# AKAMAS

# Let the machines optimize the machines: ML-driven automated performance tuning

## Conf42 SRE

### Stefano Doni (Akamas)

AKAMAS

# Agenda

**1**  **Why SREs should care about system configurations**

**2**  **A new approach: ML-driven performance tuning**

**3**  **Real-world example: optimize Kubernetes and JVM**

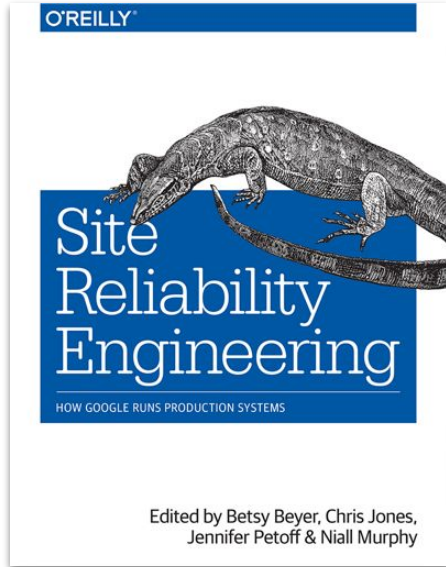**4**  **Conclusions**

**Stefano Doni**

CTO at Akamas

15 years in performance engineering

2015 CMG Best Paper Award Winner

ΛKΛMΛS

# Why SREs should care about system configurations

AKAMAS

# SREs care about efficiency and performance

*"an **SRE team** is responsible for the availability, latency, performance, efficiency, change management, monitoring, emergency response, and capacity planning of their service(s)"*
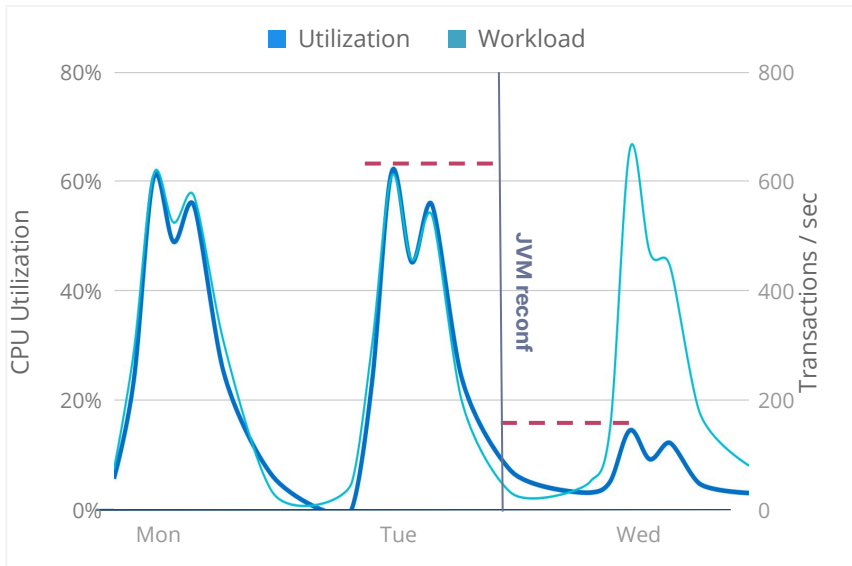
The **core SRE tenets** include:
- Pursuing maximum change velocity without violating SLOs
- Demand Forecasting and Capacity Planning
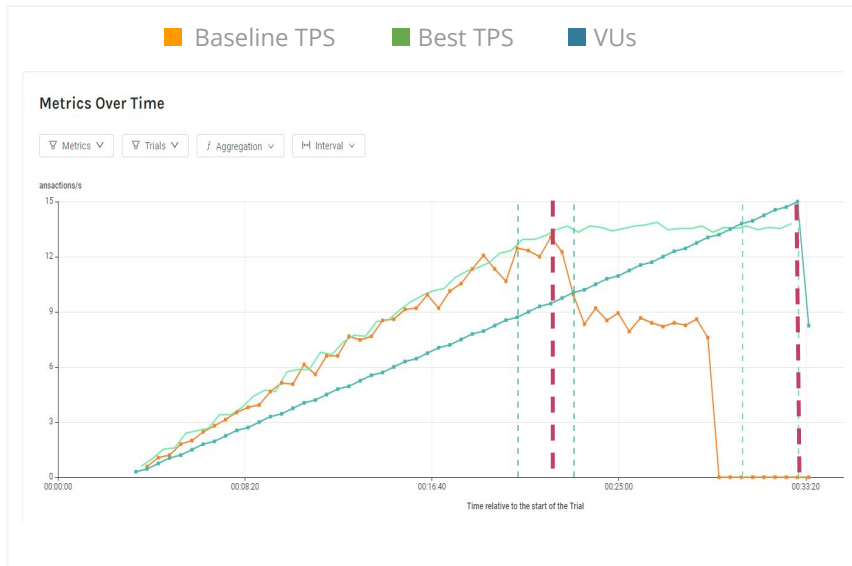- Efficiency and performance

