# Computing Paradigms for the Next-Generation Computing Landscapes

The 2nd AIoTwin Summer School

Ass. Prof. Dr. Stefan Nastic

17th September 2024

# About me

- ▶ Assistant Professor at Distributed Systems Group at TU Wien
- ▶ Founder and CEO at IntelliEdge GmbH
- ▶ Leading various commercial and research projects over the past 10+ years
- ▶ My core research interests
  - ▶ Serverless Computing
  - ▶ Edge-Cloud Continuum
  - ▶ AI & Edge AI
- ▶ Get in touch via LinkedIn or find me at nastic.at

dsg | DISTRIBUTED SYSTEMS GROUP

IntelliEdge
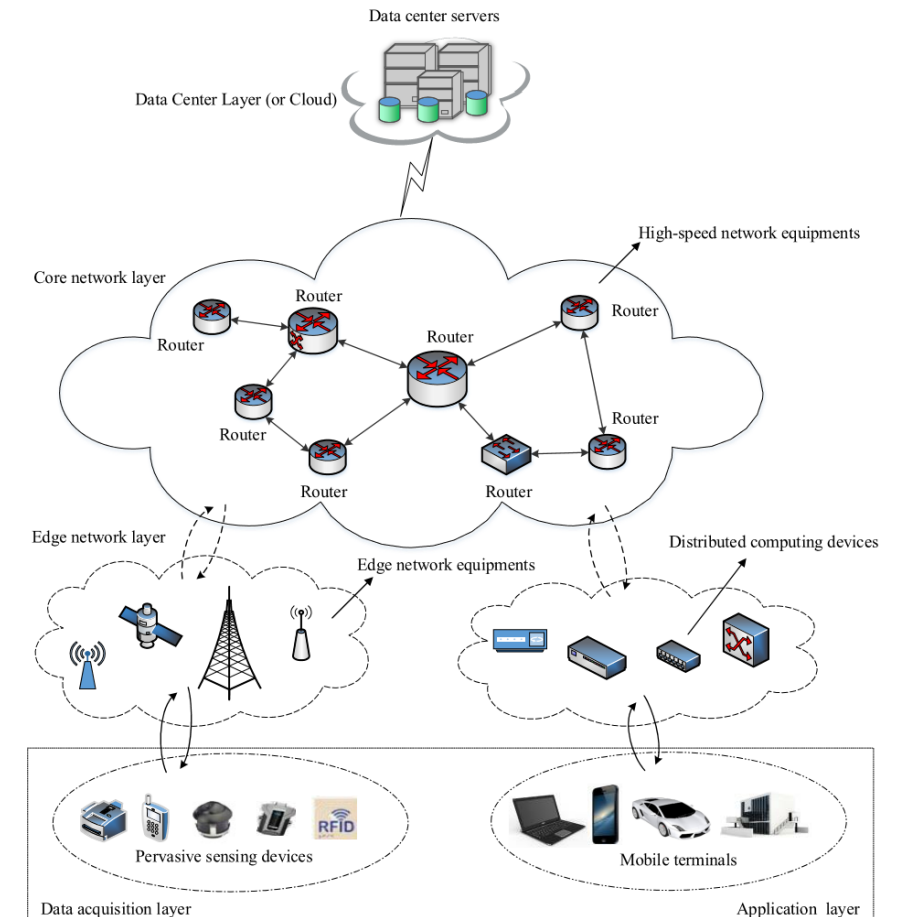
# Edge-Cloud Continuum in a Nutshell

🔊 **continuum**

/kənˈtɪnjʊəm/

*noun*

a continuous sequence in which adjacent elements are not perceptibly different from each other, but the extremes are quite distinct.
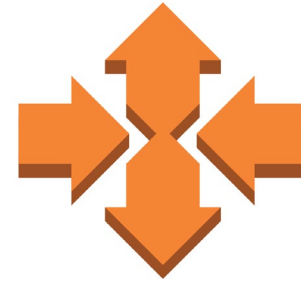"a continuum of special educational needs"

