# A Field Guide to
# Reliability Engineering at Zalando

goto; Amsterdam 2024 - Heinrich Hartmann

# 👋 I'm Heinrich - Reliability Engineer

## Experience

**zalando** — Senior Principal SRE (2021)

**CIRCONUS** — Chief Data Scientist (2015)

**universität bonn** — PhD in Mathematics (2011)

## Talking Reliability since 2015

- SRECon - <u>Statistics for Engineers</u>
- DevOps Berlin - <u>Zalando's quest to Operate 10K…</u>
- SLOConf - <u>The State of the Histogram</u>
- P99 Conf - <u>How to measure Latency</u>
- FOSDEM - <u>Latency SLOs Done Right</u>
- <u>Circllhist - A Histogram Data Structure… (arxiv)</u>

## Find me on

heinrichhartmann.com
LinkedIn, X

**Menu**

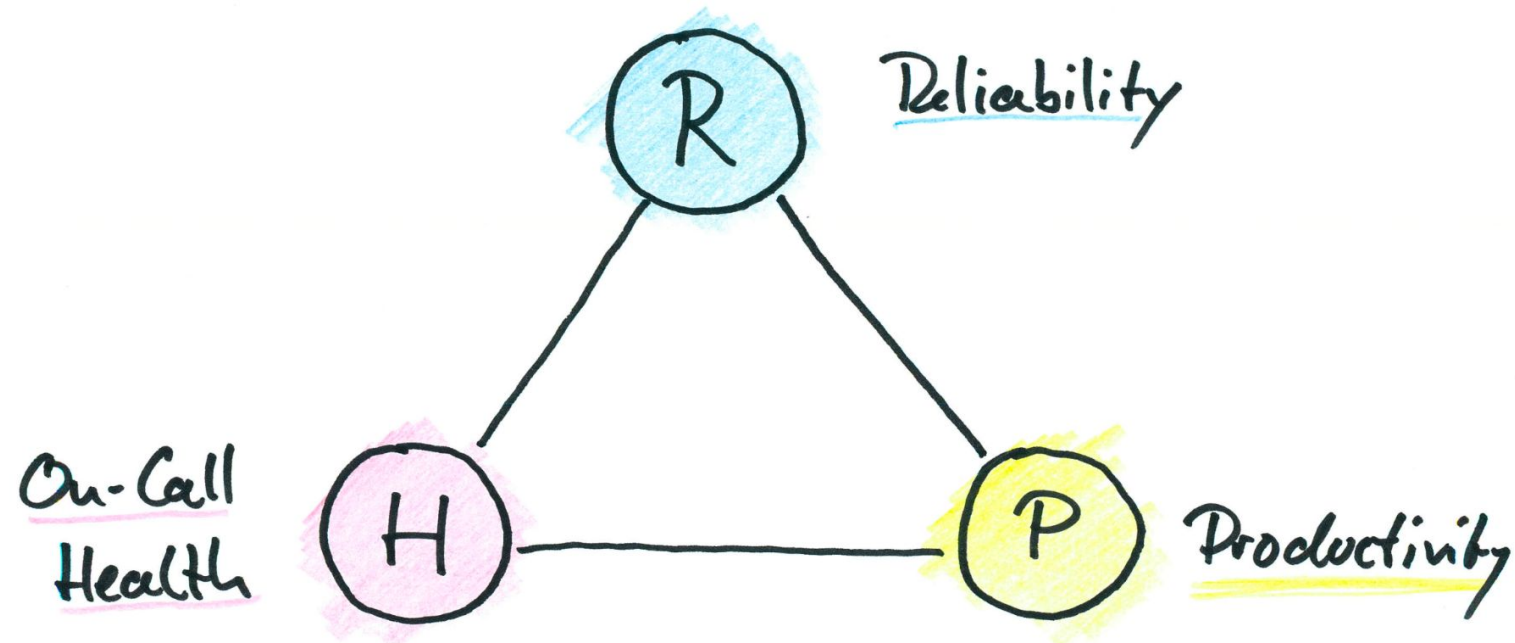goto; Amsterdam 20204. Heinrich Hartmann @ Zalando

# Principles

# Mission

**Protect the User Experience** from operational failures while keeping an eye on (1) Developer Productivity and (2) On-Call Health.

# #1 Rule of Operations

## Obsess about User Experience.

The SRE Triangle

R — Reliability
H — On-Call Health
P — Productivity

# #2 Rule of Operations

## Engineering for Reliability involves people as much as it involves technology.

# Engineering Reliability at Scale

**Small Company (~10 FTE)**

- Alerts & Dashboards
- Logging

**Medium Company (~100 FTE)**

- Incident Management
- Observability
- On-call rotations
- Playbooks
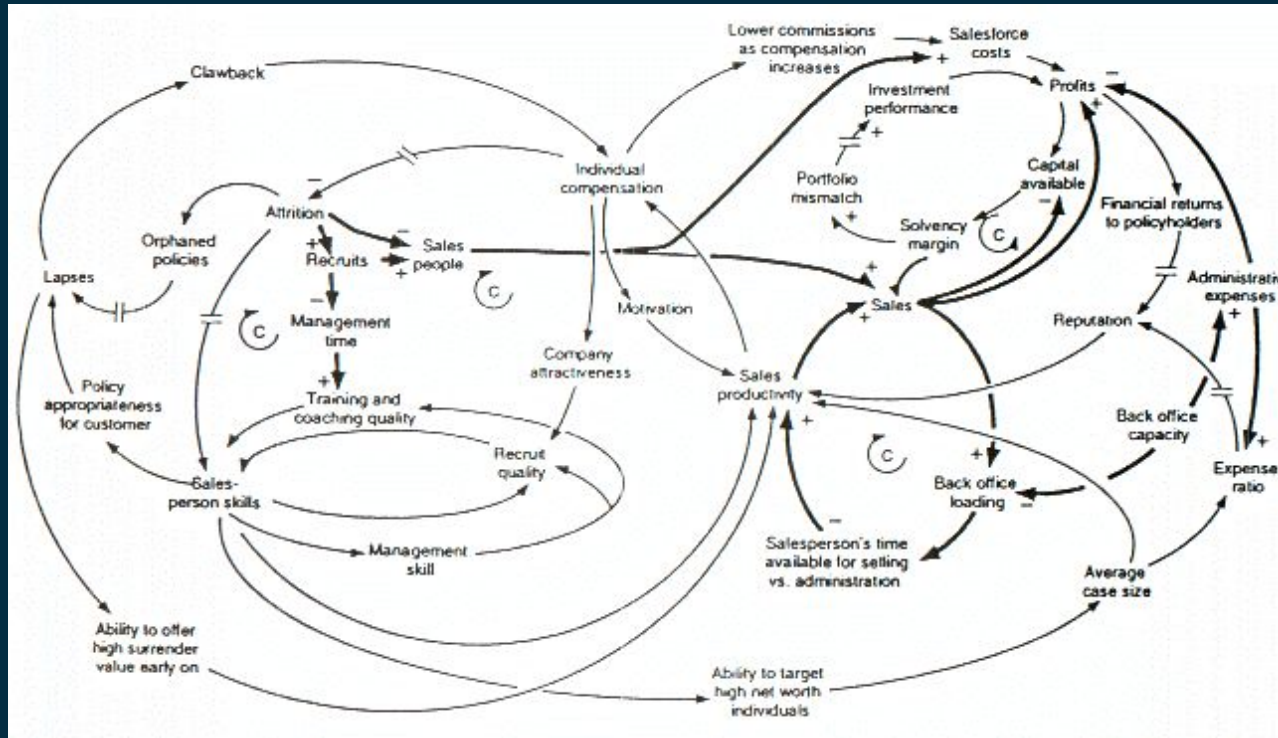- WORM Meeting

