# PERFORMANCE OPTIMIZATION

## SCALING, HORIZONTAL AND VERTICAL

# SCALING, HORIZONTAL AND VERTICAL

# SCALING PLANS

**Auto Scaling Minimum**

> Health Check monitors running instances within an Auto Scaling Group.

> If an unhealthy instance is found, it can be replaced.

**Manual Scaling**

> Specify a new minimum for your Auto Scaling Group.

> Manually invoke Auto Scaling Policies.

**Scheduled Scaling**

> Scaling functions are performed as a function of time and date.

**On Demand Scaling**

> You create a policy to scale your resources.

> Define when to scale using CloudWatch Alarms.

**For more, look here**

# RDMA

**RDMA stands for Remote Direct Memory Access**

**Zero-copy networking**

- by enabling the network adapter to transfer data from the wire directly to application memory or
- from application memory directly to the wire

**Wiki on RDMA**

# RDMA ON AZURE

Most of the HPC VM sizes (HBv2, HB, HC, H16r, H16mr, A8 and A9) feature a network interface for remote direct memory access (RDMA) connectivity

**OS**

- Linux and Windows

**Infiniband, drivers**

**MPI (Message Passing Interface)**

**RDMA network address space**

- The RDMA network in Azure reserves the address space 172.16.0.0/16. To run MPI applications on instances deployed in an Azure virtual network, make sure that the virtual network address space does not overlap the RDMA network.

# OPTIMIZING PERFORMANCE ON BIGQUERY

**Four Key Elements of Work**

- **I/O** — How many bytes did you read?

- **Shuffle** — How many bytes did you pass to the next stage?

    - Grouping — How many bytes do you pass to each group?

- **Materialization** — How many bytes did you write to storage?

- **CPU work** — User-defined functions (UDFs), functions

# AVOID INPUT / OUTPUT WASTEFULNESS

**Advice for BigQuery, your situation may be different**

- Do not SELECT *, use only the columns you need
- Filter using WHERE as early as possible in your queries
- Do not use ORDER BY without a LIMIT

# DATA SKEW

**Data skew in BigQuery**

- But applicable to many areas or data handling, e.g. Spark
- Filter your dataset as early as possible (this avoids overloading workers on JOINs)
- BigQuery will automatically attempt to reshuffle workers that are overloaded with data
    - You may have to it yourself in other situations

Skewed Data creates an
imbalance between
BigQuery worker slots
(uneven data partition sizes)

8

# CAREFUL USE OF GROUP BY

**Again, BigQuery but applicable in many SQL situations**

- Best when the number of distinct groups is small (fewer shuffles of data).
- Grouping by a high-cardinality unique ID is a bad idea.

| Row | contributor_id | LogEdits |
|-----|----------------|----------|
| 1 | 2221364 | 4 |
| 2 | 104574 | 4 |
| 3 | 73576 | 4 |
| 4 | 311307 | 4 |
| 5 | 291919 | 4 |
| 6 | 140178 | 4 |
| 7 | 181636 | 4 |
| 8 | 3661553 | 4 |
| 9 | 3600820 | 4 |
| 10 | 4737290 | 4 |
| 11 | 938404 | 4 |
| 12 | 295955 | 4 |
| 13 | 183812 | 4 |
| 14 | 1811786 | 4 |
| 15 | 8918196 | 4 |
| 16 | 561624 | 4 |
| 17 | 5338406 | 4 |

← Do not Group on an ID

# DEBUGGING, MONITORING, PERFORMANCE TUNING

**Stackdriver is a Google tool**

**However, it is cross-cloud**

- and a good tool to explain the issues

**October 2020, rebranded as Google Cloud Operations**

**Important : please take the following slides as an approach example even if you do not use the Google cloud**

# WHAT STACKDRIVER DOES

**Combines metrics, logs, and metadata**

- On Google Cloud Platform (GCP)
- On Amazon Web Services
- on-premises infrastructure
- or a hybrid cloud

**Allows to**

- understand service behaviors and issues
- from a single comprehensive view of your environment
- take action if needed

# GOOGLE STACKDRIVER

**A multicloud service**

**An example to discuss the issues**

### Error Reporting
Error notifications
Error dashboard

### Debugger
Production debug snapshots
Conditional snapshots
IDE integration

### Logging
Platform, system, and app logs
Log search/view/filter
Logs-based metrics

### Monitoring
Platform, system, and app metrics
Uptime/health checks
Dashboards
Alerts

### Trace
Latency reporting
Per-URL latency sampling

### Profiler
Low-impact profiling of
applications in production

# INCREASE APPLICATION RELIABILITY



**Monitor GCP, AWS, and Multi-Cloud Environments**

Get the insight that you need with minimal configuration. Monitor hosted services and cloud architectures.