▶ A large pool of interconnected resources

▶ Usually hierarchical/layered and globally heterogeneous

# The Promises of Serverless

**No infrastructure management**
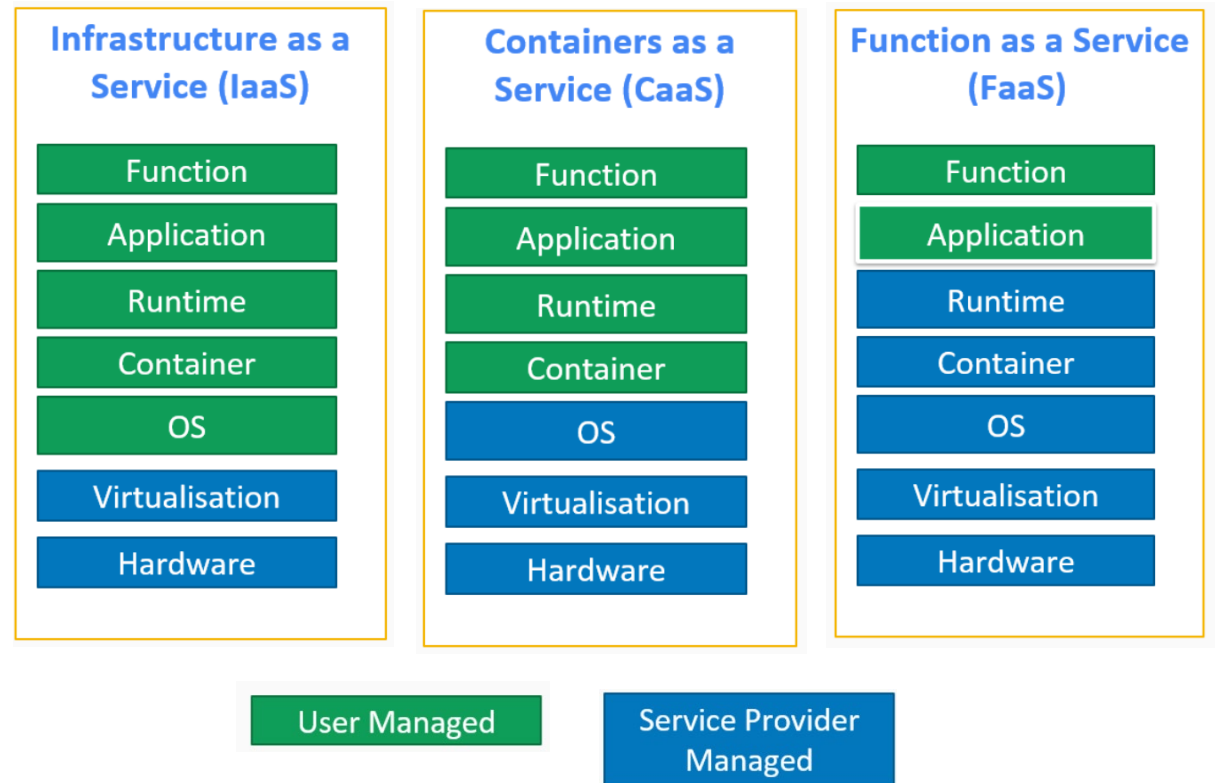
**Flexible elastic scaling**

**High reliability & availability**

**No idle capacity**

# Serverless Computing in a Nutshell

- Serverless ≠ No Servers
- Natural evolution in Cloud Computing
- Novel Function as a Service (FaaS) Paradigm
- Novel execution model

### Infrastructure as a Service (IaaS)

| Function |
| --- |
| Application |
| Runtime |
| Container |
| OS |
| Virtualisation |
| Hardware |

### Containers as a Service (CaaS)

| Function |
| --- |
| Application |
| Runtime |
| Container |
| OS |
| Virtualisation |
| Hardware |

### Function as a Service (FaaS)

| Function |
| --- |
| Application |
| Runtime |
| Container |
| OS |
| Virtualisation |
| Hardware |

User Managed

Service Provider Managed

# Function Isolation and Virtualization Models

- ▶ Containers (K-Native on Kubernetes )

- ▶ MicroVMs - Firecracker (AWS Lambda)

- ▶ V8 Isolates (Cloudflare, Deno Deploy)