**Large Company (>1k FTE)**

- WORM Cascades
- Risk Management
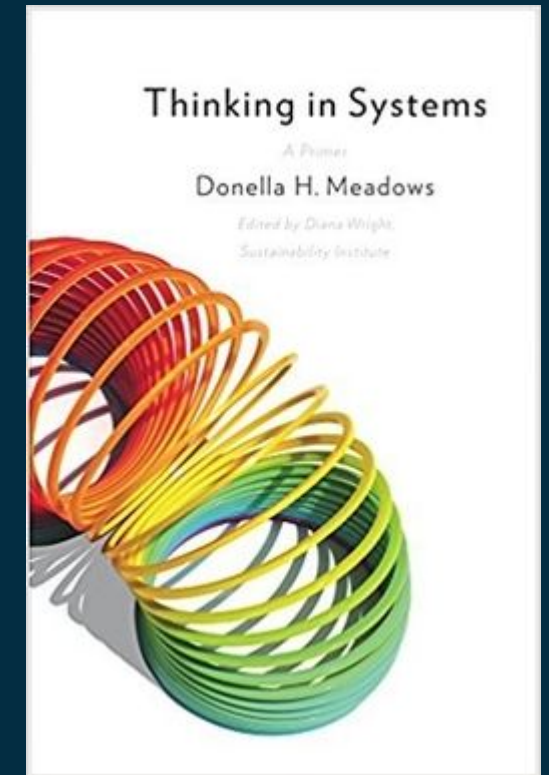- SRE Community & Guilds
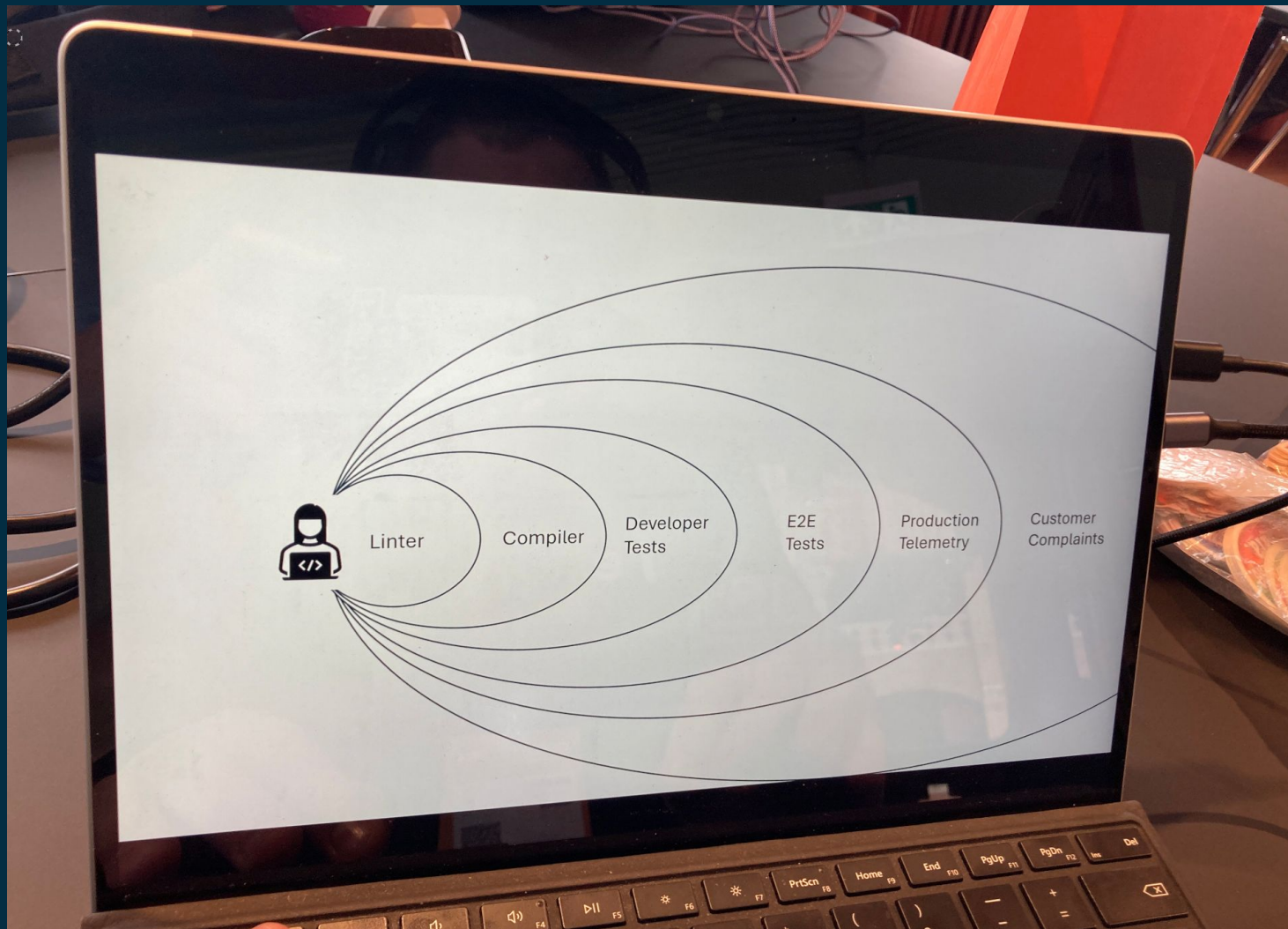
People Problems

Technical Problems

# Engineering Socio-Technological Systems

with "Systems Theory"



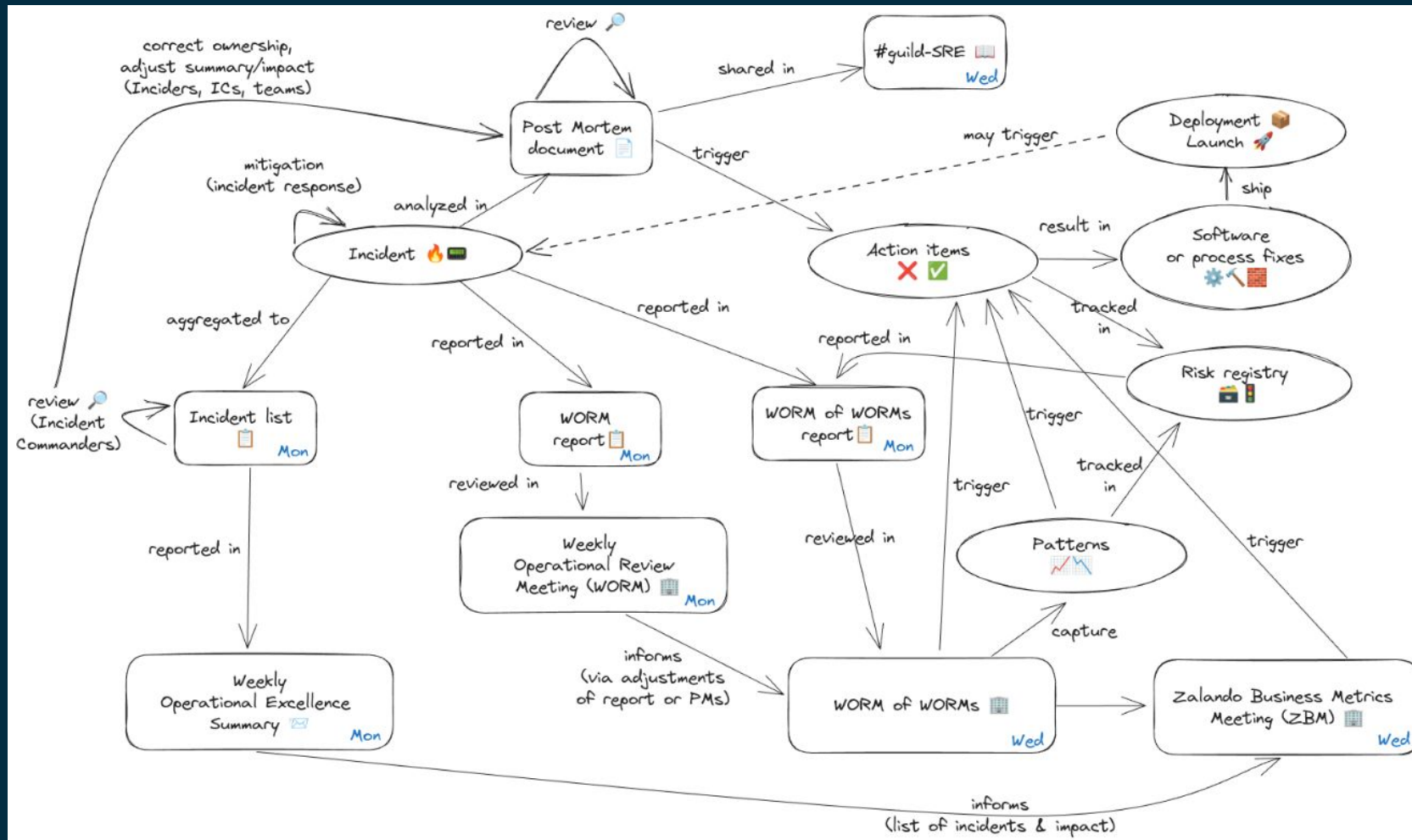Example: Causal Loop Diagram - source: wikipedia
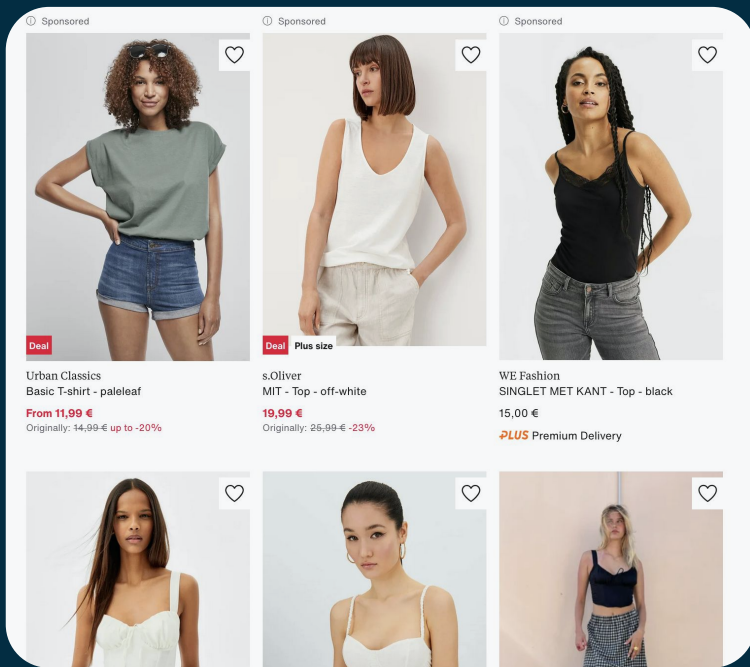
% Martin Thwaites @ Honeycomb GOTO 2024