https://sre.google/books

AKAMAS

# Tuning system configuration matters...

**performance and efficiency**



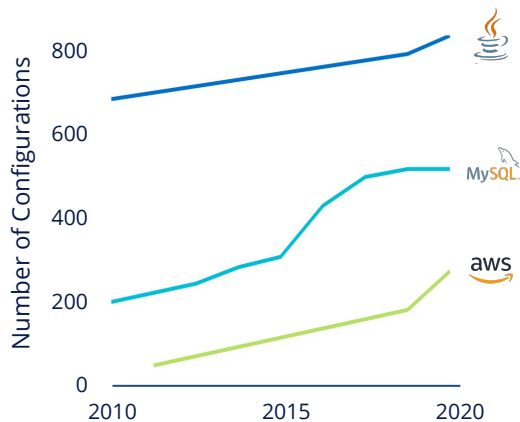*higher application performance and lower infrastructure cost*

**... and service availability**



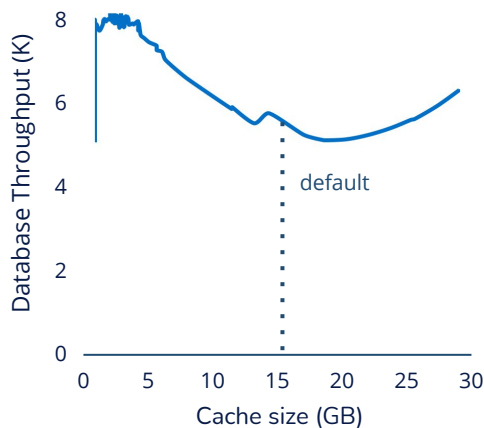*higher transaction throughput and improved service resilience*

AKAMAS

# ... but it is getting harder and harder
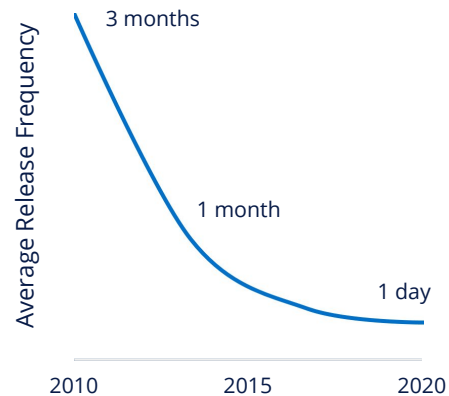
### Configuration Explosion



*properly configuring the IT stack requires analyzing thousands of configurations*

### Unpredictable Effects



*effect of changes can be counterintuitive + default values not always appropriate*

### Faster Deployments



*acceleration of release pace makes manual approach infeasible/useless*

ΛKΛMΛS

# A new approach:
# ML-driven performance tuning

AKAMAS

# Akamas ML-driven optimization
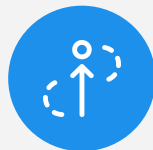


## AKAMAS

### Full-Stack

**any application, any middleware, any database, any cloud - any system**
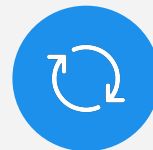
### Smart Exploration

