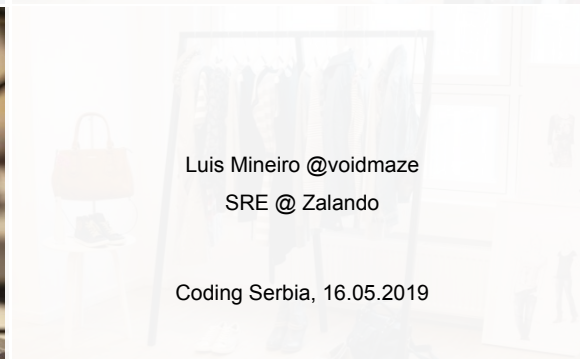zalando

# ALERTING
# MONITORING

AND ALL THAT JAZZ

Luis Mineiro @voidmaze

SRE @ Zalando

Coding Serbia, 16.05.2019

**ZALANDO AT A GLANCE**

**~ 5.4** billion EUR **revenue 2018**

**> 300 million** visits per month

**> 15,500** employees in Europe

**> 80%** of visits via mobile devices

**> 27 million** active customers

**> 400,000** product choices

**~ 2,000** brands

**17** countries

as of March 2019

zalando

# ZALANDO OFFICES

1. BERLIN **HEADQUARTERS**
2. ERFURT **TECH OFFICE**
3. MÖNCHENGLADBACH **TECH OFFICE**
4. DORTMUND **TECH HUB**
5. DUBLIN **TECH HUB**
6. HELSINKI **TECH HUB**
7. HAMBURG **ADTECH LAB**
8. LISBON **TECH HUB**

as of March 2019

# WE ARE CONSTANTLY INNOVATING TECHNOLOGY

**HOME-BREWED, CUTTING-EDGE & SCALABLE** technology solutions

help our brand to **WIN ONLINE**
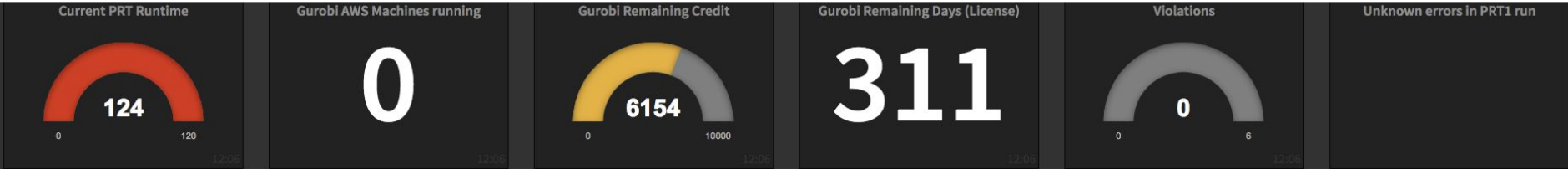
**8** international tech locations

**> 2,000** employees at

**HQs** in Berlin

zalando

## Dashboard widgets

| Current PRT Runtime | Gurobi AWS Machines running | Gurobi Remaining Credit | Gurobi Remaining Days (License) | Violations | Unknown errors in PRT1 run |
|---|---|---|---|---|---|
| 124 | 0 | 6154 | 311 | 0 | |

Search for alerts

0  Hide Widgets

Cluster Unit information has not been fetched within the last day "prtinputsteering-testing"  9h

⚠ Cluster Unit information contains too few unique entries "prtinputsteering-testing" ({minimum_num_unique_lines} < 320,000)  10h

⚠ Cluster Unit information contains too many duplicated entries "prtinputsteering-testing"  10h

⚠ AWS instance is unhealthy: {application_id} (2)  6d

⚠ Issues on AWS Instances: {issues} (2)  6d

Long running PRT task: "production" 124 min (id 1648; priority 10000) on "prt-scenario-api-production"  3m

6 open pull requests - Please review and merge them now (until we have no more than 5)!  8m

⚠ Forecast Parameter API is NOT available  10h

Undesired access to S3 buckets? ["32.0 requests for bucket prt-internationalised-output-mysqldump-experimental(undesired  1d

# Looks familiar?

# TERMINOLOGY

## MONITORING

Collecting, processing, aggregating, and displaying real-time quantitative data about a system, such as query counts and types, error counts and types, processing times, and server lifetimes.

## ALERT

A notification intended to be read by a human and that is pushed to a system such as a bug or ticket queue, an email alias, or a pager.

*SRE Book, Chapter 6: Monitoring Distributed Systems*

zalando

# MONITORING

Your monitoring system should address two questions: **what's broken**, and **why**?

The "what's broken" indicates the **symptom**; the "why" indicates a (possibly intermediate) **cause**.

"**What**" versus "**why**" is one of the most important distinctions in writing good monitoring with maximum signal and minimum noise.

*SRE Book, Chapter 6: Monitoring Distributed Systems*

zalando

# ALERTING CLASSIFICATION

| Urgency | Name | Delivery |
|---|---|---|
| Will be addressed... eventually | Report | Dashboards or nowhere (/dev/null) |
| Predicted to fail "soon" | Ticket | An issue tracker or *cough*, Email |
| Urgently and actively get the attention of a specific human | Page | A pager, cell phone or something going *beep* *beep* |

zalando

# WHAT TO ALERT ON

Alerting should be both **hard failure–centric** and **human-centric**.