goto; Amsterdam 20204. Heinrich Hartmann @ Zalando

# Reliability "Flywheel" at Zalando

# Context

- One of the leading fashion platforms in EU

- Founded in 2008

- 14.6 bn EUR Revenue / 50M+ active Customers

- 25 Countries

- 3K Tech Employees

- 3K+ Micro Services

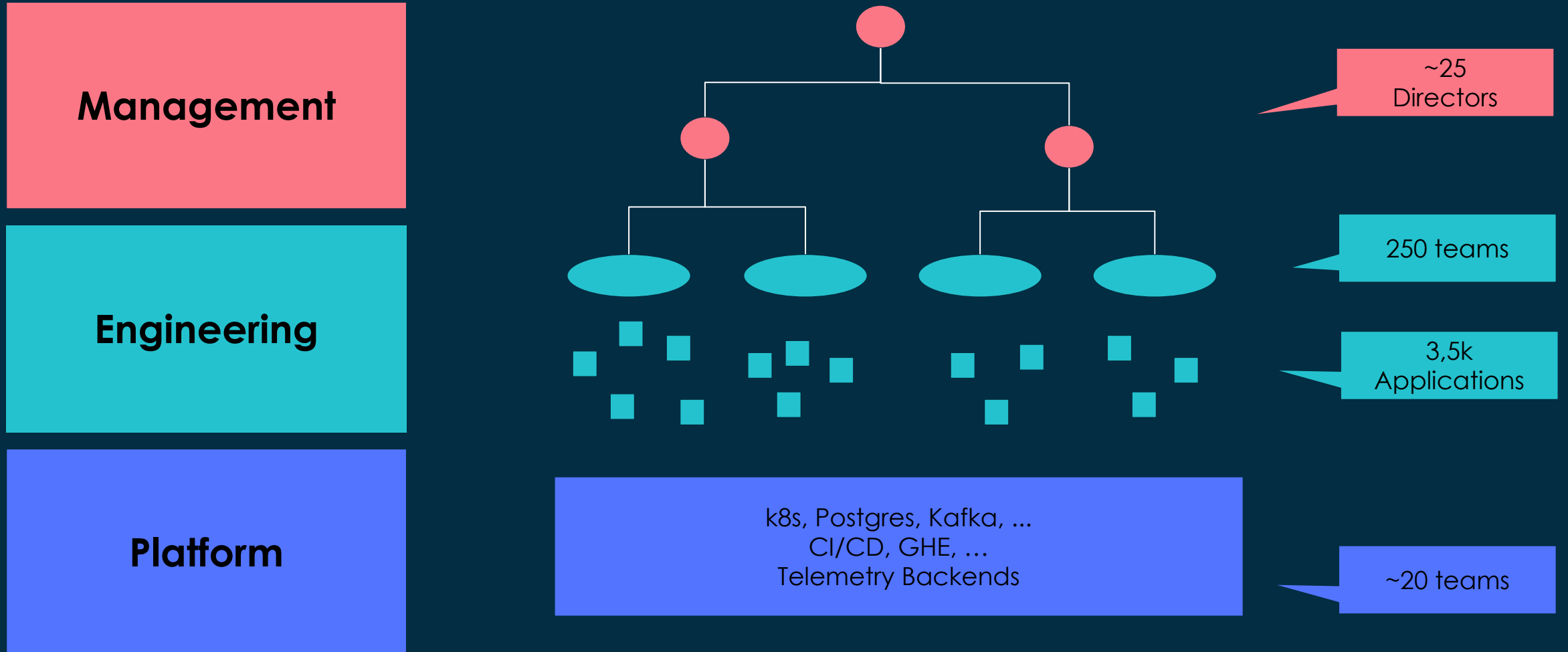Zalando
Service Graph

# Don't separate People and Technology

**Conway's law**

Technology Structures mirror People Structures.

**Law of DevOps**

You build it, you run it!

# Systems Model of Zalando



**Management** — ~25 Directors

**Engineering** — 250 teams, 3,5k Applications

**Platform** — k8s, Postgres, Kafka, ... CI/CD, GHE, ... Telemetry Backends — ~20 teams

# Where do we stand?

+ Operating "transactional" Microservices
+ Protecting the Business
+ Preparing for High-Load Events


- Understanding User Experience
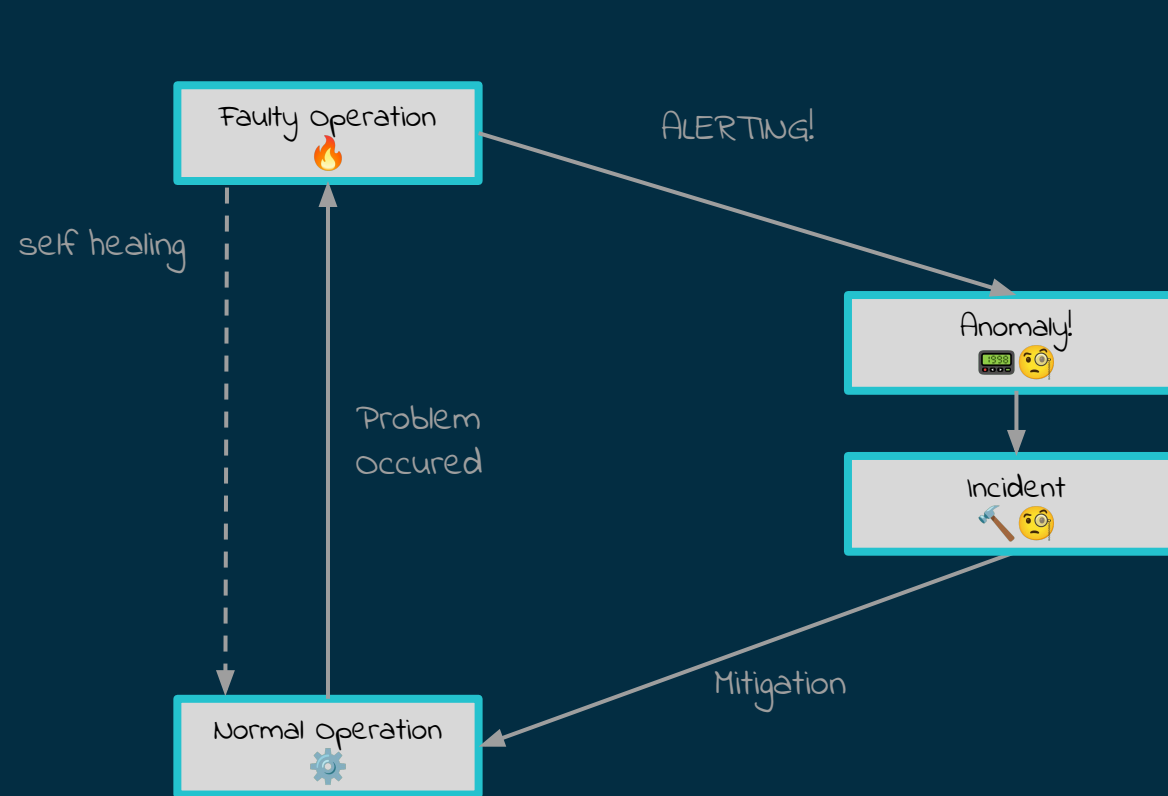- Reliability of Data Systems / Business Processes

# Operations at Zalando

# Alerting

# Why Alerting?

Reduce Time to Detect user-facing issues.