**patented ML quickly identifying optimal configurations beyond any manual tuning**
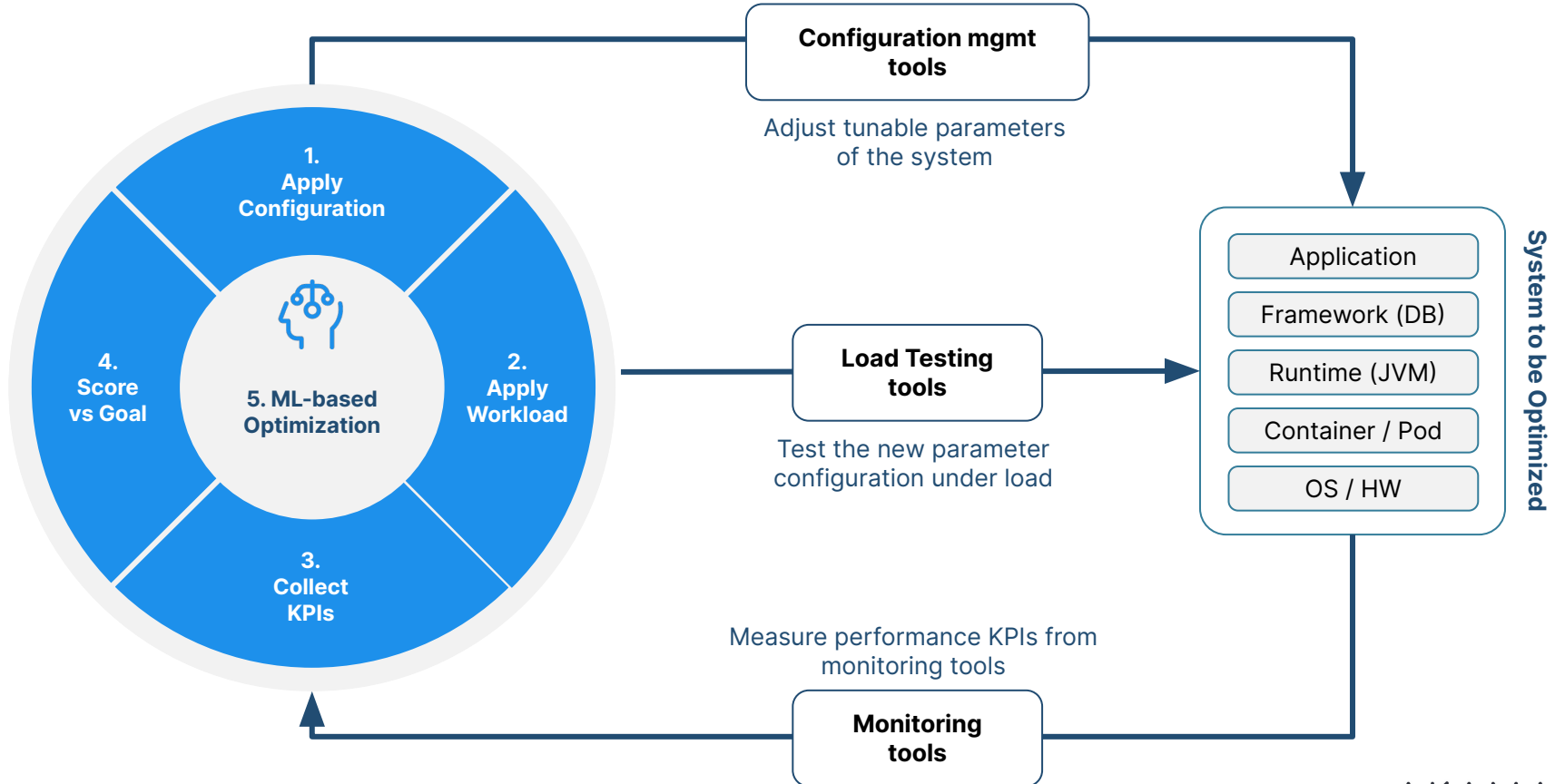
### Goal-oriented

