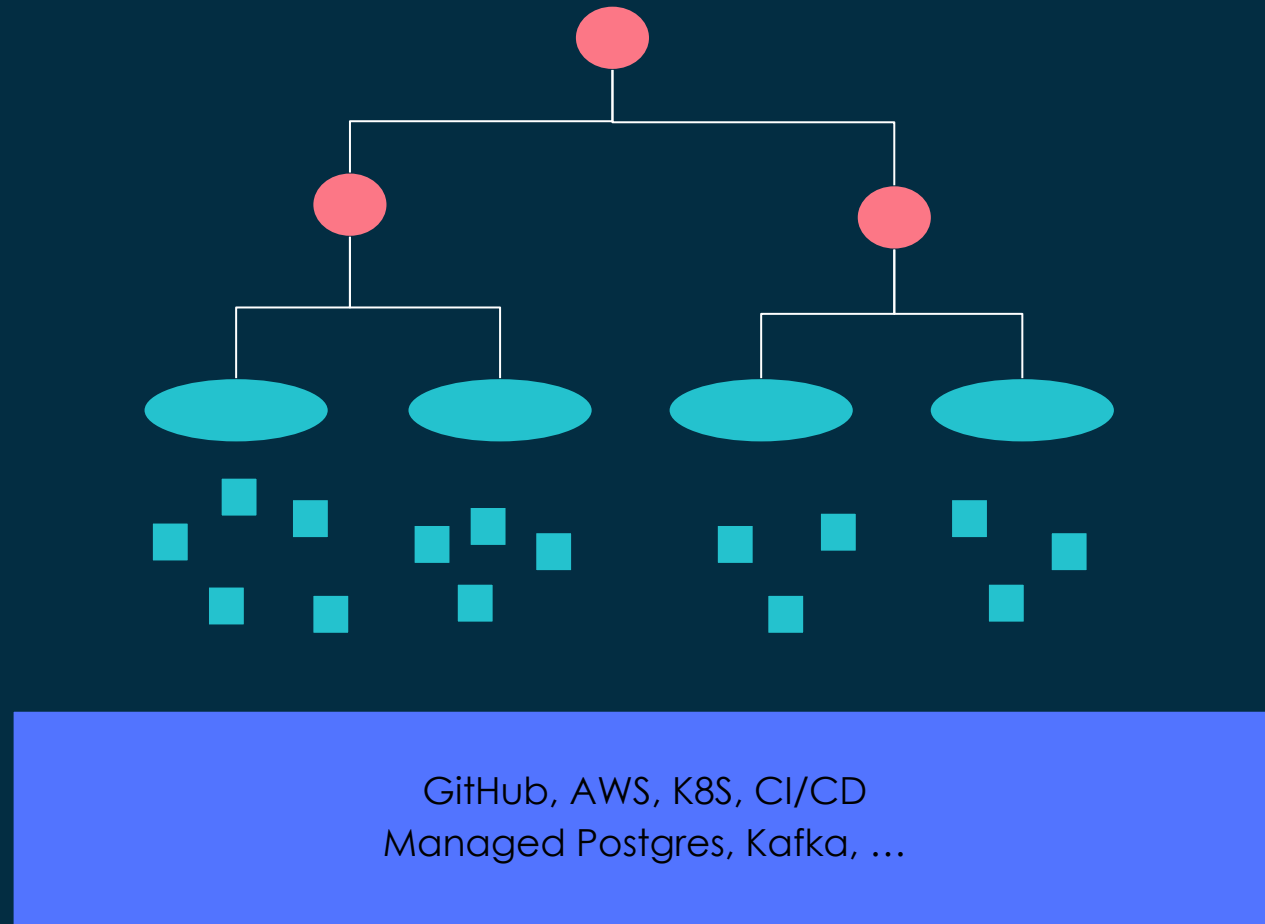
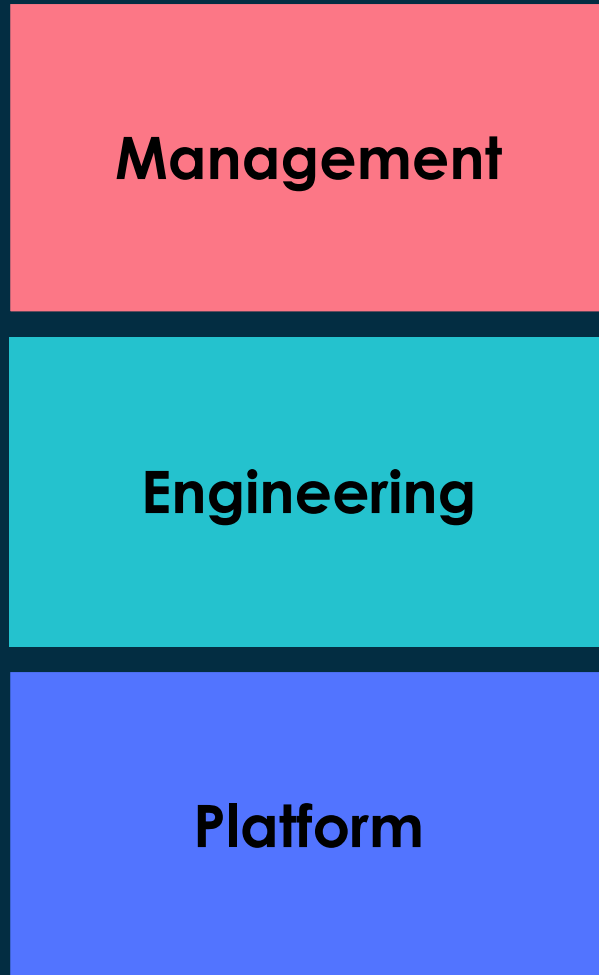
Large Company (>1k FTE)