# Alerting as Feedback Loop

# #3 Rule of Operations

**Alert on User Experience ("Symptoms")
not on Server Experience ("Causes").**

- Alert on error rates of user-facing "operations"

- Leverage SLO-based Alerting (if available)

- Don't alert on CPU Utilization

This is fine.

Incidents / #7006

11.06.2024, 09:32 (GMT+02:00)

**P1** Zalando eCommerce Platform - ~~...~~ error ratio > 0.28% over the last 6h and 30m

SEV3   adaptive-paging   alias:page/add-article-to-cart-mobile   error-rate   page-low   unified   +

**Description**

Description:
Stream Name: ~~...~~
Critical Threshold Violated: alert 'Zalando eCommerce Platform — ~~...~~ error ratio > 0.28% over the last 6h and 30m' is above 0.0028 (value is 0.002804)

Failing application detected:
coast-cart-service

Stream:
https://app.lightstep.~~...~~/49BcysGs?
anchor=1718091465&end_micros=~~...~~50000000&range=3600&start_micros=1718087550000000&utm_source=webhook

Alert:
https://app.lightstep.com/Production/moni~~...~~ion/miPhhl06j?
anchor=1718091465&end_micros=1718091150000~~...~~ge=3600&start_micros=1718087550000000&utm_source=webhook

Related playbooks:
- FS-CART-008 — ~~Restore DynamoDB tables from backup~~
- EP-CX-ASSORTMENT-PURCHASE-~~RESTRICTIONS~~ ~~...~~ ~~Update purchase restrictions~~
- FS-CART-023 — Re~~...~~it
and there are 12 more

**COMMUNICATION**

Conferences

Incident command center

Stakeholder communication

Incident status updates

Send update

**ASSOCIATED ALERTS**

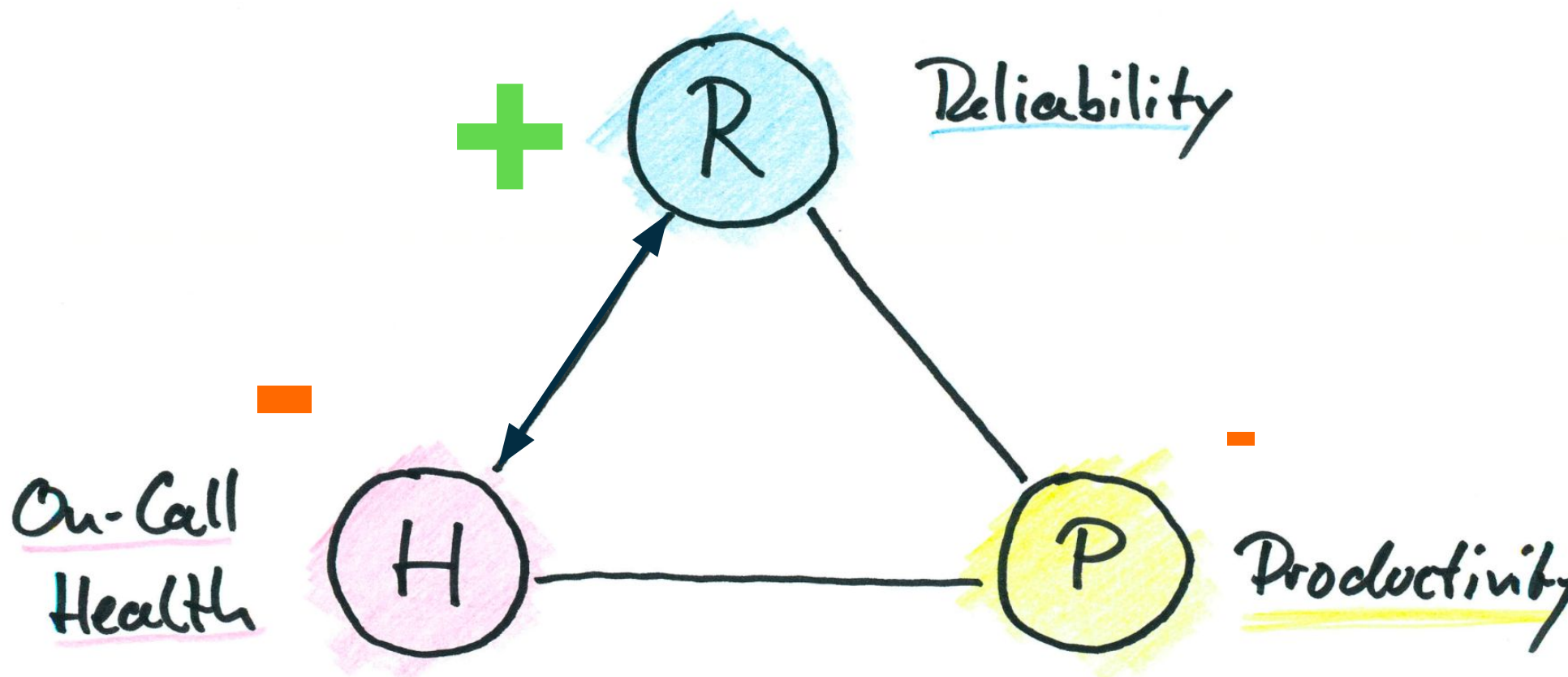See alerts

**RESPONDERS**                    + Add responder

Cart                              AWARE
Responder team

goto; Amsterdam 20204. Heinrich Hartmann @ Zalando

# Adding alerts trades Reliability of On-Call Health

goto; Amsterdam 20204. Heinrich Hartmann @ Zalando

# Review On-Call Health Weekly!

| 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 | |
| Builders | | | | | | | | | | | | |
| | - | 1 | 2 | - | 2 | 2 | 1 | 5 (-2 ▼) | 3 (+2 ▲) | 8 (+0 ▶) | 1.14 / day | ... consciously ... |
| | - | - | - | 5 | - | 7 | - | 1 (+0 ▶) | 11 (+11 ▲) | 12 (+11 ▲) | 1.71 / day | One legitimate alert, the ...sitives ...ons (e.g. ... resilience |
| Cloud ... | - | 2 | 1 | 2 | - | - | - | 4 (-7 ▼) | 1 (-2 ▼) | 5 (-9 ▼) | 0.71 / day | |
| ...undamentals | - | - | - | 1 | - | - | - | 1 (+1 ▲) | - (+0 ▶) | 1 (+1 ▲) | 0.14 / day | |
| ... | - | - | - | 1 | - | - | - | - (-3 ▼) | 1 (+0 ▶) | 1 (-3 ▼) | 0.14 / day | |
| ... | - | - | - | - | - | - | - | - (+0 ▶) | - (+0 ▶) | - (+0 ▶) | 0.0 / day | |
| ...us | - | - | - | - | - | - | - | - (+0 ▶) | - (+0 ▶) | - (+0 ▶) | 0.0 / day | |
| | - | 15 | 1 | - | - | - | - | 1 (+0 ▶) | 15 (-58 ▼) | 16 (-58 ▼) | 2.29 / day | ...to the ...es ...ay |

# Dashboards

# Why Dashboards?

- Reduce Time to Repair

- Look at them when you get alerted. Don't monitor dashboards.

- Starting point for understanding Service Health


- Every Application MUST have an Application Dashboard.