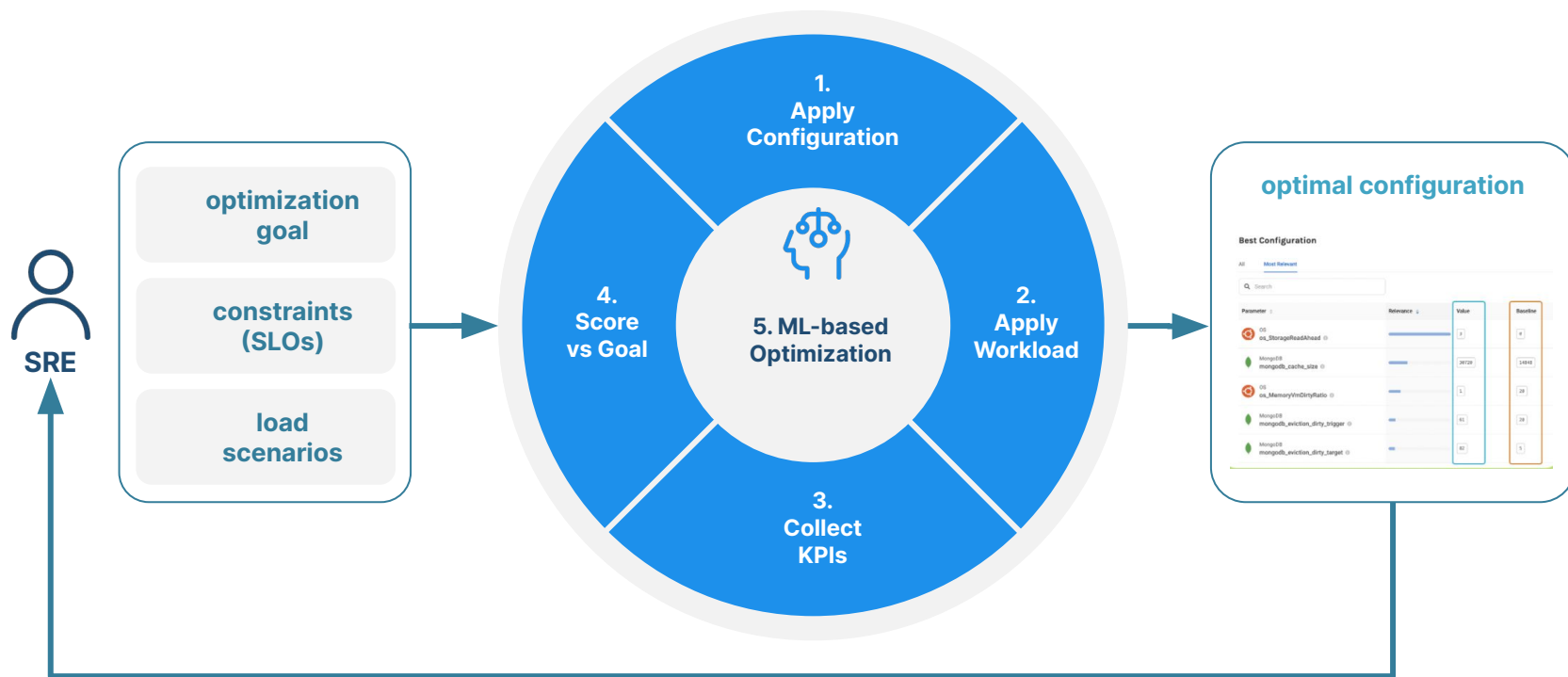**custom-defined goals translating SLOs and business & technical constraints**