- WORM Meetings
- Risk Management
- SRE Community

People Problems

Technical Problems

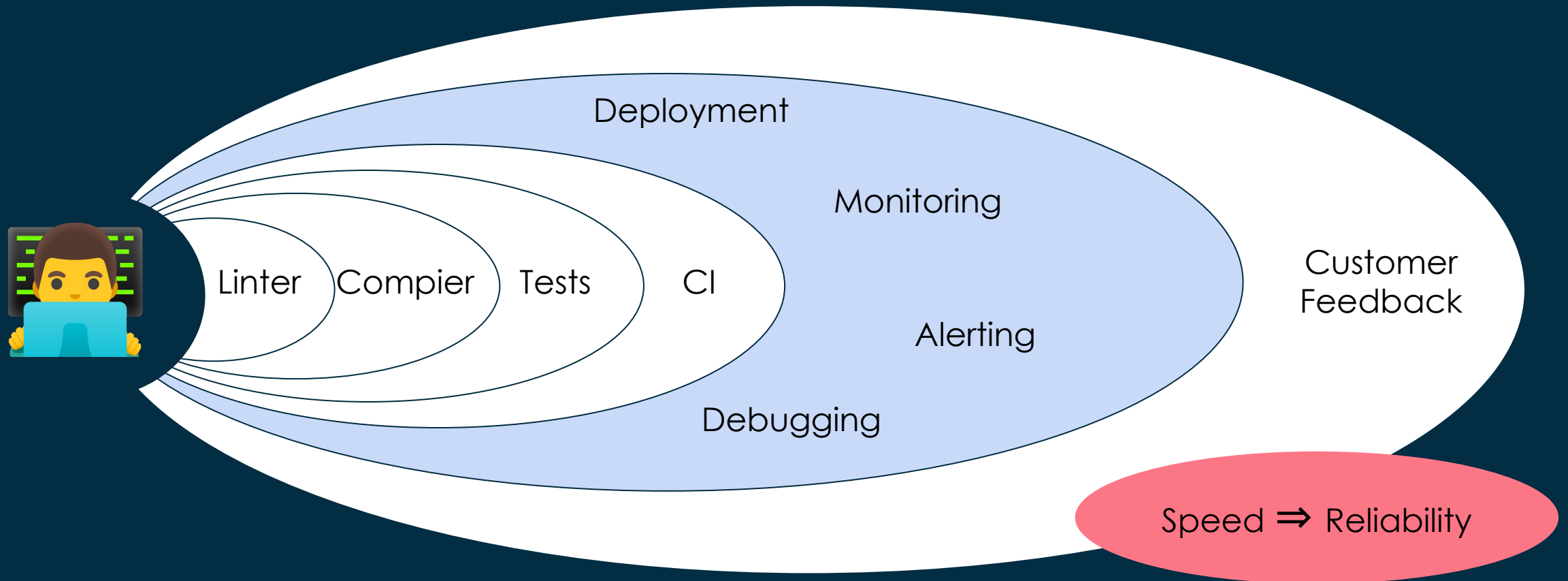
Layerd Systems Model of Zalando



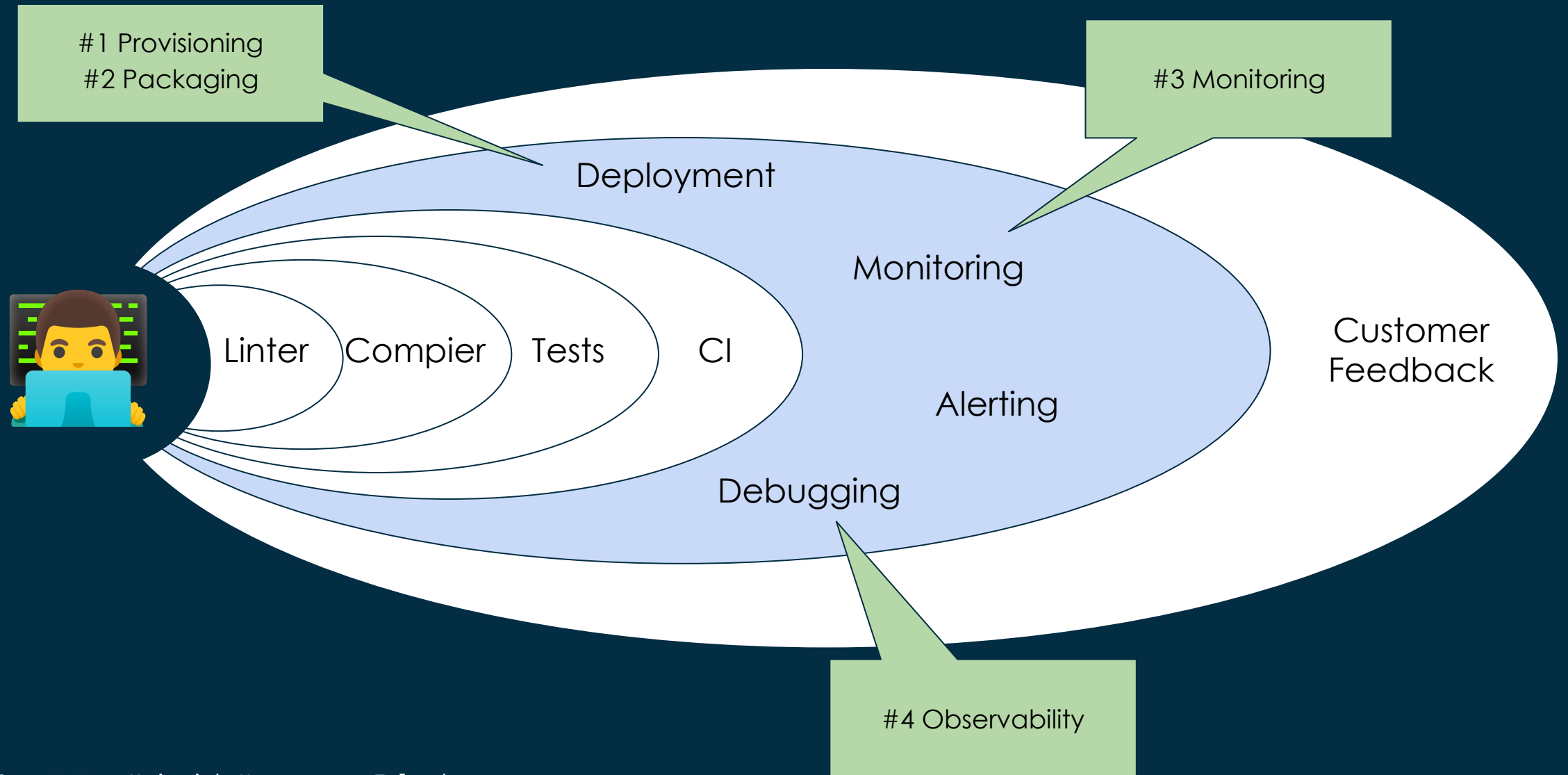
“

**Reliability Engineering is
all about Feedback Loops.**

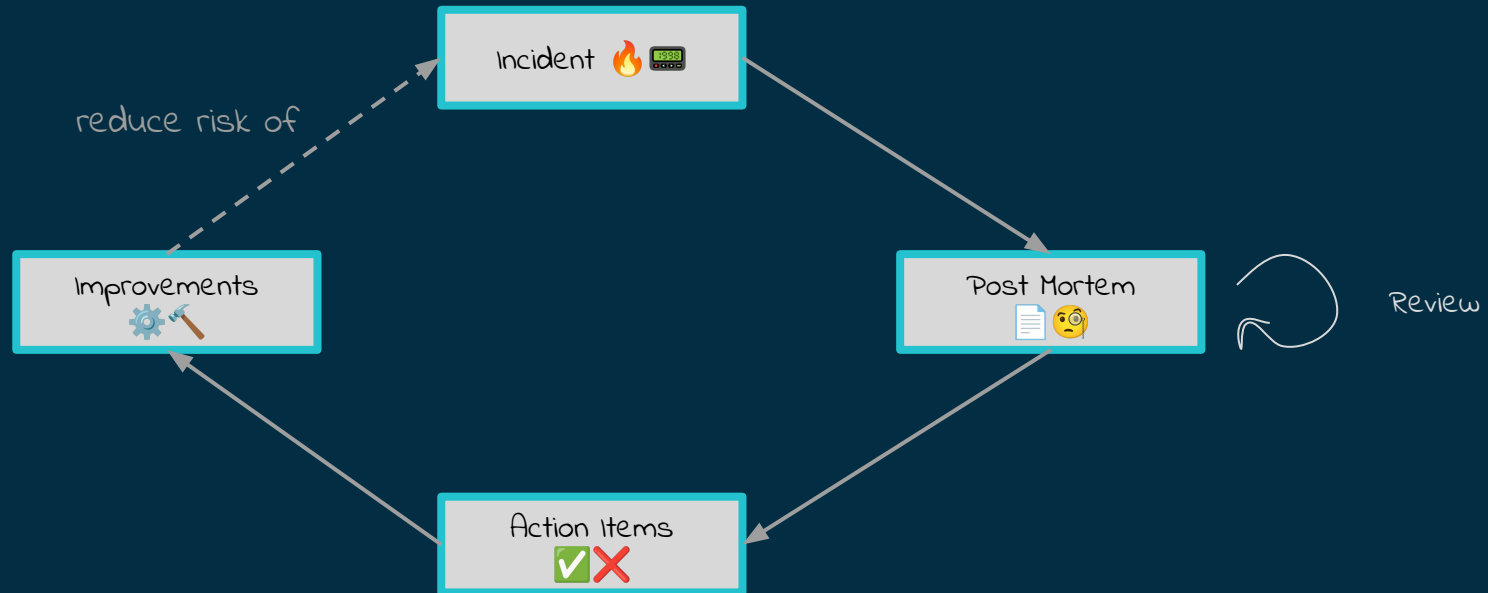
Reliability is driven by Feedback Loops



Achievements Accelerate Feedback Loop



The Incident Process is a Meta Feedback Loop



Where are
we going?



Managing for Reliability



How to enable on Management to steer for reliability?

“

You get what you inspect.

Reliability Reports

for Management on all Levels

Auto-generated Google Doc supporting
Weekly Operational Review Meetings ("WORM").

Agenda

- Incidents
- SLO
- On-Call Health
- Open Post Mortems

Reliability Report CW16-2024

Site Reliability Engineering

2024-04-15 to 2024-04-21

Overview

Between 2024-04-15 and 2024-04-21 the following reliability metrics were observed:

Incidents	Breached 7 day SLOs	GMV Loss	On Call Paging Alerts
2 (+0 ▶)	- (-1 ▼)	0 € (+0 € ▶)	4 (-11 ▼)
Last week: 2	Last week: 1	Last week: 0 €	Last week: 15

SLOs

Critical Business Operation	SLO	SLI (28 days)	SLI (7 days)	Error budget (28 days)	Notes
Configure ZMON	99.900%	99.998%	99.997%	98.09% ✓	
Log freshness	99.900%	99.907%	99.934%	7.38% ⚠	Log shipping on Kubernetes master nodes was impacted from Apr 11 to Apr 15, as discussed in previous weeks WORM. SLI has been recovering since the end of the incident.
Log freshness Test Cluster	99.500%	99.721%	99.840%	44.22% ✓	
Metric freshness	99.900%	99.927%	99.924%	27.8% ✓	
Metric freshness Test Clusters	99.500%	99.869%	99.892%	74.0% ✓	
Notify anomaly	99.900%	! 99.893%	99.980%	0% !	SLO breach occurred on 3rd of April. SLO should recover on 2nd of May.
Notify failure	99.990%	100.000%	100.000%	100.0% ✓	
Trace freshness	99.900%	99.994%	99.991%	94.93% ✓	
Trace freshness Test Cluster	99.500%	99.793%	99.793%	58.74% ✓	

Current open Post-Mortems

Post-Mortem	Severity	Repaired at	Team	Open since (working days)	Take Action
False-positive alert in ZMON	SEV3	Sat, 2024-04-13	Observability	5	Follow-up on review comments
Zalando eCommerce Platform - Log freshness - Scalyr (production clusters) - error ratio > 0.56% over the last 6h and 30m	SEV3	Mon, 2024-04-15	Observability	5	Review post-mortem and define additional follow-up action items (circuit breaker for S3 log shipping)
grafana.zalando.net is unreachable	SEV3	Tue, 2024-04-16	Observability	4	Complete post-mortem