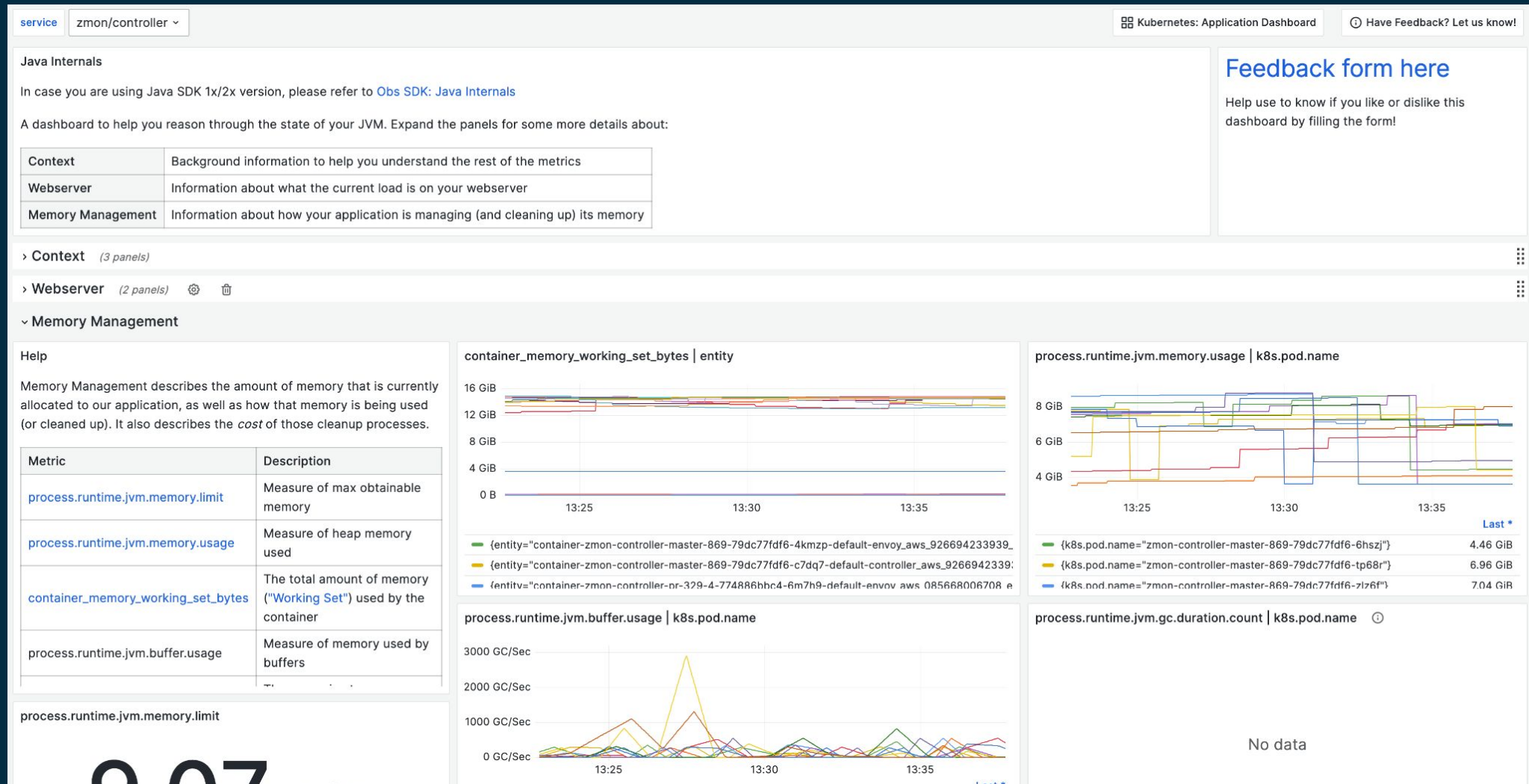- Managed Services come with Managed Dashboards.

# Managed Kubernetes Dashboard

# Managed REDIS Dashboard

# Managed JVM Internals Dashboard

# Zalando Application Dashboard Guidelines

1. Golden Signals
2. Entry Points
3. Dependencies
4. Saturation
5. Operational Insights
6. Storage

courtesy of Evgeni Sokolov & Miha Lunar

# Golden Signals Row - RED(S)

**Duration**                    **Requests**                    **Saturation**



**Errors**

w/ Evgeni Sokolov & Miha Lunar

# Entry Points Row
Golden Signals, again! - RED

w/ Evgeni Sokolov & Miha Lunar

# Saturation Row

... everything that can get saturated.
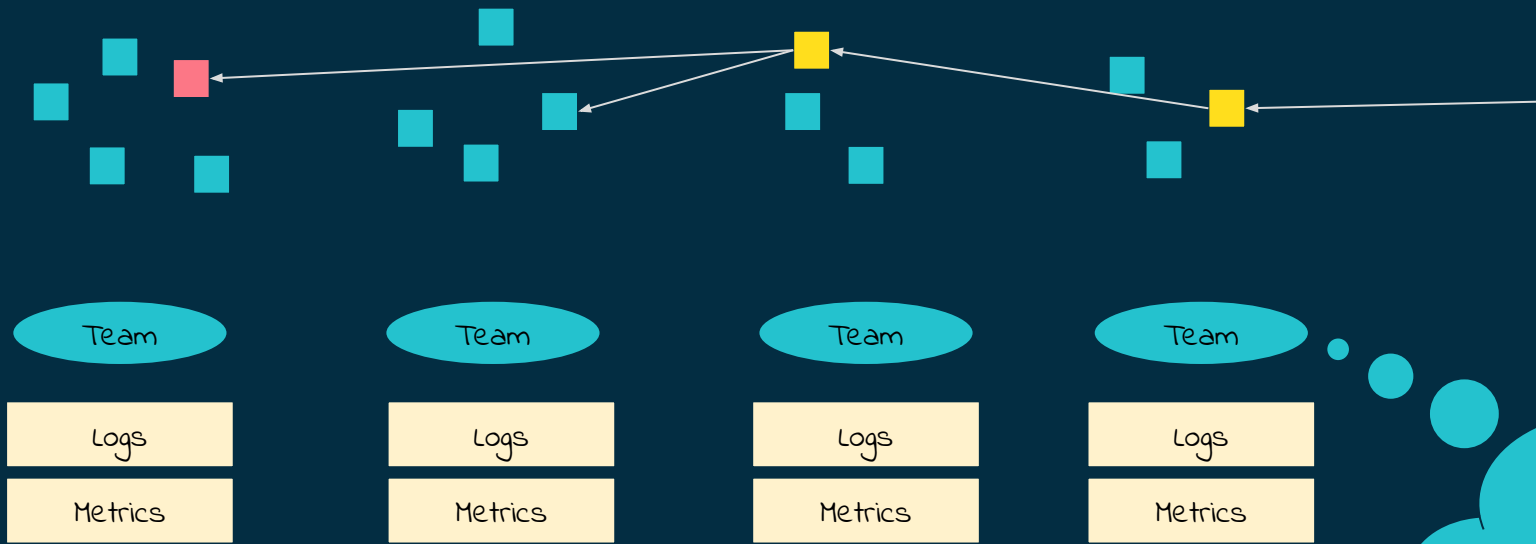
w/ Evgeni Sokolov & Miha Lunar

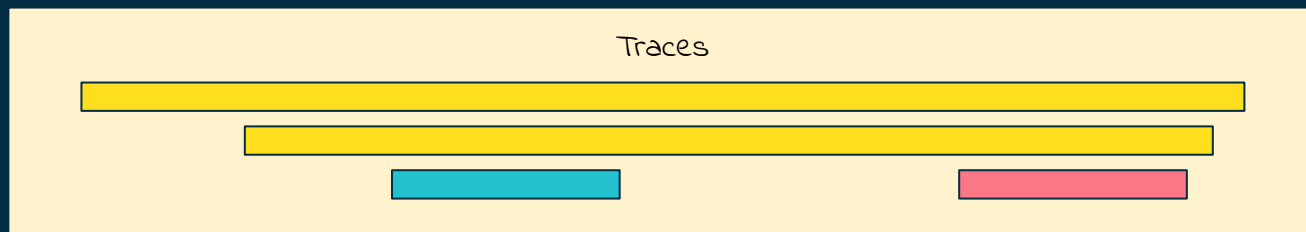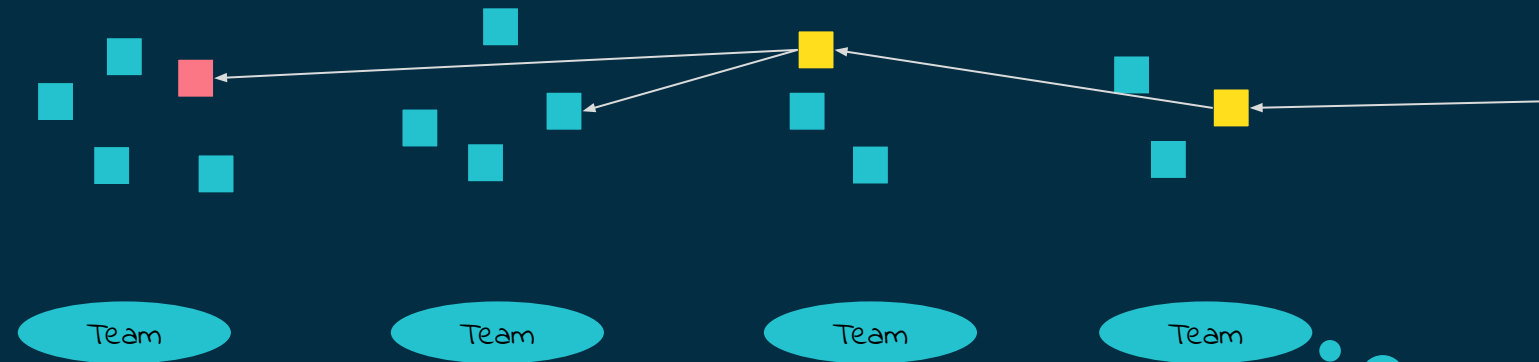# Observability