### Fully Automated

**custom workflows automating parameter changes, load testing & telemetry collection**
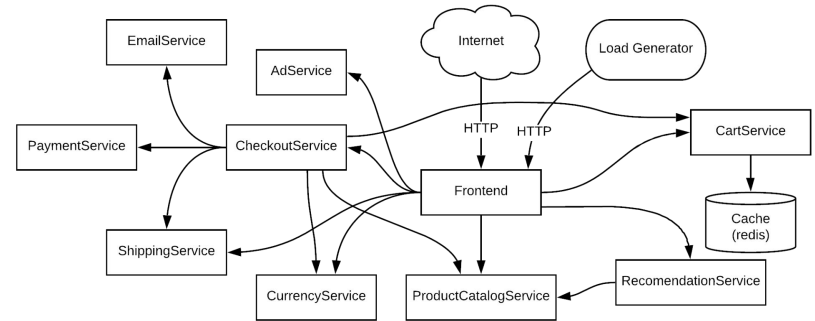
AKAMAS

# ML enables automated performance tuning...



Configuration mgmt tools

Adjust tunable parameters of the system

1. Apply Configuration

5. ML-based Optimization

2. Apply Workload

3. Collect KPIs

4. Score vs Goal

Load Testing tools

Test the new parameter configuration under load

System to be Optimized

- Application
- Framework (DB)
- Runtime (JVM)
- Container / Pod
- OS / HW

Measure performance KPIs from monitoring tools

Monitoring tools

AKAMAS

# ... and a new performance tuning process

AKAMAS

# Real world example: optimize Kubernetes and JVM

ΛΚΛΜΛS

# The target system: Online Boutique

- **Cloud-native application** by Google made of **10 microservices**

- Realistic sample web-based **e-commerce service**

- Features a **modern software stack** (Go, Node.js, Java, Python, Redis)

- Includes a Load Generator (Locust) to inject **realistic workloads**

https://github.com/GoogleCloudPlatform/microservices-demo

AKAMAS

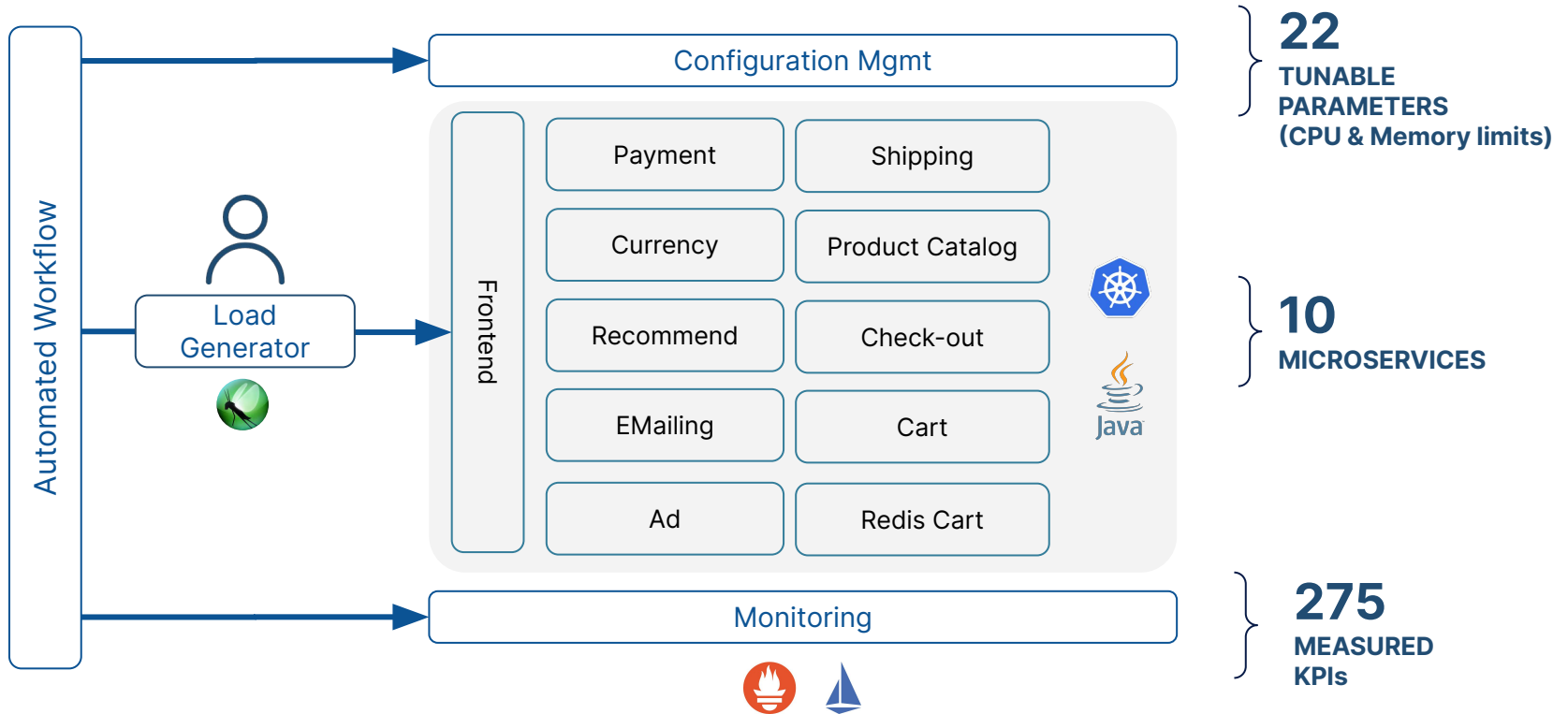# Use Case: optimizing cost of K8s microservices while ensuring reliability