**Identify Trends, Prevent Issues**

Visualize trends via flexible charts and dashboards. Identify risks using scoring, anomaly detection, and prediction.

**Reduce Monitoring Overhead**

Spend less time correlating metrics, alerts, and logs across disparate systems. Don't worry about scaling tools.

**Improve Signal-to-Noise**

Reduce false positives and alert fatigue with advanced alerting designed for modern distributed systems.

**Fix Problems Faster**

Uptime and health checks notify you quickly when endpoints become inaccessible to your users. Drill down from alerts to dashboards to logs and traces to get to the root cause quickly.

# PERFORMANCE MANAGEMENT TOOLS

Stackdriver Trace

Stackdriver Debugger

Stackdriver Profiler

# STACKDRIVER TRACE

**Distributed tracing**

# STACKDRIVER PROFILER

# FOUR GOLDEN SIGNALS

Latency    Traffic    Errors    Saturation

# PRESENTATION ON E-CAS

 **E-CAS public**

# AWS CLOUDWATCH

A monitoring service for AWS cloud resources and the applications you run on AWS

Visibility into resource utilization, operational performance, and overall demand patterns

Custom application-specific metrics of your own

Accessible via AWS Management Console, APIs, SDK, or CLI

# AMAZON CLOUDWATCH FACTS

**Monitor other AWS resources**

- View graphics and statistics

**Set Alarms**

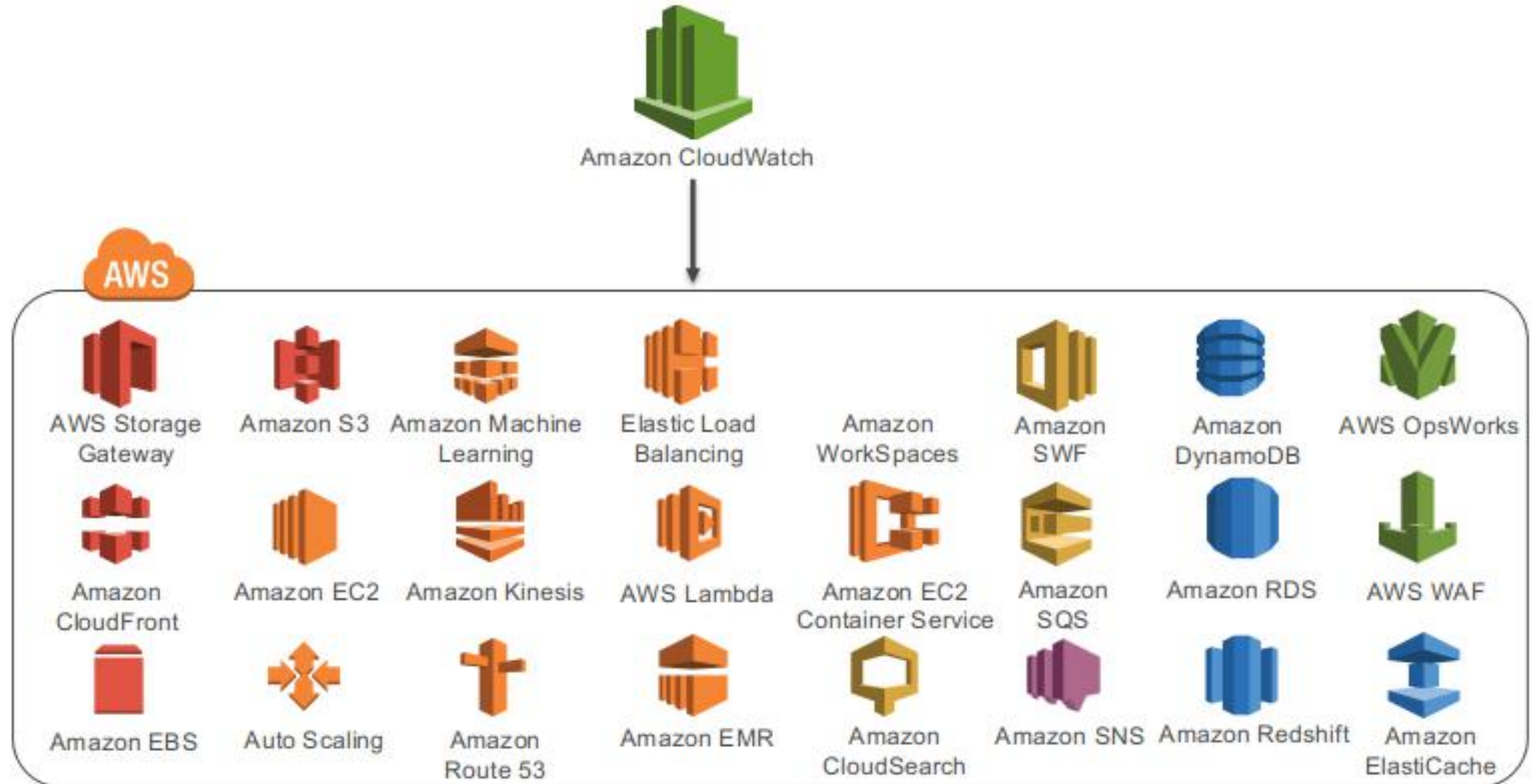# AMAZON CLOUDWATCH ARCHITECTURE

# CLOUDWATCH METRICS EXAMPLES
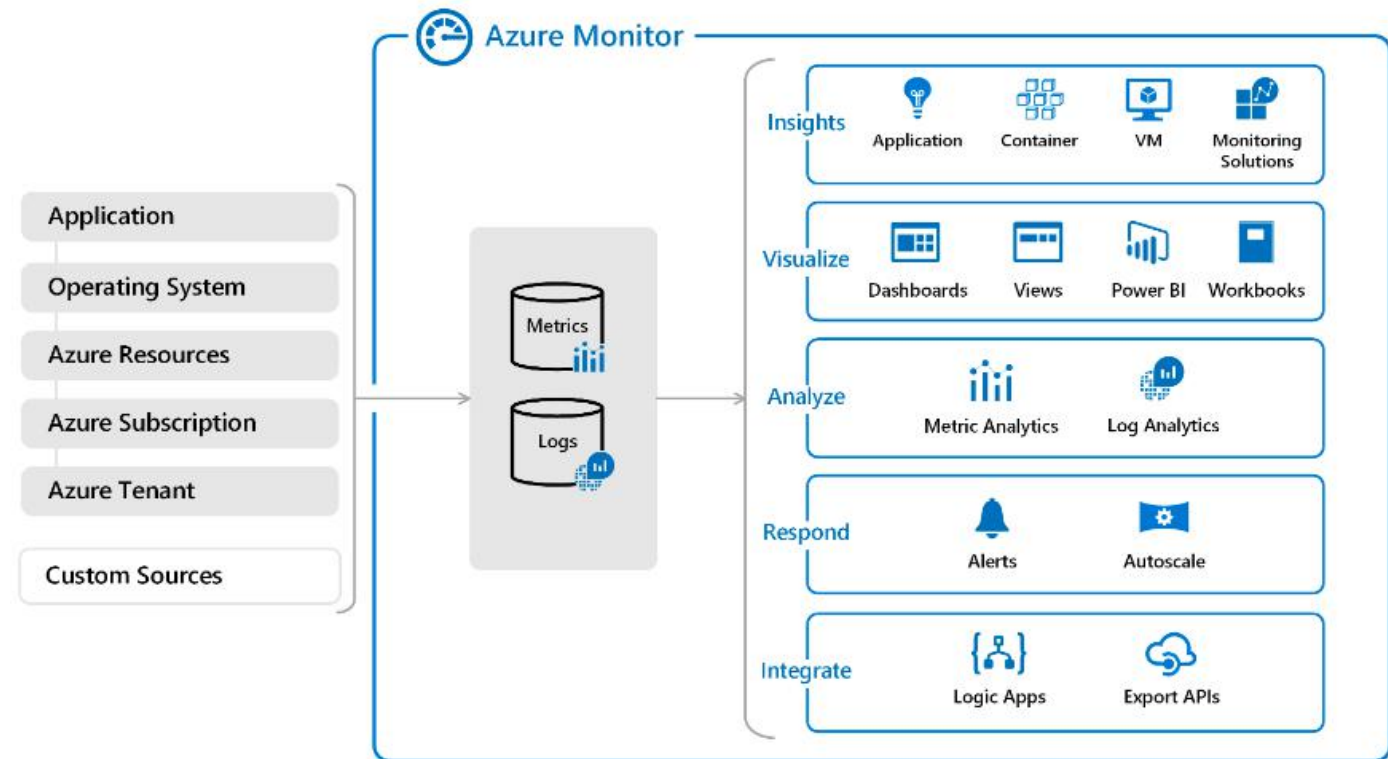
# CLOUDWATCH ALARMS

# SUPPORTED AWS SERVICES

# AZURE MONITOR

**Azure Monitor**

- A platform for collecting, analyzing, visualizing,
- Taking action based on the metric and logging data from your entire Azure and on-premises environment

**Part of a suite**

- Azure Advisor
- Azure Monitor
- Azure Service Health

# QUIZ

You want to be alerted when new recommendations to improve your cloud environment are available. Which service will do this?