# Why Observability?

- Reduce Time to Repair

- Debug failures across team boundaries

- Understand User-Experience

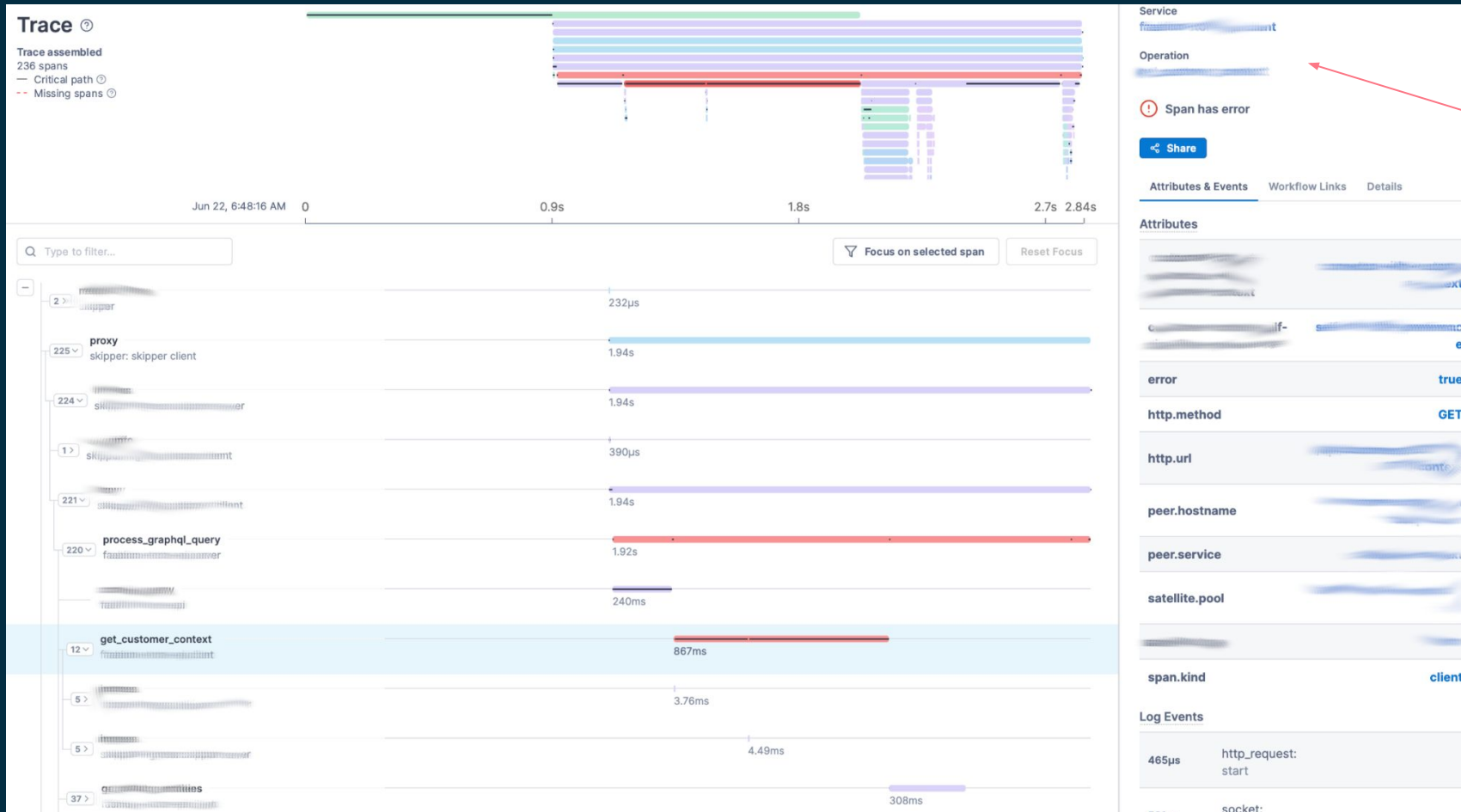- Basis for Alerting, Dashboards, Reporting, …

# Traditional Monitoring



Team

Team

Team

Team

Logs

Metrics

Logs

Metrics

Logs

Metrics

Logs

Metrics

Is **my application** healthy? **Which errors** does it throw?

# Observability



Traces

Is **the user** happy?
Which **operation** is failing?

goto; Amsterdam 20204. Heinrich Hartmann @ Zalando

# Example Trace from Zalando Front Page

# Zalando Developer Observability Guidelines

1. Use **OpenTelemetry** to instrument Applications.

2. Use **Distributed Tracing** to understand system behavior in the context of transactions (e.g. HTTP requests).

3. **Metrics** for precise counts & global resource statistics

4. Structured **Logging** for Lifecycle events

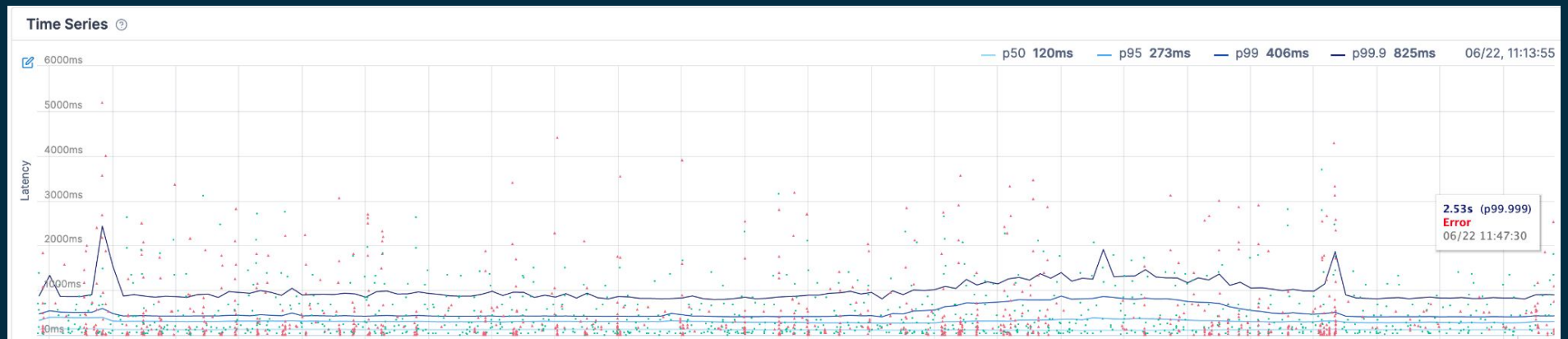# Monitor Reliability of Operations with "RED" Metrics

Operation: Reset Password

**Requests**

**Errors**

**Duration**

ServiceNow
**Cloud Observability**

# Observability SDKs
## based on Open Telemetry

```python
#!/usr/bin/env python3

import observability_sdk as obs

# Hook-up Zalando Backends
obs.initialize()

# Custom span
@obs.trace(name=..., attributes={...})
def add_to_cart():
        ...

# Custom metric
req_counter = obs.create_counter(
    name="total_requests",
    description="Total number of requests served",
    attributes = {...}
    unit="1",
    value_type=int,
)
def handle_request():
        req_counter.inc()
```

## Observability SDKs

| Language | Documentation | Implementation | Maturity Status |
|---|---|---|---|
| Java, Kotlin | on docs.zalando.net | on GHE | Supported |
| Python | on docs.zalando.net | on GHE | Supported |
| JavaScript | on docs.zalando.net | on GHE | Supported |
| Scala | on docs.zalando.net | on GHE | Beta |
| Go | on docs.zalando.net | on GHE | Alpha / ETA Q3'2023 |