## Challenge for SRE

How to provision the optimal resources to your application made of several

**Kubernetes** microservices, so that you can trust the overall service

- ➔ will sustain the expected **target load**
- ➔ while matching the defined **Service-Level Objectives** (SLOs)
- ➔ at the **minimum cost**
- ➔ while minimizing the operational effort
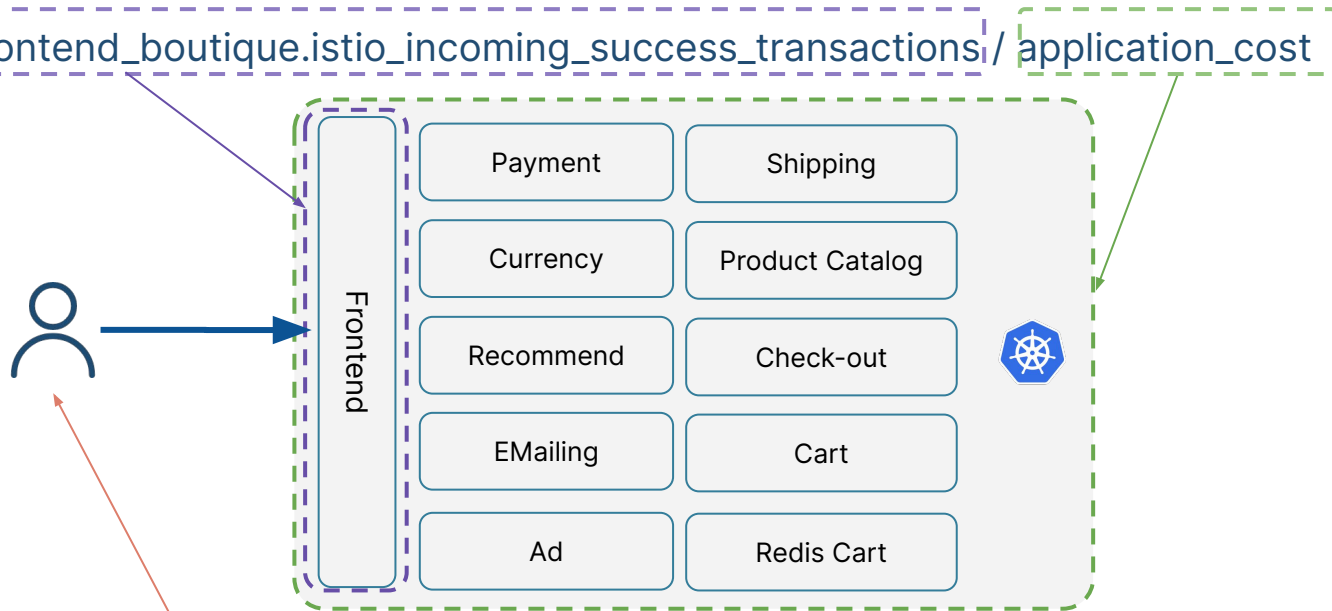- ➔ and matching delivery milestones