- A. Azure Advisor
- B. Azure Monitor
- C. Azure Service Health

# QUIZ

**Which service provides official outage root cause analyses (RCAs) for Azure incidents?**

- A. Azure Advisor
- B. Azure Monitor
- C. Azure Service Health

# QUIZ

**Which service is a platform that powers Application Insights, monitoring for VMs, containers, and Kubernetes?**

- A. Azure Advisor
- B. Azure Monitor
- C. Azure Service Health

# BACKGROUND SLIDES

# SRE

**SRE stands for Site Reliability Engineering**

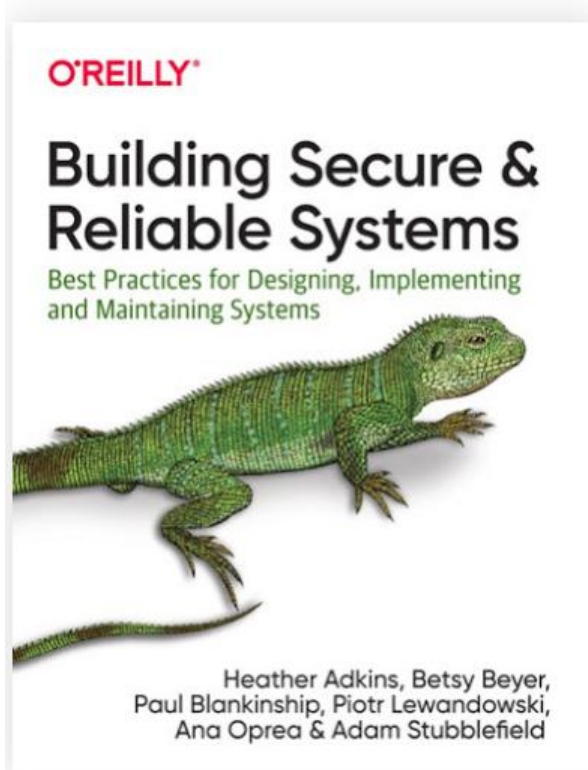**SRE is an area that was recently developed at Google**

**It studies stability and scalability**

**Its tenets are**

- A significant portion of a software system's lifespan is spent in use, not in design or implementation.
- This calls in question the conventional wisdom that insists that software engineers focus primarily on the design and development of large-scale computing systems?
- SRE teaches how and why the DevOps commitment to the entire lifecycle can enable the company to successfully build, deploy, monitor, and maintain some of the largest software systems in the world.
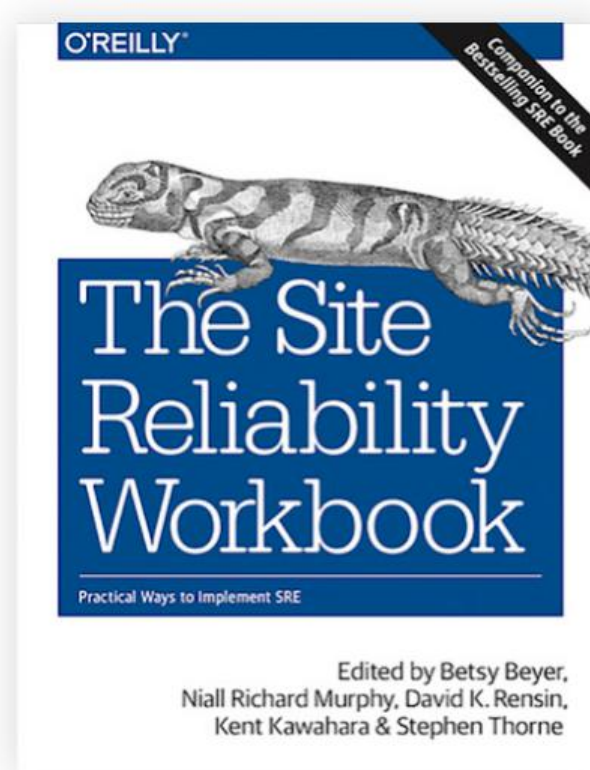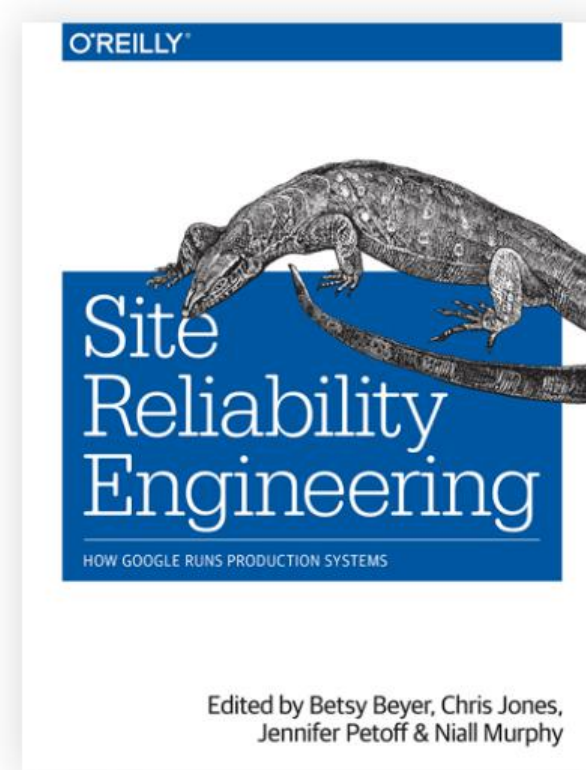
# SRE BOOKS