**OpenTelemetry**

goto; Amsterdam 20204. Heinrich Hartmann @ Zalando

# SLOs

# Why SLOs?

- Provide Top-Down understanding of Reliability provided to the user

- Steer engineering investments into Reliability

- Quantify impact of incidents

- … also derive high-quality alerting rules

# #4 Rule of Operations

**SLIs quantify the reliability of a <u>User Experience</u>.**

**SLOs are Reliability targets for <u>managerial steering</u>.**

# Zalando SLOs on Business Operations

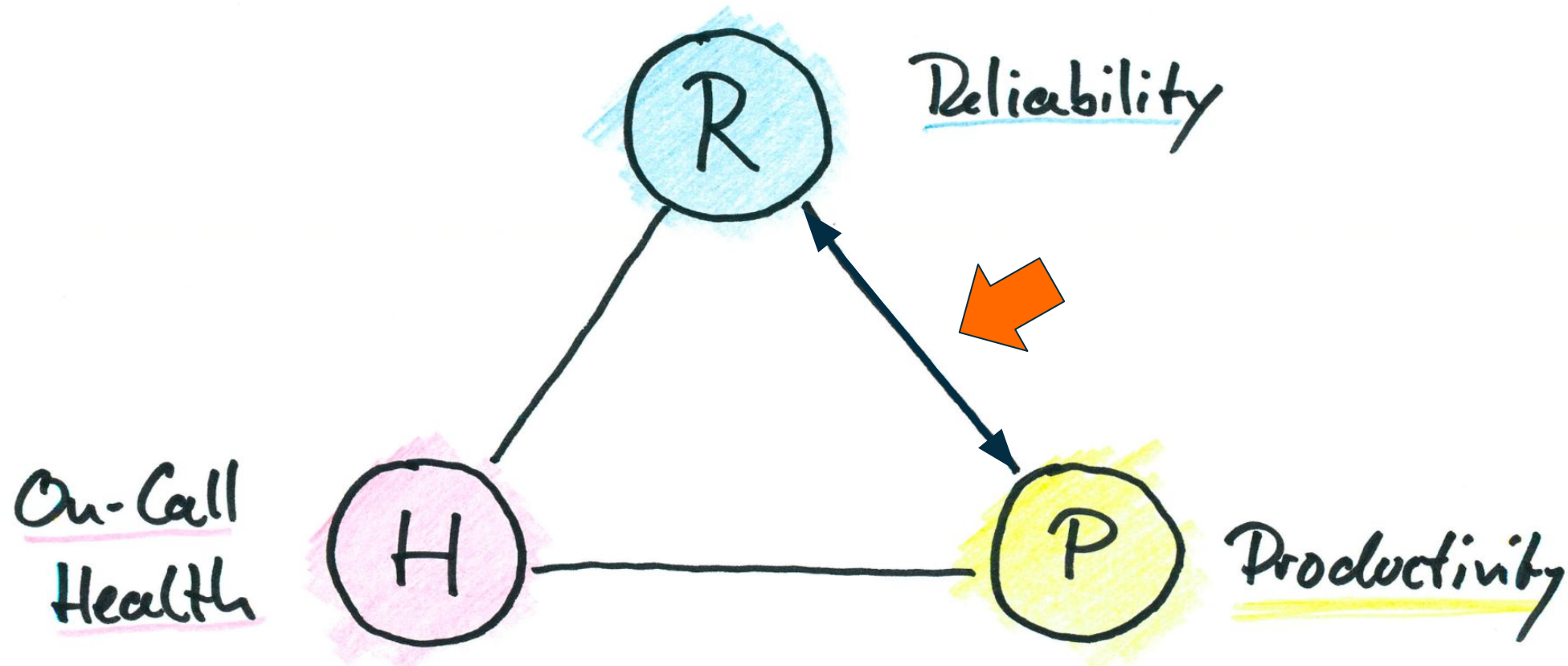# SLO Table Reviewed by Management



| 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% ⚠️ | ▓▓▓▓▓▓▓▓▓▓ to ▓▓▓▓ |
| 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% | ▓▓▓▓ | ▓▓▓▓ | ▓3% ✅ | |
| S▓▓ registration | 99.950% | ▓▓▓▓ | ▓▓▓▓▓ | ▓96 ✅ | |
| SS▓▓▓▓▓ on | 99.950% | 9.▓ | ▓▓▓▓% | ▓ ✅ | |
| S▓▓▓▓▓▓ ation | 99.950% | ▓% | ▓▓▓▓% | ▓▓▓ ✅ | |
| ▓▓ Step up authen▓▓ | 99.950% | 9▓ | ▓% | ▓▓▓% ✅ | |
| **Demand / Home & Content Visibility** | | | | | |
| ▓▓▓▓▓▓▓▓▓ | 99.000% | ▓▓▓% | ▓% | ▓% ✅ | |
| ▓▓▓▓▓▓▓▓▓ | 99.000% | ▓▓▓ | ▓▓▓ | ▓ ✅ | |

# SLOs are used to Prioritize Engineering Investments



goto; Amsterdam 20204. Heinrich Hartmann @ Zalando

# SLOs are also used to tune Alerting Sensitivity



goto; Amsterdam 20204. Heinrich Hartmann @ Zalando

# Decouple Alerting/Reporting SLOs to get more value!
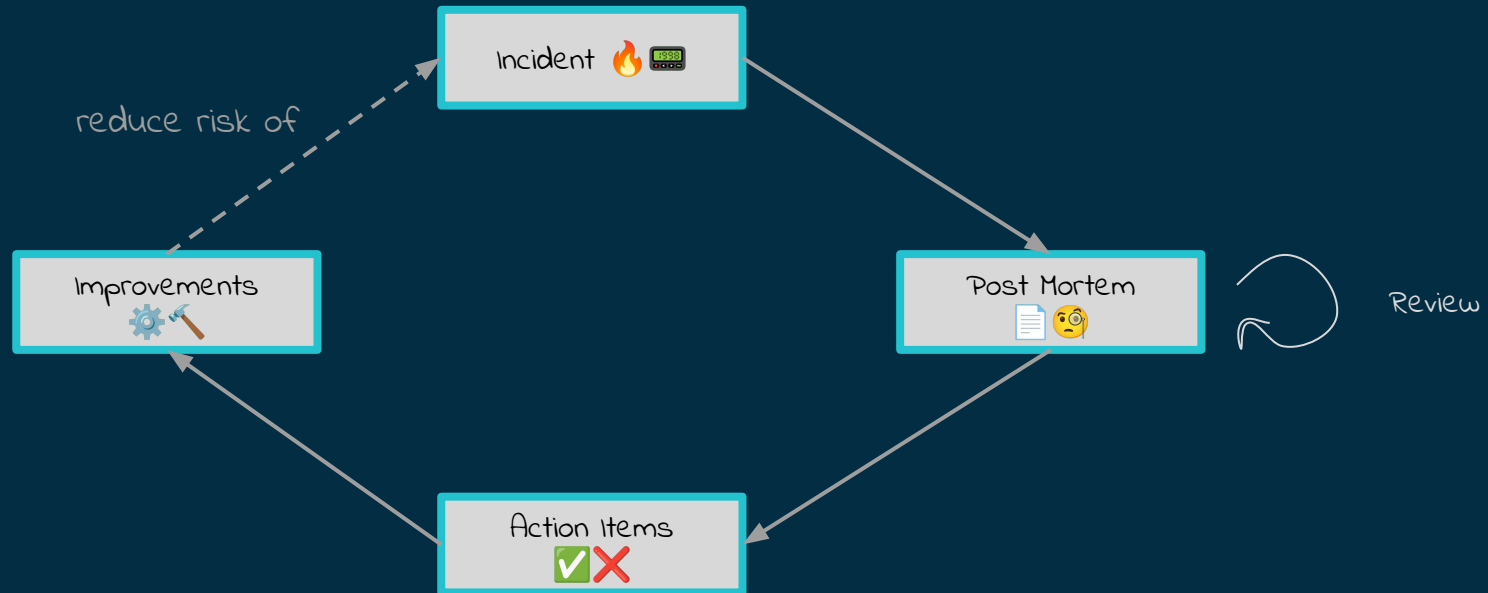


goto; Amsterdam 20204. Heinrich Hartmann @ Zalando

# Incident Process

# #5 Rule of Operations

**Past Failures lead the way towards future Reliability.**

# Incident Process as Feedback Loop



Incident 🔥📟

Post Mortem 📄🧐

Review

Action Items ✅❌

Improvements ⚙️🔨

reduce risk of

# Zalando Incident Process



Incident-Bot **App** · 22 Sept, 11:28 · Edited

ℹ️ **Incident**: [training room] cannot open cart
**Severity**: SEV3
**Status**: closed 🟢
**Involved Teams**: Size and Fit ,Cart
**Owning Team**: esre-txn
**Application**: size-advice-service
**Links**: Incident · Chat thread

This message will be automatically updated. Updates can take up to a minute to materialize.