Incident Table shows where we failed the User

Title	Metadata	Impact	Take action
Zalando eCommerce Platform Engineering / Builder Infrastructure			
website unavailable	<p>Owning Team: Cloud Governance</p> <p>Detected at: Wed, 2024-10-02 09:27</p> <p>TTD: 1h 40m</p> <p>TTR: 35m</p>	<p>Severity: SEV3</p> <p>Categories: Third-Party</p> <p>Impact: GMV Loss: 0 EBIT Loss: 0 Affected Customers: 0 Affected Partners: 0 Affected Employees: 0 Involved Applications: - b2</p>	<p>Effect: still unclear Cause: Unknown, the website is managed by an external agency Action: Clarify ownership of the website</p>
data-governance availability zone b in crashloop backoff	<p>Owning Team: Observability</p> <p>Detected at: Sun, 2024-10-06 08:28</p> <p>TTD: 9m</p> <p>TTR: 22m</p>	<p>Severity: SEV3</p> <p>Categories:</p> <p>Impact: GMV Loss: - EBIT Loss: - Affected Customers: Affected Partners: Affected Employees: Involved Applications:</p>	<p>Effect: crashloop backoff loading sequentially no space one. Cause: investigation ongoing Action: Investigation of the incident is ongoing and the root cause is still under investigation.</p>
ZMON backend 'zmon.zalando.net' slow/unavailable leading to stale entities and false positive alerts	<p>Owning Team: Observability</p> <p>Detected at: Wed, 2024-10-02 12:55</p> <p>TTD: 16m</p> <p>TTR: 47m</p>	<p>Severity: SEV3</p> <p>Categories:</p> <p>Impact: GMV Loss: - EBIT Loss: - Affected Customers: 0 Affected Partners: 0 Affected Employees: Involved Applications: - zmon</p>	<p>Effect: Creation of false positive alerts as well as delayed metrics ingestion. Cause: The incident was triggered by the traffic switch of ZMON entity management requests from the ZMON Controller to a new dedicated entity management service. The root cause for the observed effects is still under investigation. Action: Traffic was switched back to ZMON Controller. Additional testing is being conducted to understand why the traffic switch caused unexpected effects (sudden increase in total entity count) that were not observed before, when the new service was running in shadow mode.</p>

Teams report on

1. Impact
2. Cause
3. Actions

SLOs provides top-down view on Reliability

Critical Business Operation	SLO	SLI (28 days)	SLI (7 days)	Error budget (28 days)	Notes
Builder Infrastructure					
Configure ZMON	99.900%	99.993%	99.992%	93.99% ✓	
Log freshness	99.900%	99.904%	99.966%	4.51% ⚠	
Log freshness Test Cluster	99.500%	99.853%	99.921%	70.64% ✓	
Metric freshness	99.900%	! 99.752%	! 99.665%	0% !	
Metric freshness Test Clusters	99.500%	99.523%	! 99.441%	4.75% ⚠	
Notify anomaly	99.900%	99.986%	99.968%	86.64% ✓	
Notify failure	99.990%	100.000%	100.000%	100.0% ✓	
Trace freshness	99.900%	99.990%	99.983%	90.82% ✓	
Trace freshness Test Cluster	99.500%	99.992%	99.984%	98.47% ✓	
Write to Nakadi	99.990%	99.999%	99.999%	94.21% ✓	
Customer Domain					
SSO Login	99.950%	99.999%	99.999%	93% ✓	
SSO Registration	99.950%	99.999%	99.999%	86% ✓	
SSO Logout	99.950%	99.999%	99.999%	93% ✓	
SSO Step up authentication	99.950%	99.999%	99.999%	93% ✓	
SSO Step up authentication	99.950%	99.999%	99.999%	93% ✓	
Demand / Home & Content Visibility					
	99.000%	99.999%	99.999%	93% ✓	
	99.000%	99.999%	99.999%	93% ✓	

Teams report on

1. Impact
2. Cause
3. Actions



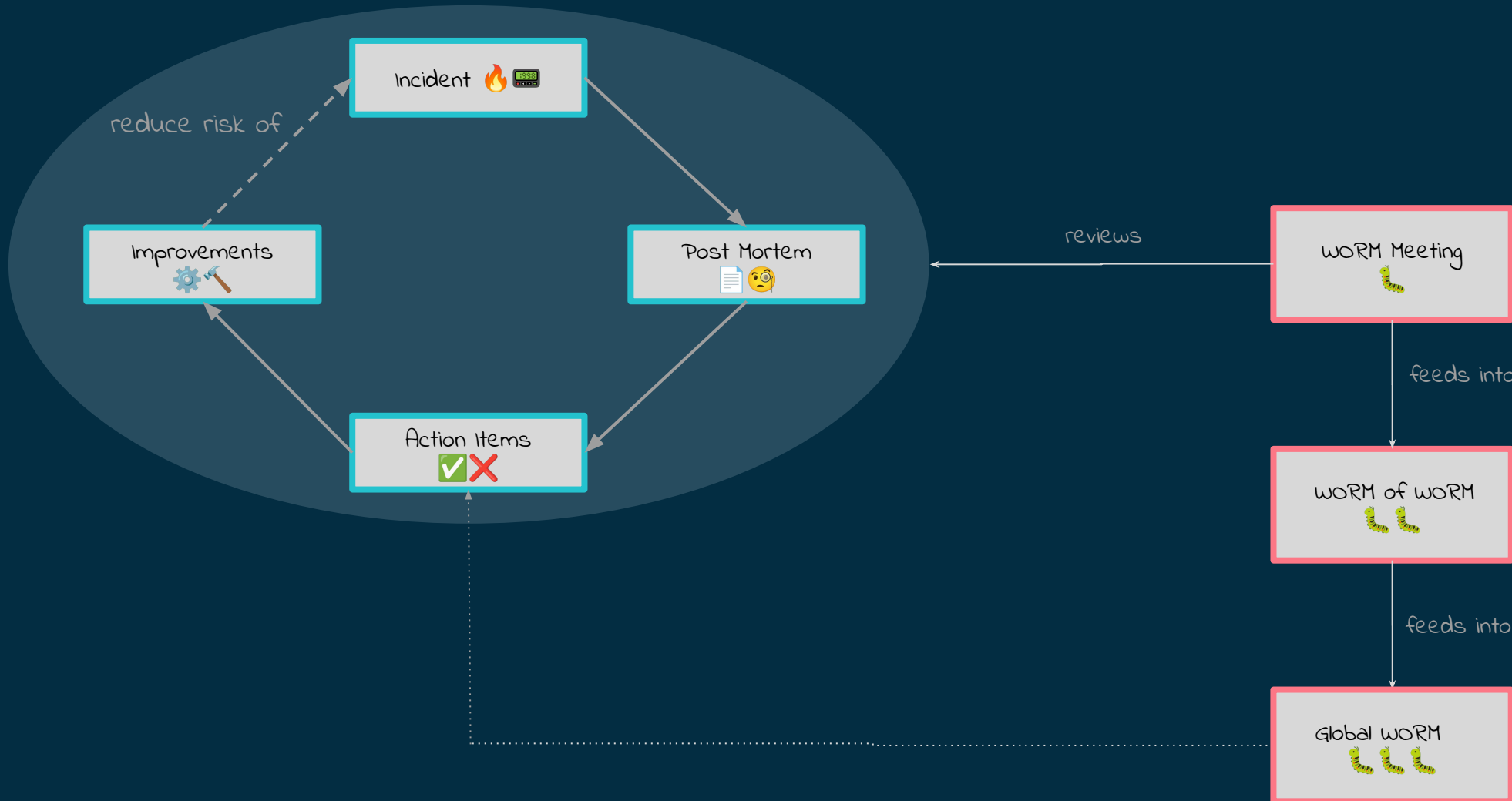
We check Alerting loads for All On-Call teams

On-Call Team	Paging alerts / day							Paging alerts				Context
	Mon 20	Tue 21	Wed 22	Thu 23	Fri 24	Sat 25	Sun 26	within working hours	Off hours	Total	Average	
[Redacted]	-	1	2	-	2	2	1	5 (-2 ▼)	3 (+2 ▲)	8 (+0 ►)	1.14 / day	[Redacted] consciously
[Redacted]	-	-	-	5	-	7	-	1 (+0 ►)	11 (+11 ▲)	12 (+11 ▲)	1.71 / day	One legitimate alert, the positives (e.g. cause resilience
[Redacted]	-	2	1	2	-	-	-	4 (-7 ▼)	1 (-2 ▼)	5 (-9 ▼)	0.71 / day	
[Redacted]	-	-	-	1	-	-	-	1 (+1 ▲)	- (+0 ►)	1 (+1 ▲)	0.14 / day	
[Redacted]	-	-	-	1	-	-	-	- (-3 ▼)	1 (+0 ►)	1 (-3 ▼)	0.14 / day	
[Redacted]	-	-	-	-	-	-	-	- (+0 ►)	- (+0 ►)	- (+0 ►)	0.0 / day	
[Redacted]	-	-	-	-	-	-	-	- (+0 ►)	- (+0 ►)	- (+0 ►)	0.0 / day	
[Redacted]	-	15	1	-	-	-	-	1 (+0 ►)	15 (-58 ▼)	16 (-58 ▼)	2.29 / day	[Redacted] to the