- ▶ WASM runtime (Fastly)
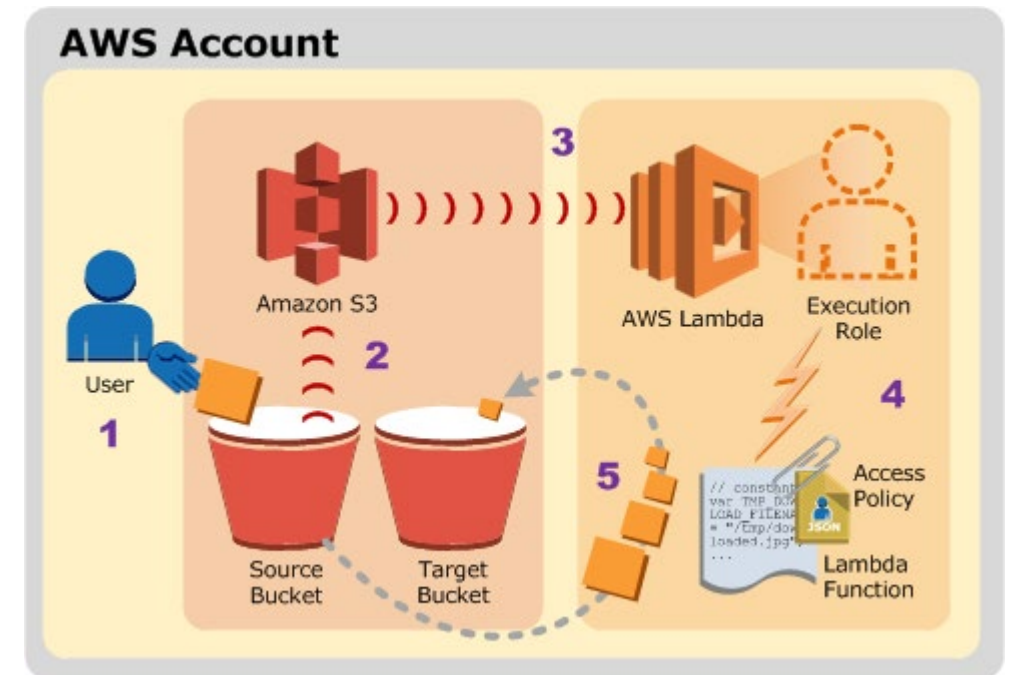
- ▶ Unikernels (Mainly academia)

# Function as a Service (FaaS)

- A type of serverless computing where developers can deploy and run code in response to events or requests

- Application is decomposed into "triggers" and "actions" (functions)
  - Developers write functions and upload them to a FaaS platform
  - Functions are triggered by events, such as a user request or a change in data

- Example: Thumbnail image processing

# Popular Serverless Platforms

- ▶ Cloud-First
  - AWS Lambda
  - Azure Functions
  - Google Cloud Functions
- ▶ Edge-First
  - Cloudflare Workers
  - Fastly Compute@Edge
- ▶ Open source
  - OpenWhisk
  - OpenFass

# Serverless Beyond the FaaS

## Serverless "support BaaS services"

- API Gateways for synchronous (request-based) communication

- Object stores for persistent state management

- Database change streams for reacting to data changes (DynamoDB Streams)

- Function orchestration (Step Function)

**T. Larcher, P. Gritsch, S. Nastic and S. Ristov, "BAASLESS: Backend-as-a-Service (BaaS)-Enabled Workflows in Federated Serverless Infrastructures," in IEEE Transactions on Cloud Computing**

## Serverless managed Cloud services

- Serverless relational database (Amazon Aurora)

- Serverless container management (AWS Fargate)

- Serverless warehouse (Google BigQuery)

# Serverless Beyond the Cloud

- Serverless was born in the Cloud and remains largely limited to the Cloud
- True potential of Serverless is unlocked in the *Edge-Cloud continuum*
  - Localized computation (Edge locations vs datacenters)
  - Low network latency due to proximity to the user
  - Flexible geo-restrictions, compliance, etc.

## Serverless Compute Fabric

- Edge-Cloud native BaaS services
- Distributed and decentralized messaging solutions (Emma)
- Lightweight and secure virtualization and isolation (WASM)
- Edge-native state management (Durable Objects)

Nastic, S., Dustdar, S., Philipp, R., Alireza, F., & Pusztai, T. (2022). A Serverless Computing Fabric for Edge & Cloud. In 4th IEEE International Conference on Cognitive Machine Intelligence (CogMi).

# The Promises of Serverless

**No infrastructure management**

**Flexible elastic scaling**

**High reliability & availability**

**No idle capacity**

# Many Opportunities …

- Fine-grained pay-per-use model with preserved SLOs

- Optimal resource usage with Scale-to-Zero

- Rapid elasticity & auto-scaling due to lightweight virtualization

- Native cloud offloading

- True utility-based resource consumption

- Disaggregation of compute and storage

- Stateless workflows and life-cycle

- Effortlessly parallelizing computing

- Extremely bursty workloads

- CUP-bound applications

- Developing scalable Edge-native solution

- True event-driven applications

- ….