*Distributed Systems Observability e-Book, Chapter 2: Monitoring and Observability*

Symptoms are a better way to capture more problems more comprehensively and robustly with less effort - "**symptom-based monitoring**," in contrast to **"cause-based monitoring"**.

*Rob Ewaschuk, "My Philosophy on Alerting"*

Keep alerting simple, **alert on symptoms**. Aim to **have as few alerts as possible**, by alerting on symptoms that are associated with end-user pain rather than trying to catch every possible way that pain could be caused.

*Prometheus Best Practices, https://prometheus.io/docs/practices/alerting/*

zalando

# SERVICE LEVEL OBJECTIVES

You should pick SLOs that represent the **most critical aspects of the user experience**.

*Google Cloud Platform Blog, Building good SLOs - CRE life lessons*

Start by thinking about (or finding out!) **what your users care about**, not what you can measure.

Choose **just enough SLOs to provide good coverage** of your system's attributes. Defend the SLOs you pick: if you can't ever win a conversation about priorities by quoting a particular SLO, it's probably not worth having that SLO.

*SRE Book, Chapter 4 - Service Level Objectives*

zalando

**ALERTING STRATEGY**

# What to alert on:

"hard failure–centric and human-centric"

zalando

**ALERTING STRATEGY**

# What to alert on:

"hard failure–centric and human-centric"

"symptom-based monitoring"

zalando

**ALERTING STRATEGY**

# What to alert on:

"hard failure–centric and human-centric"

"symptom-based monitoring"

"alert on symptoms"

zalando

**What to alert on:**

"hard failure–centric and human-centric"

"symptom-based monitoring"

"alert on symptoms"

"symptoms that are associated with end-user pain"

zalando

# ALERTING STRATEGY

## Service Level Objectives:

"most critical aspects of the user experience"

zalando

# ALERTING STRATEGY

## Service Level Objectives:

"most critical aspects of the user experience"

"what your users care about"

zalando

**ALERTING STRATEGY**

"hard failure–centric and human-centric"

"symptom-based monitoring"

"alert on symptoms"

"symptoms that are associated with end-user pain"

**=**

"most critical aspects of the user experience"

"what your users care about"

zalando

# ALERTING STRATEGY

## What to alert on:

"Keep alerting simple"

zalando

**ALERTING STRATEGY**

## What to alert on:

"Keep alerting simple"

"Aim to have as few alerts as possible"

zalando

**Service Level Objectives:**

"just enough [...] to provide good coverage"

zalando

**ALERTING STRATEGY**

"Keep alerting simple"

"Aim to have as few alerts as possible"

**=**

"just enough SLOs to provide good coverage"

zalando

Service Level Objective = Symptom + Threshold

zalando

# Page only when your SLO is missed or in danger of being missed

zalando

# ALERTING CHECKLIST

1. Does this rule detect **an otherwise undetected condition** that is urgent, actionable, and actively or imminently **user-visible**?

zalando

# ALERTING CHECKLIST

1. Does this rule detect **an otherwise undetected condition** that is urgent, actionable, and actively or imminently **user-visible**?

2. **Will I ever be able to ignore this alert**, knowing it's benign?

zalando

# ALERTING CHECKLIST

1.  Does this rule detect **an otherwise undetected condition** that is urgent, actionable, and actively or imminently **user-visible**?

2.  **Will I ever be able to ignore this alert**, knowing it's benign?

3.  Does this alert **definitely indicate** that users are being **negatively affected**?

# ALERTING CHECKLIST

1. Does this rule detect **an otherwise undetected condition** that is urgent, actionable, and actively or imminently **user-visible**?

2. **Will I ever be able to ignore this alert**, knowing it's benign?

3. Does this alert **definitely indicate** that users are being **negatively affected**?

4. **Can I take action in response to this alert**?

zalando

# ALERTING CHECKLIST

1. Does this rule detect **an otherwise undetected condition** that is urgent, actionable, and actively or imminently **user-visible**?

2. **Will I ever be able to ignore this alert**, knowing it's benign?

3. Does this alert **definitely indicate** that users are being **negatively affected**?

4. **Can I take action in response to this alert**?

5. **Are other people getting paged for this issue**?

*SRE Book, Chapter 6: Monitoring Distributed Systems*

zalando

**"Load average is high"**

zalando

**"Cassandra node is down"**

zalando

**"EC2 instance is unhealthy"**

zalando

# CREDIT

The majority of these slides were inspired or contained references to the excellent work from many industry experts and publications:

**People:**

- Rob Ewaschuk
- Björn Rabenstein
- Cindy Sridharan
- Charity Majors
- And many more...

**Publications:**

- Site Reliability Engineering (Book)
- The Site Reliability Workbook (Book)
- Distributed Systems Observability (e-Book)

zalando

**ХВАЛА**

# QUESTIONS?

Don't miss my next talk tomorrow at 11:30
"Are we all on the same page? Let's fix that"

Luis Mineiro @voidmaze

We're Hiring!
**https://jobs.zalando.com**

zalando