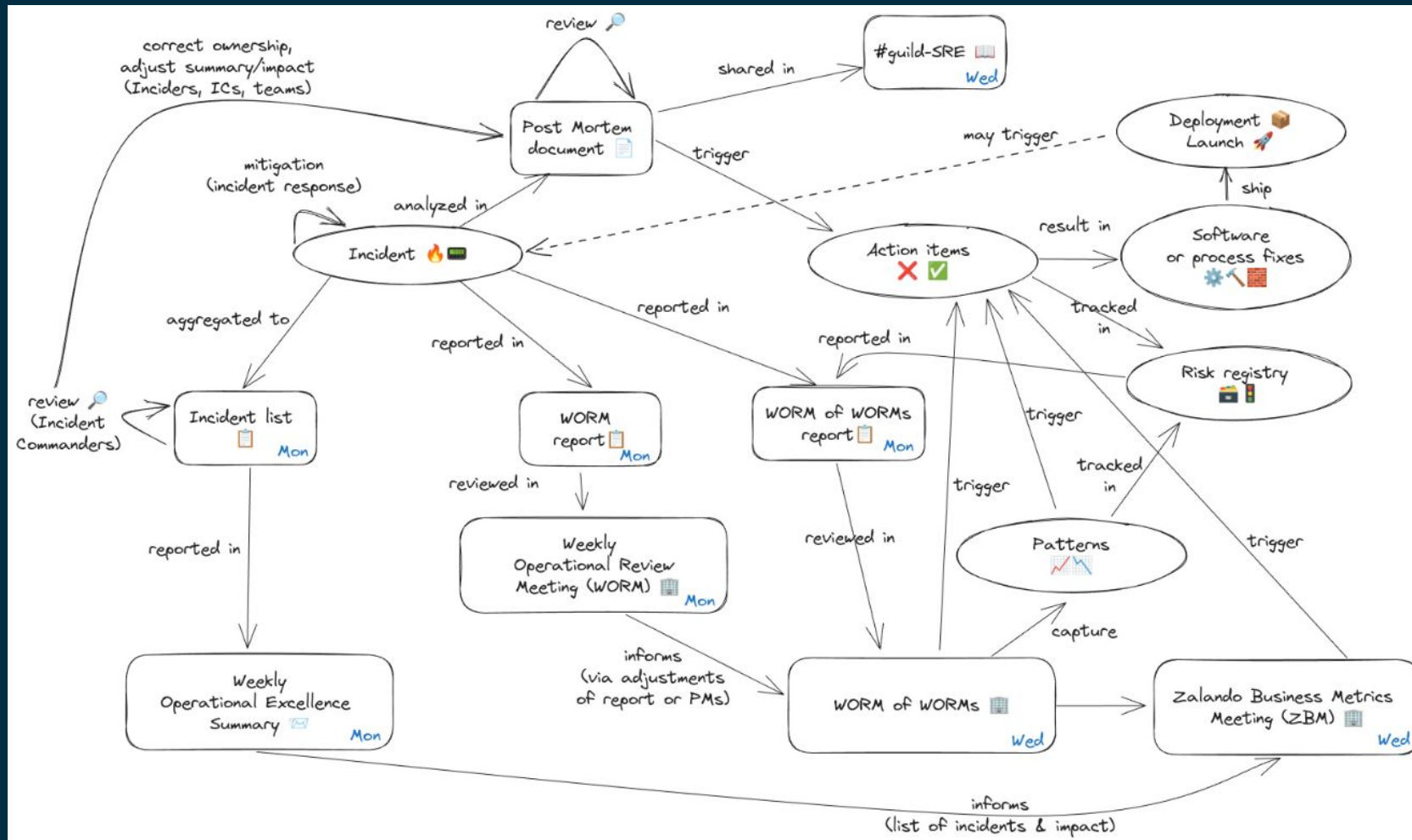
What caused the alerts?

What are we doing to prevent this?

Weekly Operational Review Meeting (WORM) Cascade

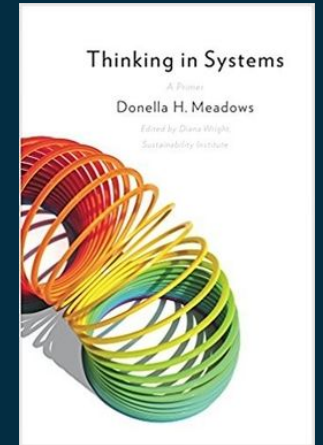


Further Augmentations of the Incident Process



The Quest Continues

- How to steer Management Attention on Reliability?
 - Expand coverage of Reliability Reports
 - Apply “Systems Engineering” to Reliability
- How to drive Cross-Organisational Reliability Initiatives?
 - SRE/Champion Model
- How to increase value from Incident Process?
 - Post Motem AI Capabilities