# Serverless Horror Stories
## Examples from the Wild: A runaway function

- ► What is a runaway function?
- ► There are no default safeguards in place to prevent runaway functions

**Beware "RunOnStartup" in Azure Functions – a serverless horror story**

By Tom in Uncategorised

📅 6th September 2018    💬 6 Comments

The curious case of the spiralling AWS Lambda bill

We Burnt $72K testing Firebase + Cloud Run and almost went Bankrupt [Part 1]

**Google Cloud**
**Summary for Mar 1, 2020–Mar 31, 2020**

| | |
|---|---|
| Starting balance as of Mar 1, 2020 | $0.00 |
| Total new activity | $71,393.86 |
| Total payments received | $0.00 |
| Ending balance in USD | $71,393.86 |

. But sometimes costs go wrong.
ch situation - ever increasing costs,
uled to execute once per hour!

Sudeep Chauhan
About me
Dec 08, 2020 · 9 mins read

Share this!

" *T* *his is the story of how close we came to shutting down before even launching our first product, how we survived, and the lessons we learnt.*

In March, 2020, when COVID hit the world, our startup Milkie Way too was hit with a big blow and almost shut down. **We burnt $72,000 while exploring and internally testing Cloud Run with Firebase within a few hours.**

# Serverless Horror Stories
## Examples from the Wild: A service misconfiguration

- Serverless functions don't live in isolation ("support services")
- Serverless says nothing about configuring & managing other edge-cloud services
- Arguably more complex ecosystem

### Serverless: 15% slower and 8x more expensive
Posted: 2019-09-18 Last updated: 2019-09-26

...ging the API we have at CardGames.io and try using ...verless has been a hot topic in the tech world for the ...stinating wanted to keep my tech skills up to date

### Lessons learned from combining SQS and Lambda in a data project

In June 2018, AWS Lambda added Amazon Simple Queue Service (SQS) to supported event sources, removing a lot of heavy lifting of running a polling service or creating extra SQS to SNS mappings. In a recent project we utilized this functionality and configured our data pipelines to use AWS Lambda functions for processing the incoming data items and SQS queues for buffering them. The built-in functionality of SQS and Lambda provided us serverless, scalable and fault-tolerant basis, but while running the solution we also learned some important lessons. In this blog post I will discuss the issue of valid messages ending up in dead-letter queues (DLQ) and correctly configuring your DLQ to catch only erroneous messages from your source SQS queue.

#### ꞔ segment
Popular   Podcasts   Growth & Marketing   Engineering   Company

Product   Pricing   Customers   Docs   Company

### The million dollar engineering problem

Achille Roussel, Rick Branson on March 14th 2017

For an early startup, using the cloud isn't even a question these days. No RFPs, provisioning orders, or physical shipments of servers. Just the promise of getting up and running on "infinitely scalable" compute power within minutes.

But, the ability to provision thousands of dollars worth of infrastructure with a single API call comes with a *very large hidden cost*. And it's something you won't find on any pricing page.

# Main Challenges with Serverless in Edge-Cloud Continuum

## Performance Challenges

- Startup/scheduling latency ("cold start")

- Lack of performance isolation ("noisy neighbors")

- Inconsistent performance due to hardware heterogeneity (also present in Cloud only solutions!)

## Data Management Challenges

- Function execution latency due to lack of data locality

- State management for stateful computations

- Efficient and cost-effective caching solutions

- Particularly important for emerging Edge AI workloads

Nastic, S., Dustdar, S., Philipp, R., Alireza, F., & Pusztai, T. (2022). A Serverless Computing Fabric for Edge & Cloud. In 4th IEEE International Conference on Cognitive Machine Intelligence (CogMi).

# Main Challenges with Serverless in Edge-Cloud Continuum

## Reliability Engineering Challenges

- Dealing with function failures beyond simple retries

- Network partitioning can render functions useless (detached storage)

- SLO-aware provisioning of functions (beyond memory and CPU)

- Interoperability and portability of functions

## Software Development Challenges

- A misconfiguration of services and resources

- Bad coding practices

- Wrong or incomplete error handling mechanism

- Testing of the functions, beyond the unit tests

- Resource quotas & limits

Nastic, S., Dustdar, S., Philipp, R., Alireza, F., & Pusztai, T. (2022). A Serverless Computing Fabric for Edge
& Cloud. In 4th IEEE International Conference on Cognitive Machine Intelligence (CogMi).