🤝 **All conversation about incidents should happen in inline threads, please keep the main thread clean.** 🤝

5 replies  22 Sept, 14:21

---

Zalando confidential

## Post-Mortem Document

Title: {{TITLE}}
Severity: {{SEVERITY}}
Ticket: {{TINY_ID}}
Owner:
Driver & Authors:
Reviewer:
Categories:  select a category ▾  __copy for multiple categories__
Status:  PM in progress ▾

Documentation: How to Write a Post Mortem? | Post Mortem Reference | Post Mortem Checklist | Examples

### Summary

On __DATE__ between __IMPACT_STARTED__ and __RESOLUTION__ __CUSTOMER_GROUP__ experienced __DEGREDATION__ for a business impact of __BUSINESS_IMPACT__ This was triggered by __TRIGGER__ and repaired through __INTERVENTION__. Action items include __ACTION_ITEMS__. The incident surfaced __LESSONS_LEARNED__.

{{DESCRIPTION}}

### Impact

Customer Impact
- Markets impacted: __market__
- Propositions impacted: __proposition__
- Customer experience during the incident: __description__
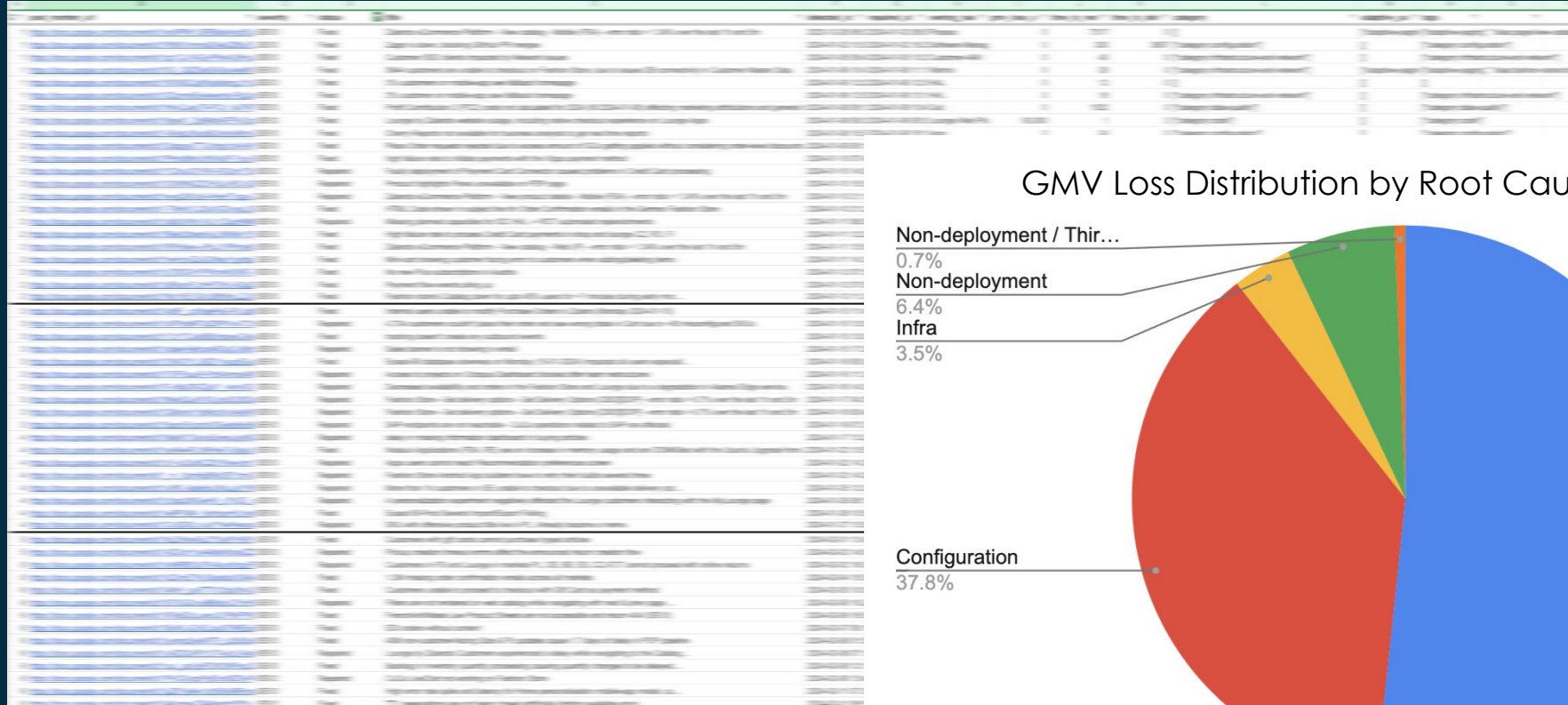
Business Impact
- __description__
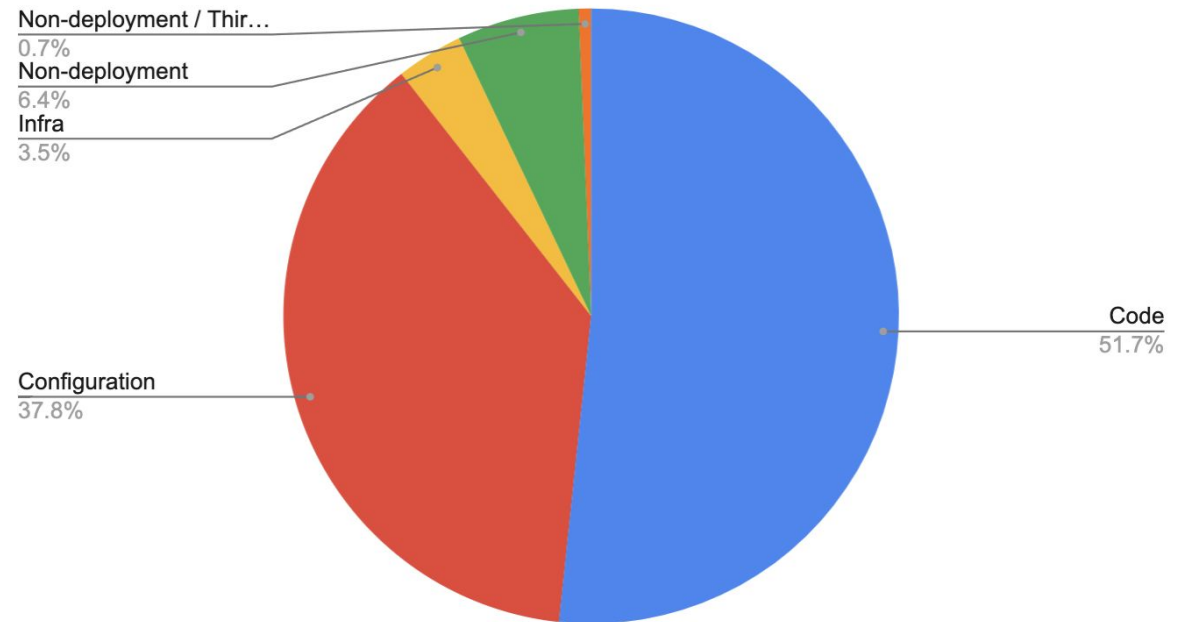
Internal Impact
- __description__

**1. Impact**

**2. Root Cause**

**3. Action Items**

# Zalando Severity Definitions

| | | |
|---|---|---|
| **SEV1** | **Example Incidents**<br>● Order Drop<br>● AWS Zone Outage | Ownership: Vice President |
| **SEV2** | **Example Incidents**<br>● Payments processor degraded<br>● Order confirmation emails delayed | Ownership: Director |
| **SEV3** | **Example Incidents**<br>● Users don't receive voucher<br>● Lounge users see not personalised articles | Ownership: Head of Engineering |

# Incident Insights every Quarter



GMV Loss Distribution by Root Cause in Q?/20??

Non-deployment / Thir…
0.7%
Non-deployment
6.4%
Infra
3.5%

Code
51.7%

Configuration
37.8%

# Weekly Operational Review Meeting

# #6 Rule of Operations

You get what you inspect.

# Reliability Reports
Supporting WORM Meetings on all Levels

Auto Generated Google Doc

## WORM Agenda

- Incident Review -> Patterns?
- SLO Review
- Open Post Mortems
- On-Call Health

# Zalando WORM Cascade

# Rules of Operations

1. Obsess about User Experience.

2. Engineering for Reliability involves People & Technology.

3. Alert on User Pain ("Symptoms") not Server Pain ("Causes").

4. SLIs quantify the reliability of a User Experience.

5. Past Failures lead the way towards future Reliability.

6. You get what you inspect.

## Thank you!

> Heinrich@HeinrichHartmann.com
#Let's talk Reliability! 💚

goto; Amsterdam 20204. Heinrich Hartmann @ Zalando