Mobile Observability



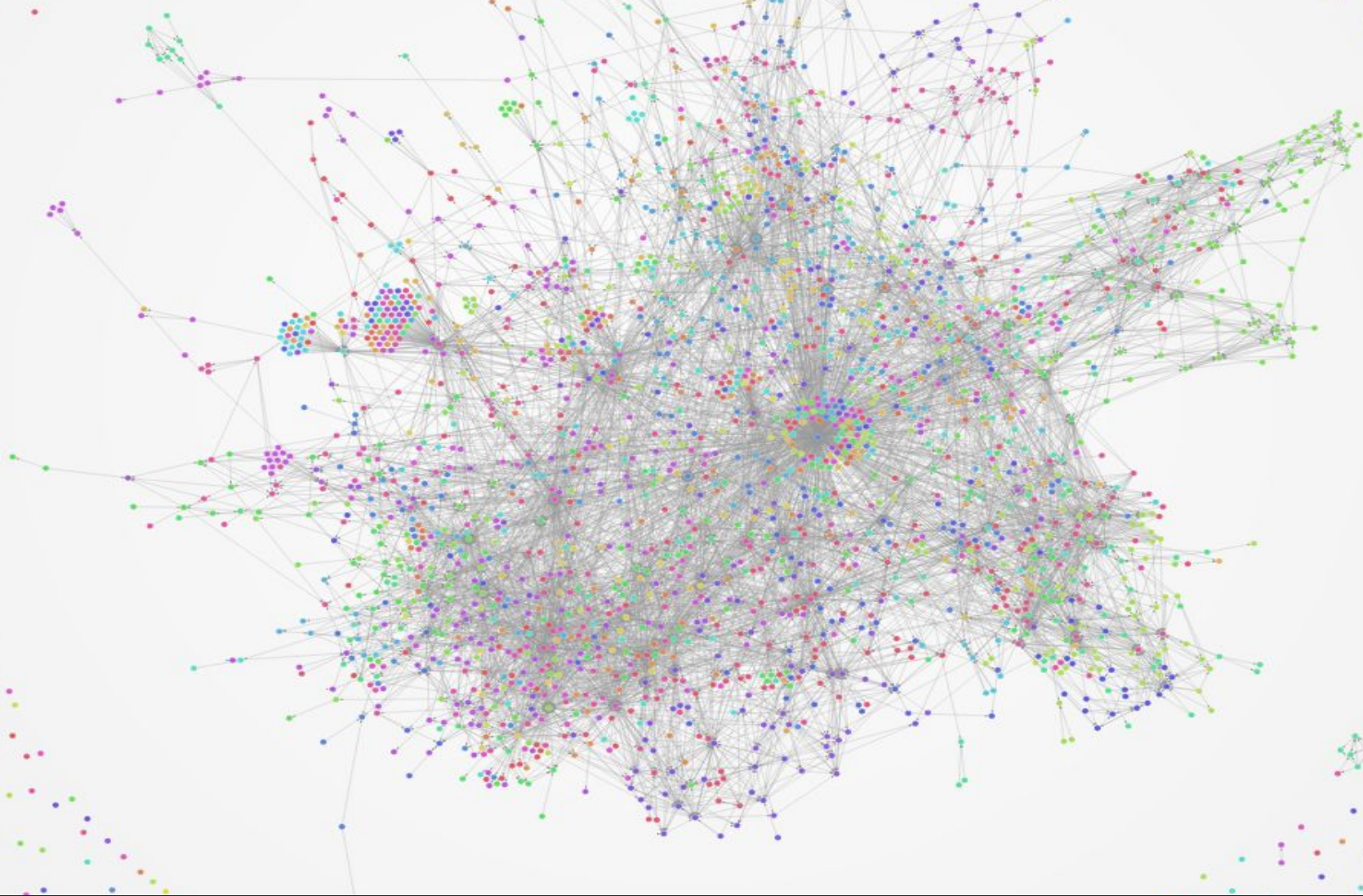
Incident Story

The undetected Order Drop

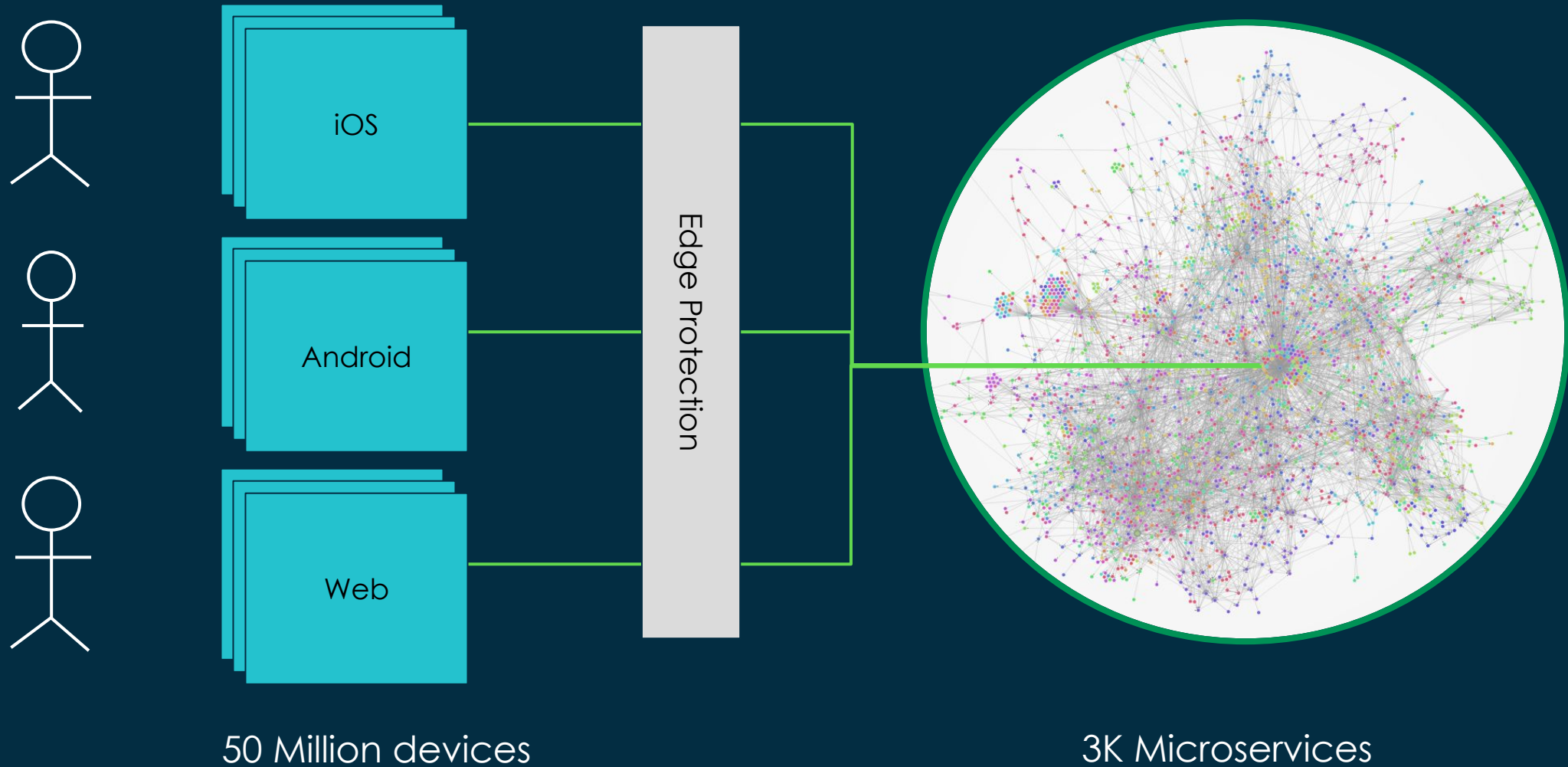
Incident Story

The lurking Add-To-Cart Failure

Is this even the right system to look at?

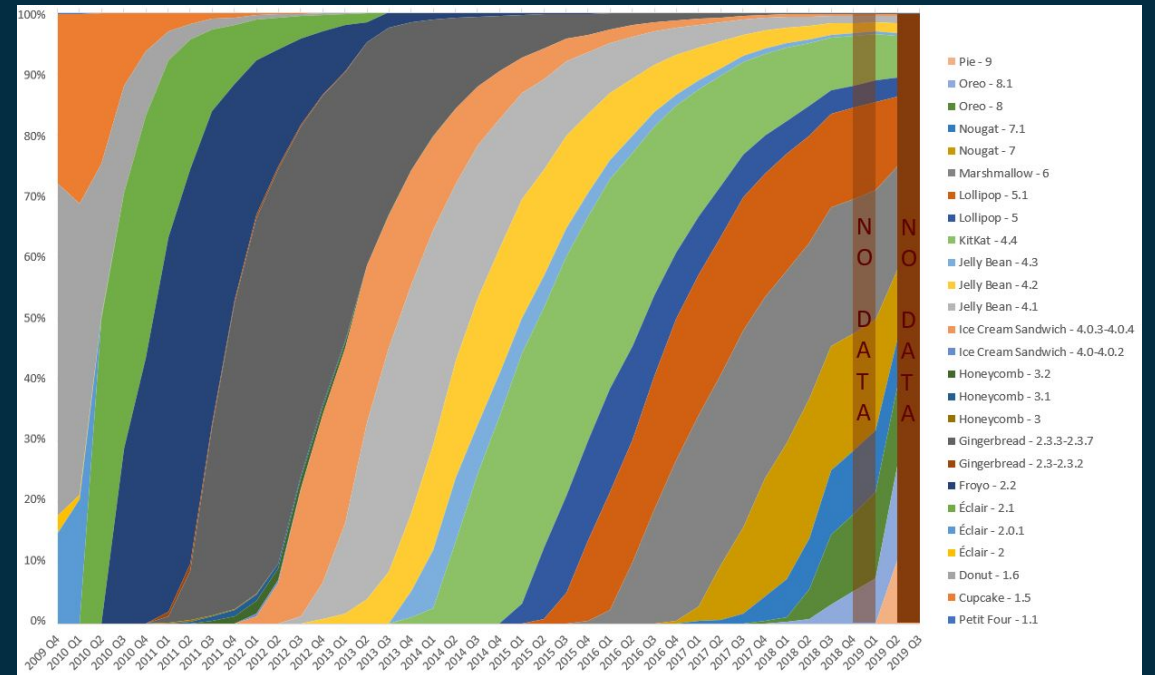


To Protect UX, we have to look at the full system!



Challenges in Mobile Environment

1. **GLACIAL** deployment speed
 - One release every 2 weeks
 - 4 weeks to reach 80% penetration
 2. Fragmented Platforms
 3. Cellular Networks
 4. Legal constraints
 5. Automated UI Testing is HARD (i.e. missing)
 6. Available Telemetry Data very limited
- => No DevOps culture in Mobile teams



Android versions deployed in the field (src: [wikipedia](https://en.wikipedia.org/wiki/Android_version_history))

Where are we investing?

* Distributed Tracing on Mobile + Web Clients

- Observability SDKs for Client Platforms

- Browser (done)
- Mobile (WIP) iOS / Android



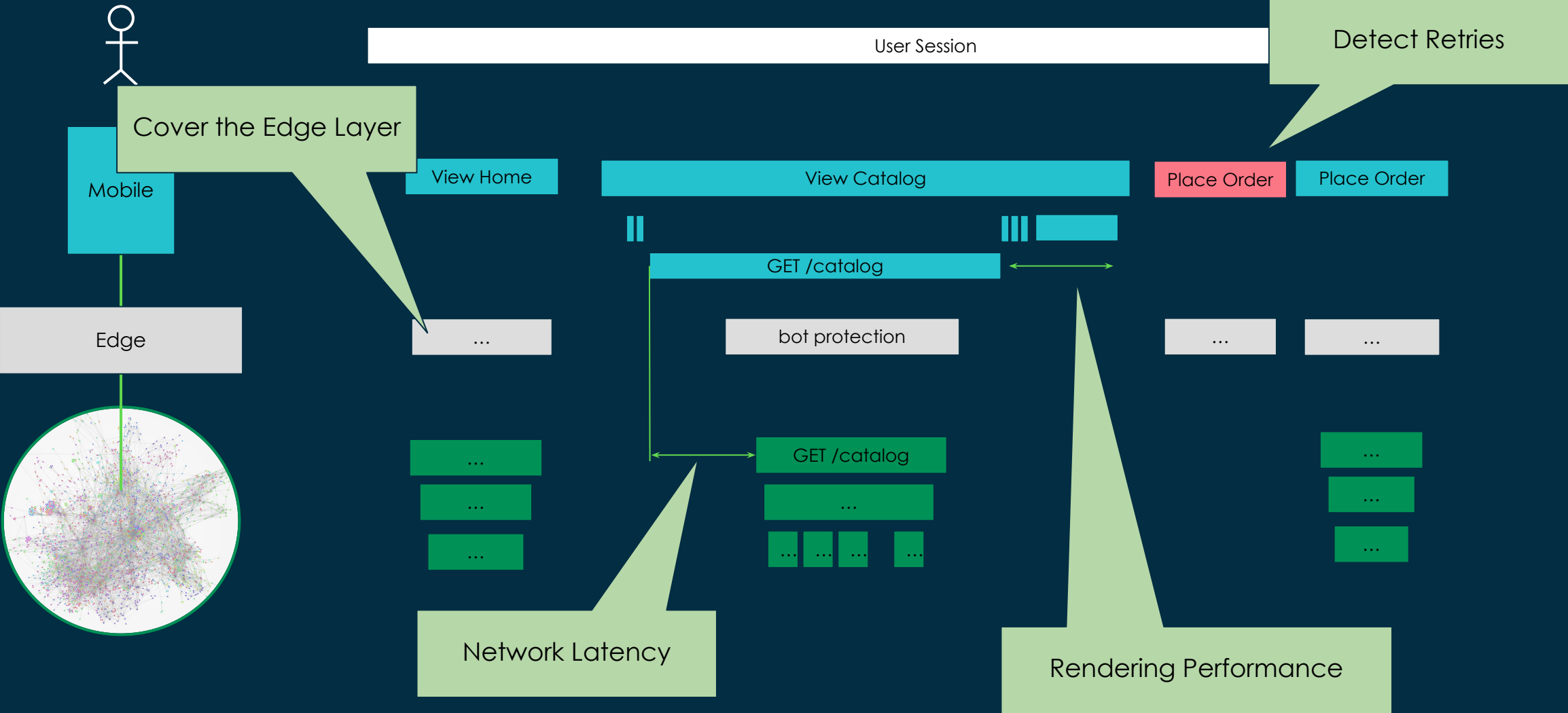
- Client Side SLOs

... complementing server-side measurements.