# Project Polaris

- ► Enabling SLO- and reliability- awareness in the Edge-Cloud

- ► Evolved into an ecosystem of tools and framework for the serverless Edge-Cloud continuum

- ► Polaris is hosted by Linux Foundation as a SIG of a broader Project Centaurus

- ► Very active GitHub: https://github.com/polaris-slo-cloud

Nastic, S., Morichetta, A., Pusztai, T., Dustdar, S., Ding, X., Vij, D. and Xiong, Y., 2020. SLOC: Service level objectives for next-generation cloud computing. IEEE Internet Computing, 24(3), pp.39-50.

# Polaris SLO Framework

**Monitoring** → **SLO Script** → **Elastic Strategies**

Pusztai, T., Morichetta, A., Pujol, V.C., Dustdar, S., Nastic, S., Ding, X., Vij, D. and Xiong, Y., 2021, September. SLO script: A novel language for implementing complex cloud-native elasticity-driven SLOs. In 2021 IEEE International Conference on Web Services (ICWS) (pp. 21-31).

# SLO Controllers

- Developed as Monitoring, SLO, and Scaling Controllers
- Reference implementation extends Kubernetes
- Orchestrator/platform independent control loop

Pusztai, T., Morichetta, A., Pujol, V.C., Dustdar, S., Nastic, S., Ding, X., Vij, D. and Xiong, Y., 2021, September. A novel middleware for efficiently implementing complex cloud-native SLOs. In 2021 IEEE 14th International Conference on Cloud Computing (CLOUD) (pp. 410-420).

# Predictive Monitoring Controllers Results



(a) Prediction

(b) Prediction

(c) Prediction for t+3

# SLO User Tool Suite

## SLO Framework UI



## SLO Script CLI



**Tooling demo video on YouTube: https://www.youtube.com/watch?v=JVZ4hB2AmGs**
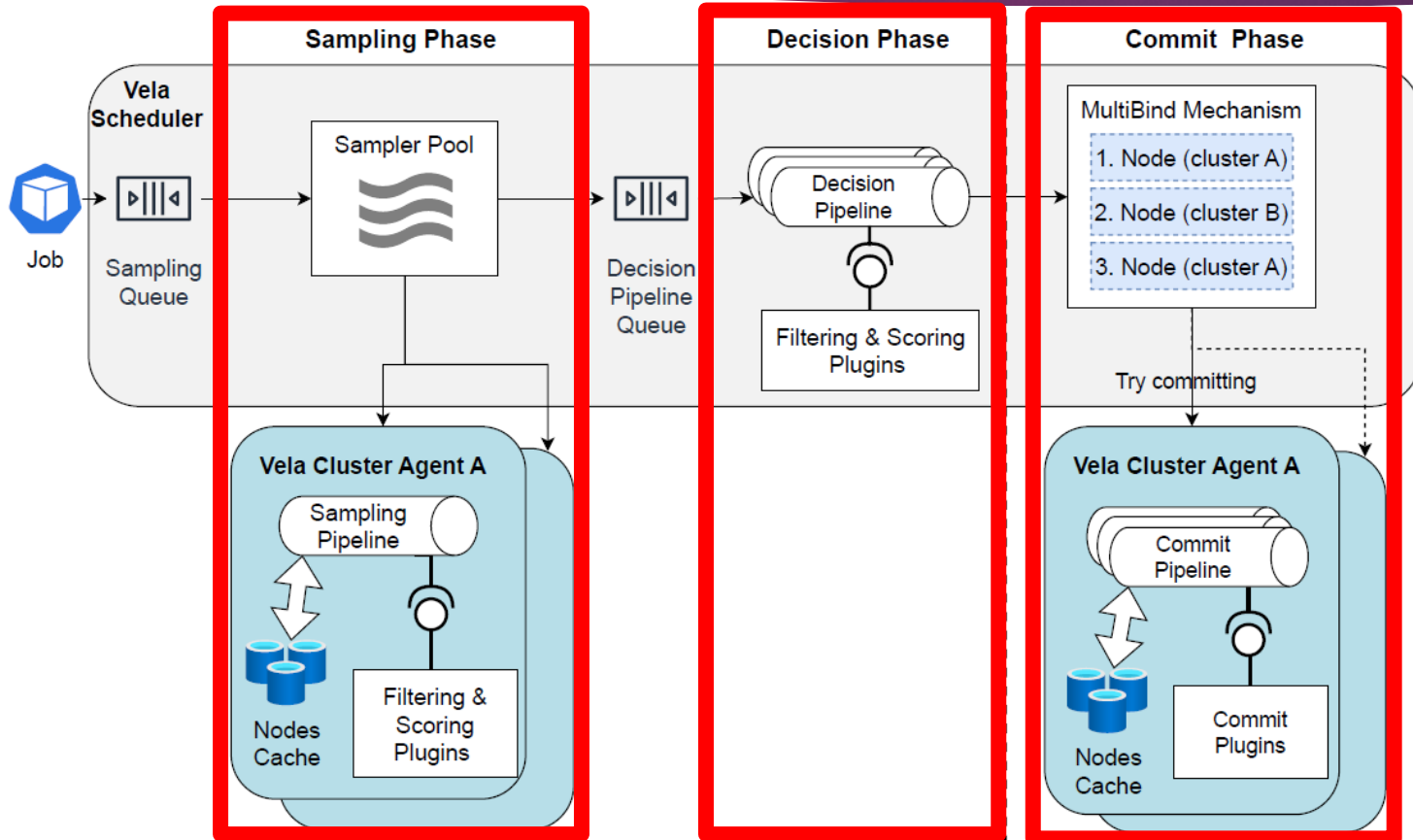
# Polaris Scheduler for Edge-Cloud Continuum