**Free on Google site here**

# SRE 1

**Introduction**

- The Sysadmin Approach to Service Management
- Google's Approach to Service Management: Site Reliability Engineering
- Tenets of SRE
- Demand Forecasting and Capacity Planning
- Efficiency and Performance

ADS **SRE**

# SRE 2

**Principles**

- Embracing Risk
- Managing Risk
- Motivation for Error Budgets
- Benefits

**Service Level Objectives**

- Service Level Terminology
- Indicators in Practice
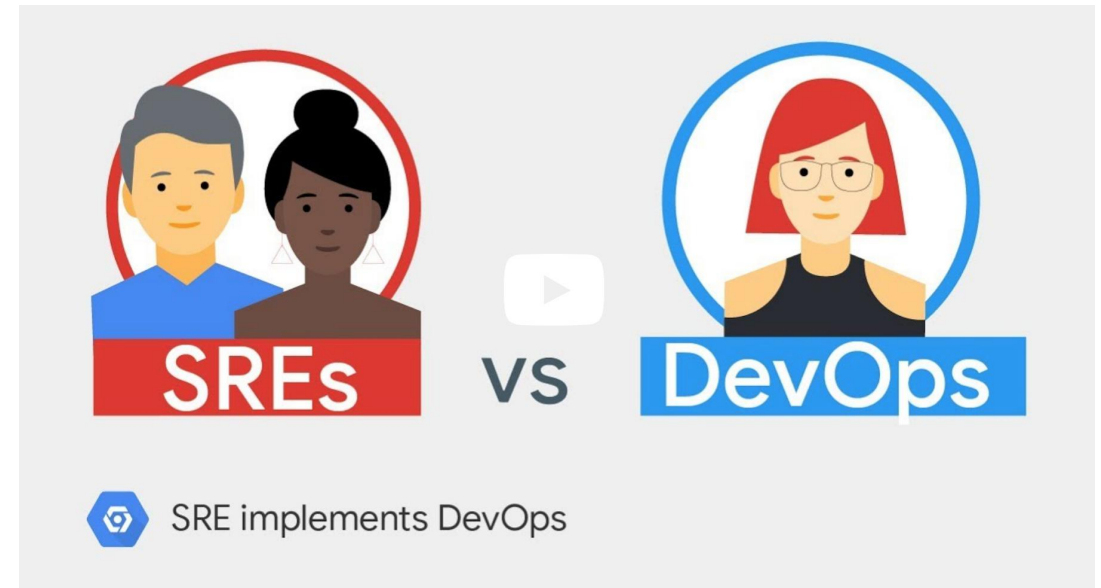- What Do You and Your Users Care About?
- Agreements in Practice

# SRE 3

**Eliminating Toil**

- Toil Defined
- Why Less Toil Is Better

**Monitoring Distributed Systems**

- Why Monitor?
- Setting Reasonable Expectations for Monitoring
- Symptoms Versus Causes
- Black-Box Versus White-Box
- As Simple as Possible, No Simpler
- Bigtable SRE: A Tale of Over-Alerting
- Gmail: Predictable, Scriptable Responses from Humans



SREs VS DevOps

SRE implements DevOps

# CONGRATS ON COMPLETION