Goal: Expand Distributed Tracing to the Client!



Goal: Expand Distributed Tracing to the Client!

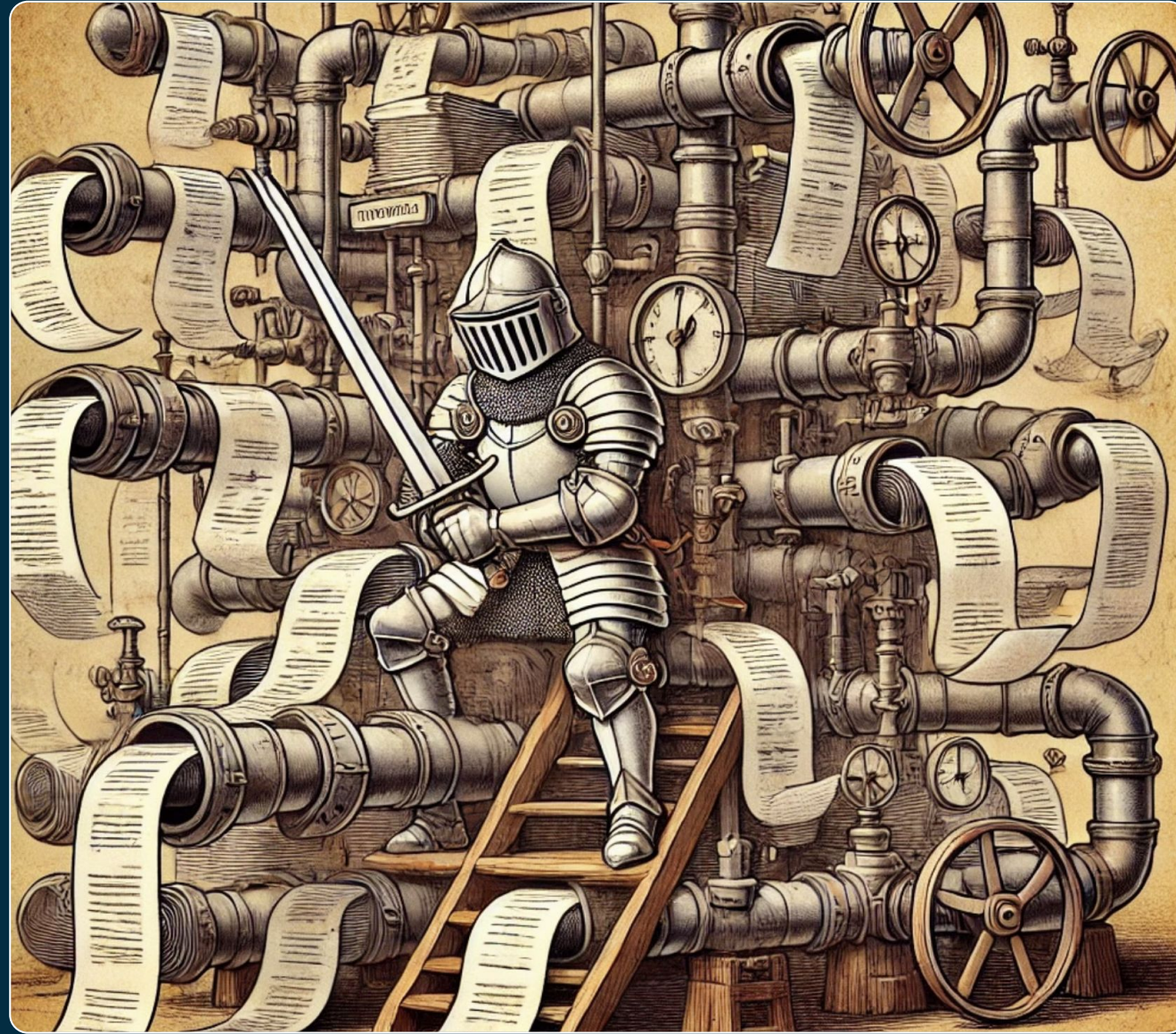


Benefits from full Mobile Tracing Coverage

- Detect issues on the Client + Edge Layers
- Understand #impacted users
- Understand Retry behavior of users
- Understand Network latency
- Understand UI Performance
- Understand degradations of Business KPIs
- Detect general UX degradations



Data Operations



The “EURO” Incident in 2022

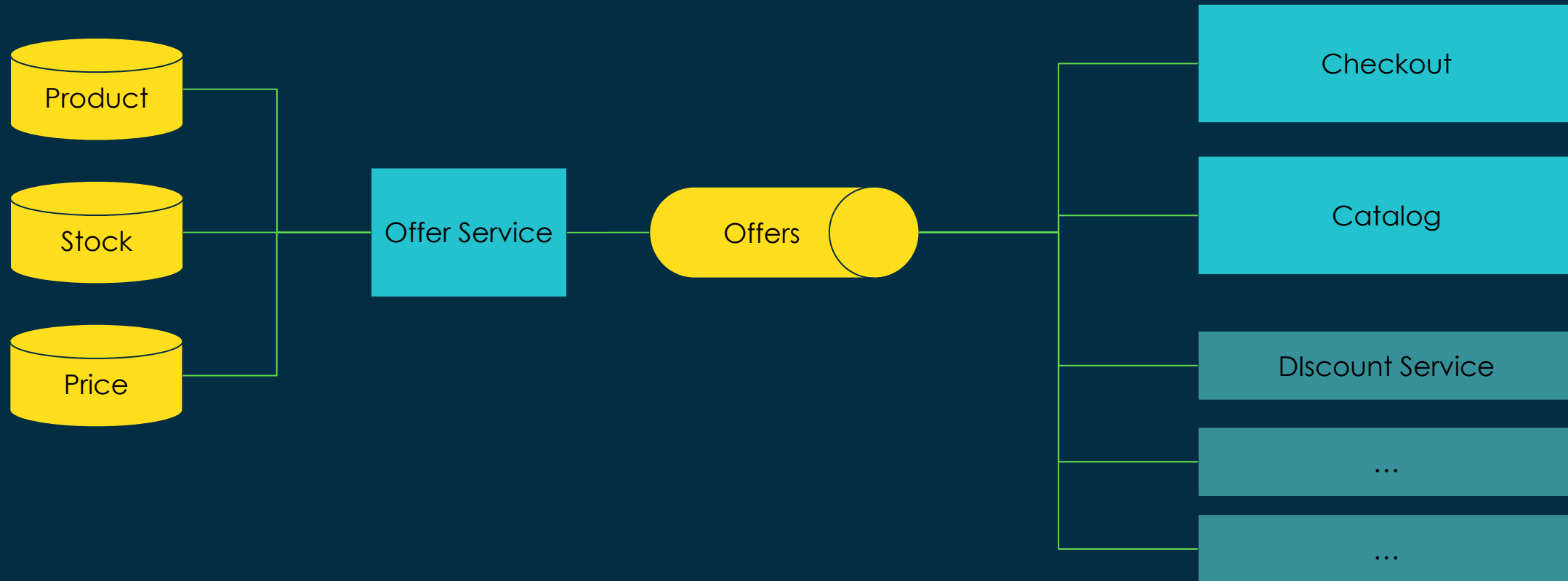
The Culprit

```
{
  "offer_id": "OF12345678",
  "product": {
    "product_id": "PR98765432",
    "name": "Men's Classic Leather Jacket",
    "category": "Men's Clothing > Jackets",
    "brand": "UrbanStyle",
    "description": "A premium classic leather jacket for men",
    "material": "Leather",
    "color": "Black",
    "size": ["S", "M", "L", "XL"],
    "images": [
      "https://alando.com/images/products/PR98765432_1.jpg",
      "https://alando.com/images/products/PR98765432_2.jpg"
    ],
    "tags": ["leather", "men's fashion", "jackets", "urban style",
"winter"]
  },
  "price": {
    "currency": "EURO",
    "current_price": 129.99,
    "original_price": 159.99,
    "discount": {
      "percentage": 18,
      "description": "Autumn Sale"
    }
  },
  "stock": {
    "available": true,
    "quantity": 45,
    "locations": [
```

This should be "EUR".