▶ Plugins-based, Edge-sensitive, and SLO-aware scheduling framework for Edge-Cloud Continuum



Latency and bandwidth variance values are used to assess the QoS requirements (stability of the network connections)

Nastic, S., Pusztai, T., Morichetta, A., Casamayor Pujol, V., Dustdar, S., Vii, D., & Xiong, Y. (2021). Polaris scheduler: Edge sensitive and slo aware workload scheduling in cloud-edge-iot clusters. In IEEE 14th International Conference on Cloud Computing (CLOUD) (pp. 206–216).

# Dealing with the Scale of the Edge-Cloud Continuum



- The Edge-Cloud Continuum is very large in scale
- We need to work with a sample of the infrastructure nodes
- We introduced a 3-phase distributed scheduler

Pusztai, T., Nastic, S., Casamayor Pujol, V., Raith, P., Dustdar, S., Vij, D., & Xiong, Y. (2023). Vela: A 3-Phase Distributed Scheduler for the Edge-Cloud Continuum. In The 11th IEEE International Conference on Cloud Engineering (IC2E 2023)

# Scheduler Global Scalability Results

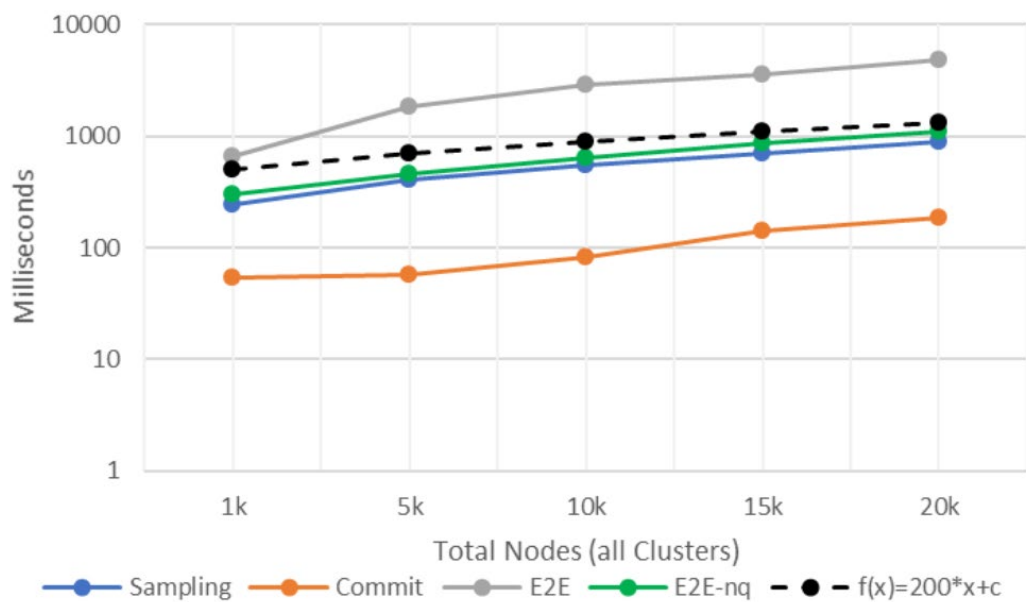## Scalability with Respect to Infrastructure