**SRE**

AKAMAS

# The reference architecture



**22** TUNABLE PARAMETERS (CPU & Memory limits)

**10** MICROSERVICES

**275** MEASURED KPIs

Configuration Mgmt

Automated Workflow

Load Generator

Frontend

Payment | Shipping
Currency | Product Catalog
Recommend | Check-out
EMailing | Cart
Ad | Redis Cart

Monitoring

AKAMAS

# The optimization goals & constraints

**GOAL:**
**MAXIMIZE** frontend_boutique.istio_incoming_success_transactions / application_cost

| Frontend | Payment | Shipping |
|---|---|---|
| | Currency | Product Catalog |
| | Recommend | Check-out |
| | EMailing | Cart |
| | Ad | Redis Cart |

**CONSTRAINTS:** loadgenerator_locust.locust_fail_ratio <= 2% **AND**
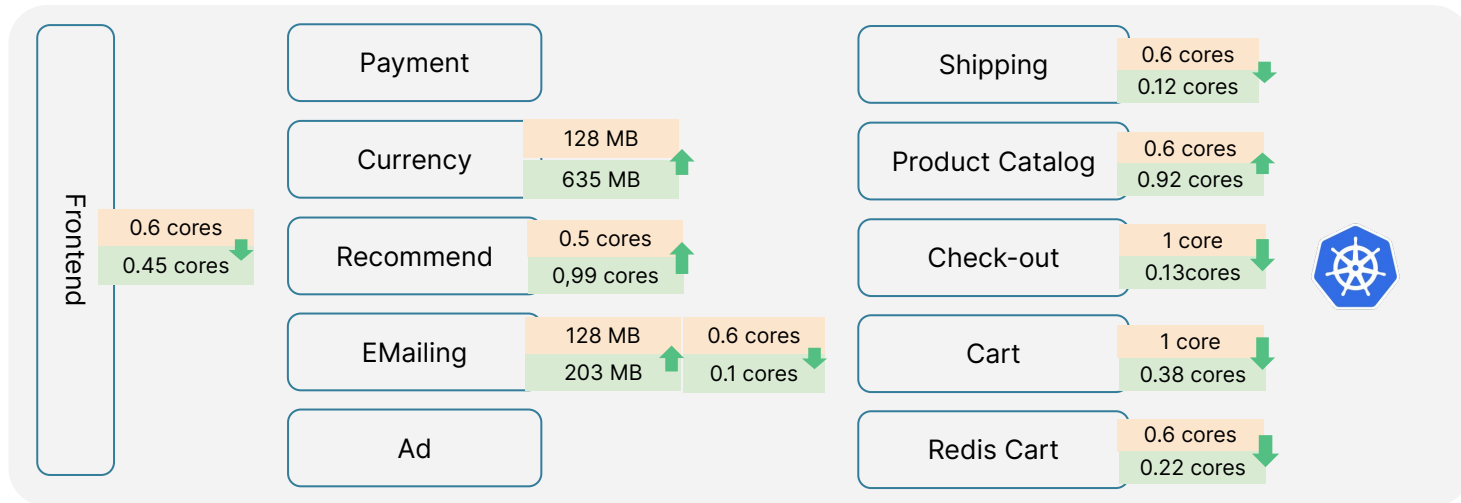frontend_boutique.istio_incoming_response_time_90pct <= 500ms

AKAMAS

# Best configuration found by ML in 24H improves cost efficiency by 77%



© 2021 Akamas • All Rights Reserved • Confidential

ΛKΛMΛS

# Best config: optimal resources assigned to microservices

**10**
**TOP IMPACT PARAMETERS**

Baseline

Best

Frontend

0.6 cores
0.45 cores

Payment

Currency
128 MB
635 MB

Recommend
0.5 cores
0,99 cores

EMailing
128 MB
203 MB
0.6 cores
0.1 cores

Ad

Shipping
0.6 cores
0.12 cores

Product Catalog
0.6 cores
0.92 cores

Check-out
1 core
0.13cores

Cart
1 core
0.38 cores

Redis Cart
0.6 cores
0.22 cores

- decreased CPU limits set for almost all containers
- increased CPU assigned to 2 microservices
- all these changes to achieve max cost efficiency and match SLOs