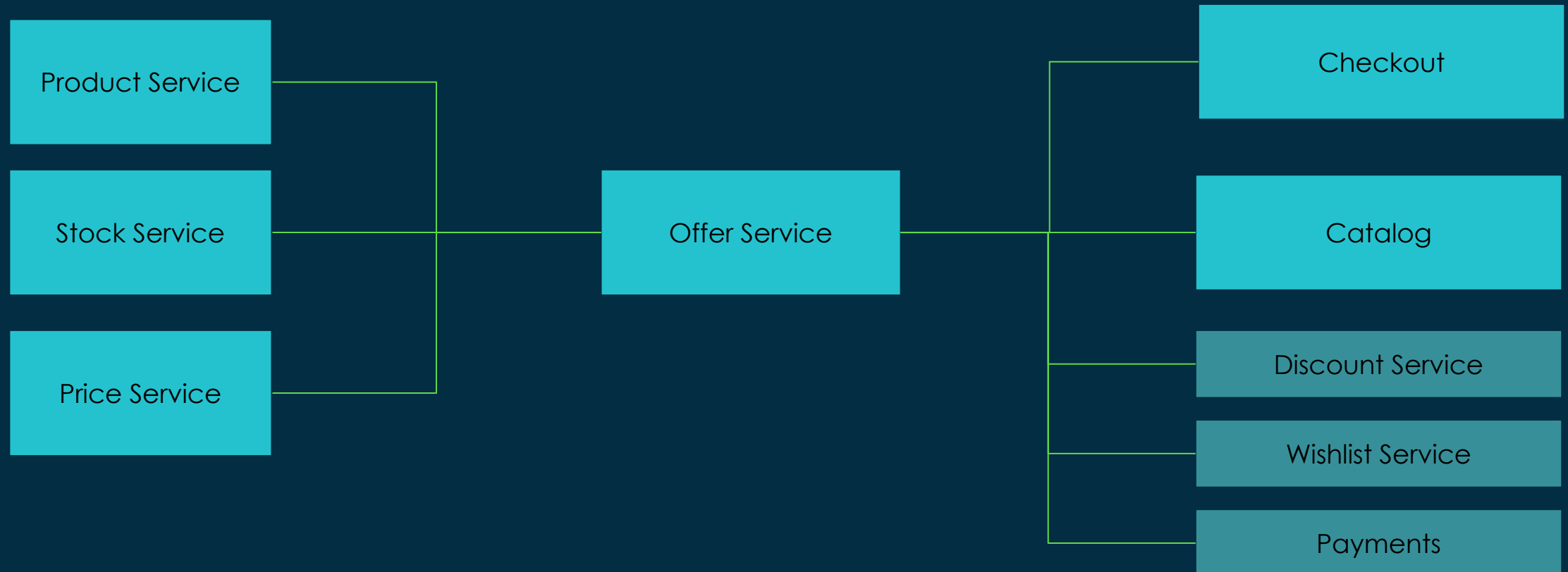
The EURO Incident - Data Architecture



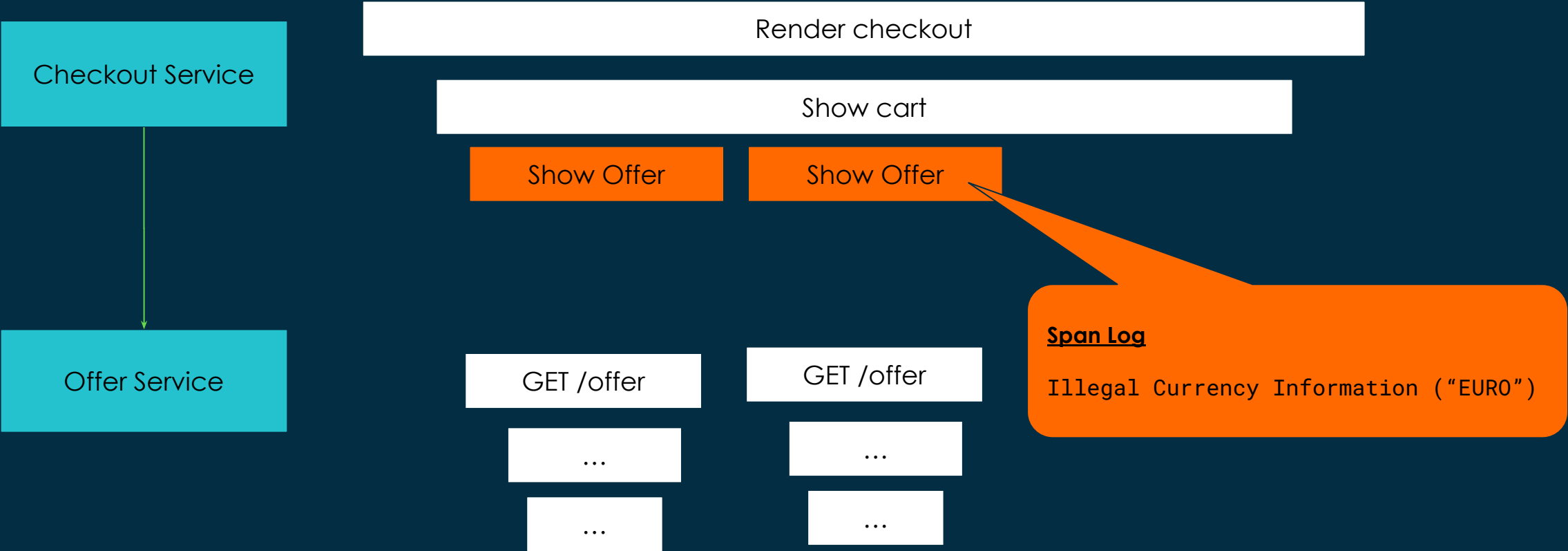
The EURO Incident - Take Aways

- Significant delay (30min) between cause (bad deployment) and impact.
- Significant delay (3h) between fix (rollback) and mitigation of symptoms.
- Limited help from Telemetry:
 - Metrics (Throughput, Backlog, Latency) - Not useful.
 - Tracing - Carried error information but no Causality Information

The “EURO Incident” with REST Architecture



Debuggin “EURO Incident” with Tracing



Data is of growing importance for ...

1. AI
2. Business Processes
3. Business Intelligence

Patterns from ~70 Data Incidents in 2024

Backlogs / Delays / Capacity

Data Quality

Unclear Handover

Unstructured Data

Schema Changes

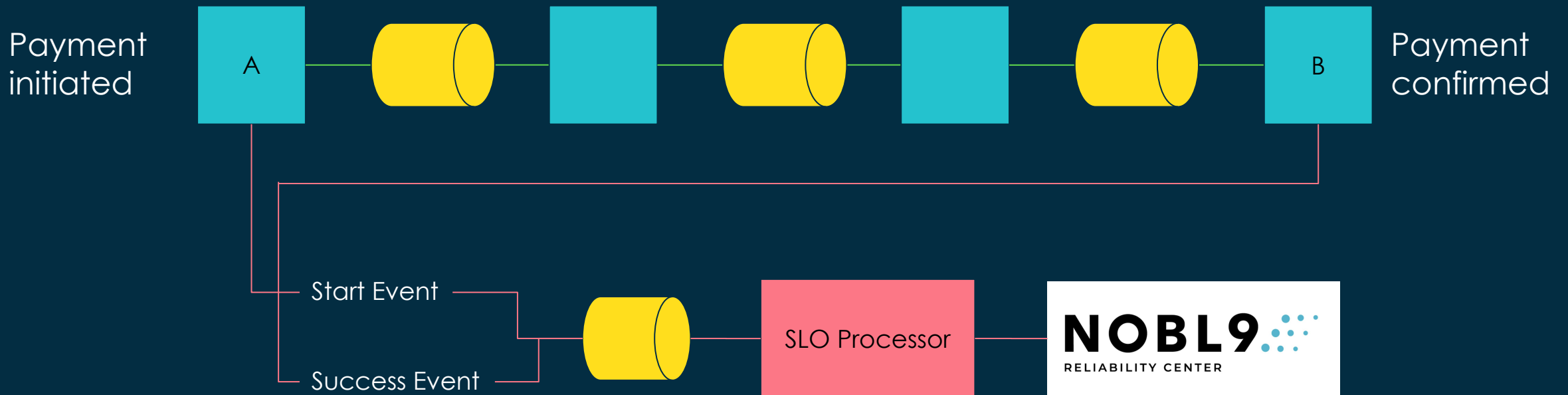
Long Data Chains

Unavailable Datasets

Integration w/ External Data Sources

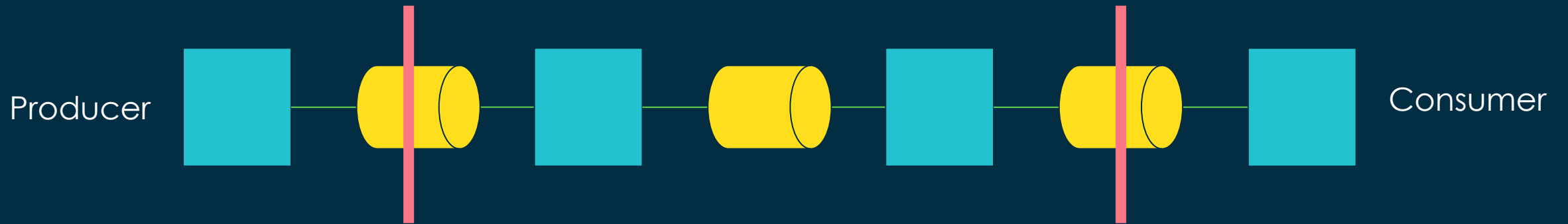
Investment Area - Data SLOs

- Check reliability of data processing pipelines end-2-end
- Check availability of datasets



Investment Area - Data Contracts

- Allow consumers to articulate expectations to data
- Detect Data Quality Problems earlier in the chain



Investment Area - Data Lineage

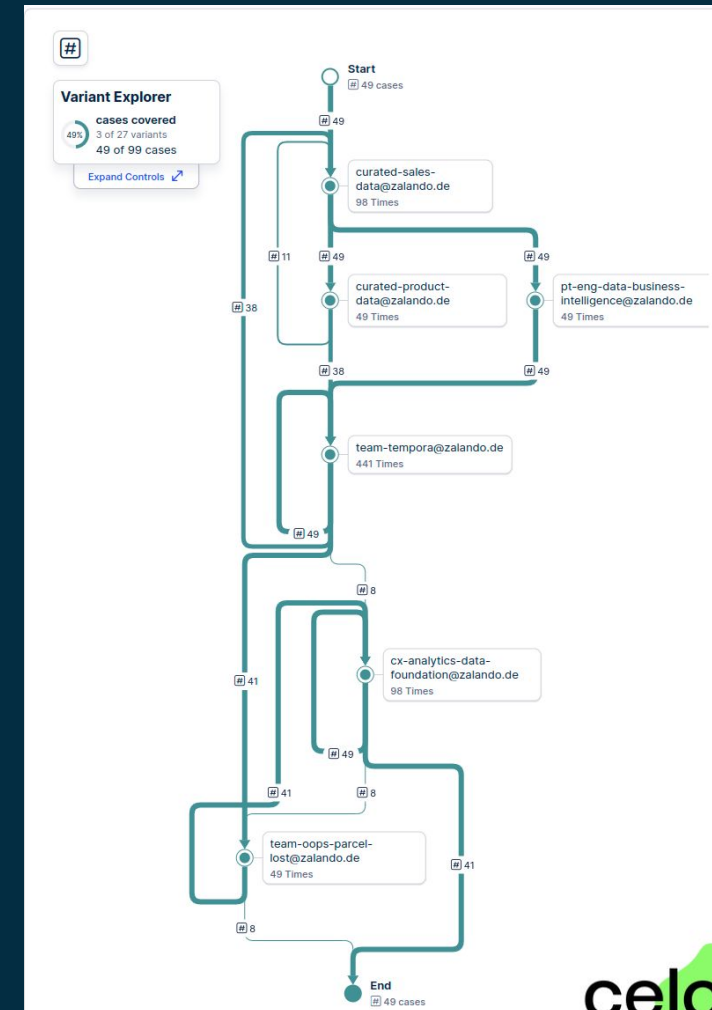
- Map dependency chain of data products
- Upstream - Who produced/processed this data?
- Downstream - Who is depending on my data?

OpenLineage

AN OPEN FRAMEWORK FOR DATA LINEAGE COLLECTION AND ANALYSIS

Experiment:

- Monitor timeliness of batch delivery for Business Analytics.
- Leverage Process Mining tool - Celonis.



celonis

Question: Do we have the right abstractions?

Data System Integration

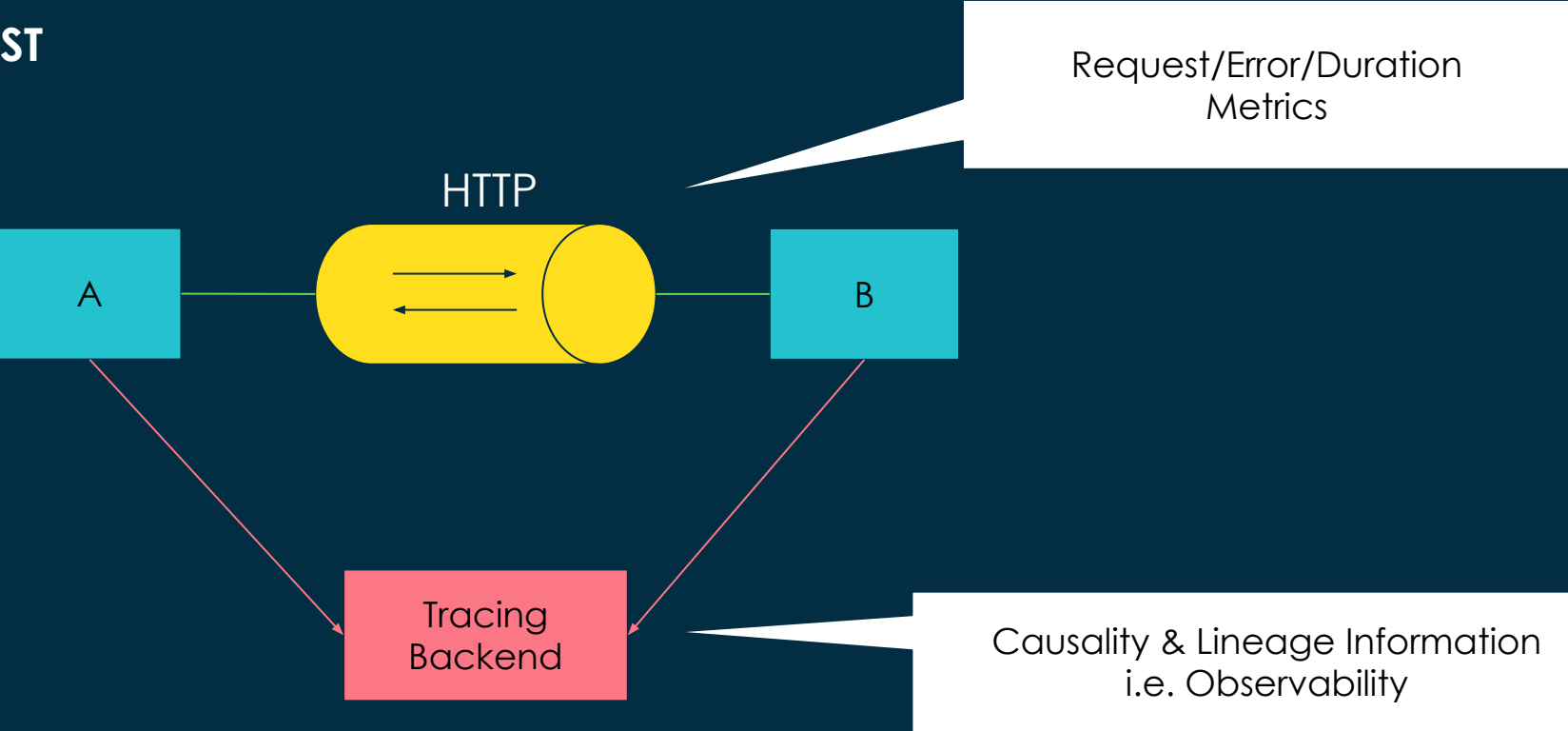


This looks like
half of
TCP-Connection



Question: Do we have the right abstractions?

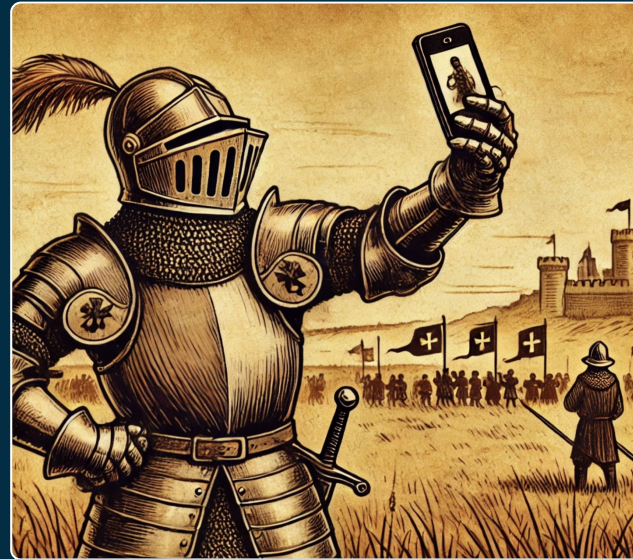
Microservice Integration / REST



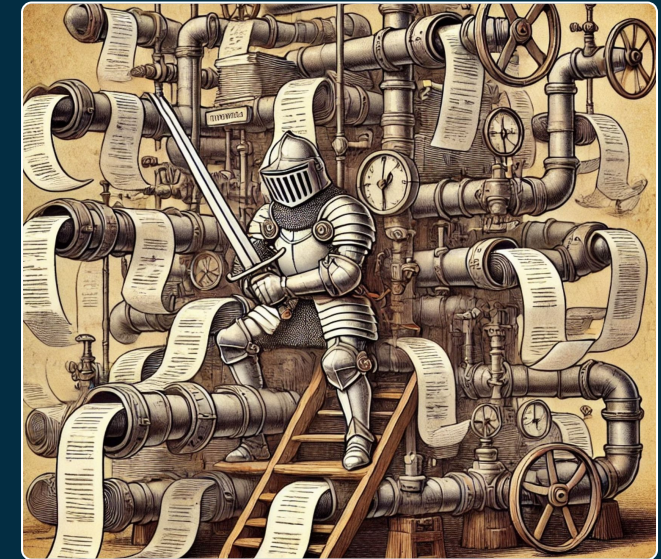
Three Frontiers of Reliability Engineering



Managing
for
Reliability



Mobile
Observability



Data
Operations

The SRE Triangle