ΛKΛMΛS

# Best config: higher performance & efficiency for the overall service

## Baseline vs Best: Service throughput



**+19%**
**TPS**

Baseline, Trial 1, frontend_boutique.istio_incoming_success_transactions
Exp 31, Trial 1, frontend_boutique.istio_incoming_success_transactions

## Baseline vs Best: Service p90 response time



**-60%**
**Response Time**

Baseline, Trial 1, frontend_boutique.istio_incoming_response_time_90_ms
Exp 31, Trial 1, frontend_boutique.istio_incoming_response_time_90_ms

AKAMAS

# Use Case: maximizing service performance & efficiency with JVM tuning

**SRE**

## Challenge for SRE

How to ensure a reliable product launch, by properly configuring JVM options,

so that you can trust the overall service

- will sustain the expected **target load**

- while matching the defined **Service-Level Objectives** (SLO)

- at the **minimum cost**

- while minimizing the operational effort

- and staying aligned product launch milestones

AKAMAS

# The reference architecture

# The optimization goals & constraints

**GOAL:**
**MAXIMIZE** `ad.istio_incoming_success_transactions`

Frontend

| Payment | Shipping |
| Currency | Product Catalog |
| Recommend | Check-out |
| EMailing | Cart |
| Ad | Redis Cart |

**CONSTRAINTS:** `ad.transaction_response_time <= 100ms`

AKAMAS

# Best config:
# +28% throughput, and meeting SLOs



Best configuration
**+28%**
*Peak Throughput matching SLO:* **95** *TPS*

Baseline configuration

*Peak Throughput matching SLO:* **74** *TPS*

***SLO*** *breaking at 100ms*

AKAMAS

# Best config: optimal JVM options

## 8
### TOP IMPACT PARAMETERS

| Parameter | Relevance | Best | Baseline |
|---|---|---|---|
| jvm jvm_newSize | | 550 MB (+83.3%) | 300 MB |
| jvm jvm_GCTimeRatio | | 100 (+1%) | 99 |
| jvm jvm_concurrentGCThreads | | 1 threads (-87.5%) | 8 threads |
| jvm jvm_gcType | | Parallel | G1 |
| jvm jvm_maxHeapSize | | 901 MB (+252%) | 256 MB |
| jvm jvm_maxTenuringThreshold | | 6 (-60%) | 15 |
| jvm jvm_parallelGCThreads | | 3 threads (-62.5%) | 8 threads |
| jvm jvm_survivorRatio | | 100 (+1,150%) | 8 |

- increased max heap memory
- changed Garbage Collector type
- decreased number of Garbage Collector threads
- adjusted heap regions & object aging thresholds

AKAMAS

# Conclusions

AKAMAS

# Key takeaways

**Tuning modern applications** for increasing their efficiency, performance and reliability is a **complex problem** that represent a **relevant toil** for SRE teams

A new approach leveraging fully-automated **ML-based optimization** enables SRE teams to ensure applications will have **higher performance & reliability**

Leveraging this new **ML-based optimization** approach, SRE teams can **reduce the operational toil** and **stay aligned to release milestones**

ΛKΛMΛS

# Contacts

**Italy HQ**
Via Schiaffino 11
Milan, 20158
+39-02-4951-7001

**USA East**
211 Congress Street
Boston, MA 02110
+1-617-936-0212

**USA West**
12655 W. Jefferson Blvd
Los Angeles, CA 90066
+1-323-524-0524

**Singapore**
5 Temasek Blvd
Singapore 038985

**LinkedIn**
@akamaslabs

**Twitter**
@AkamasLabs

**Email**
info@akamas.io

AKAMAS

# BACKUP SLIDES