Fig. 2: Mean Scheduling Times (ms) at `percentageOfClustersToSample` = 50% and `nodesToSampleBp` = 4% for Total Nodes.
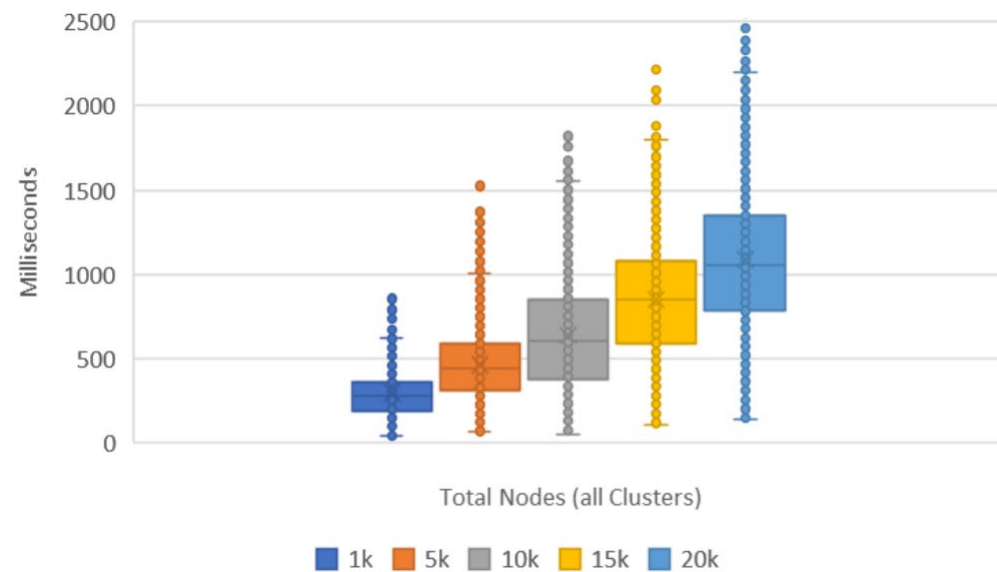


Fig. 5: End-to-End Times (ms), without Sampling Queue at `percentageOfClustersToSample` = 50% and `nodesToSampleBp` = 4% for Total Nodes.

# Scheduler Global Scalability Results

## Scalability with Respect to Workload Throughput
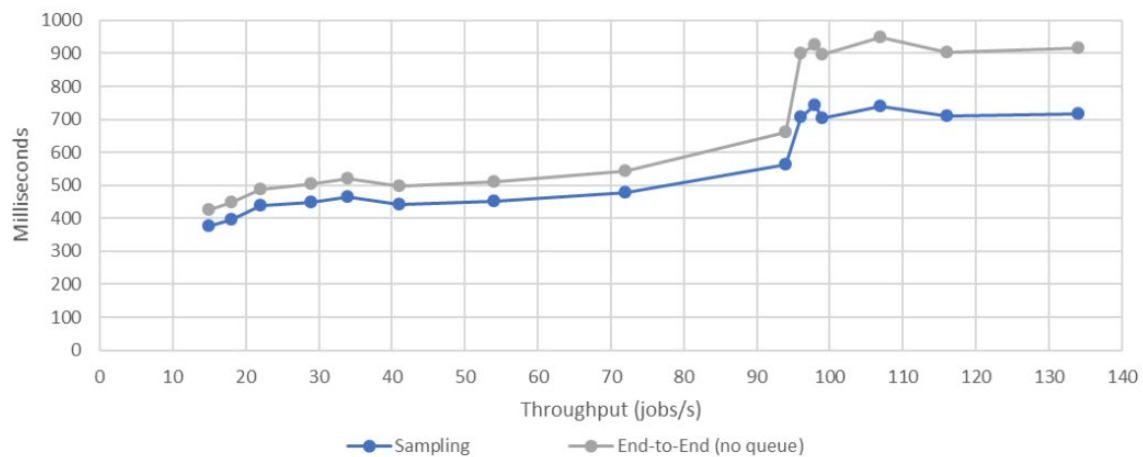


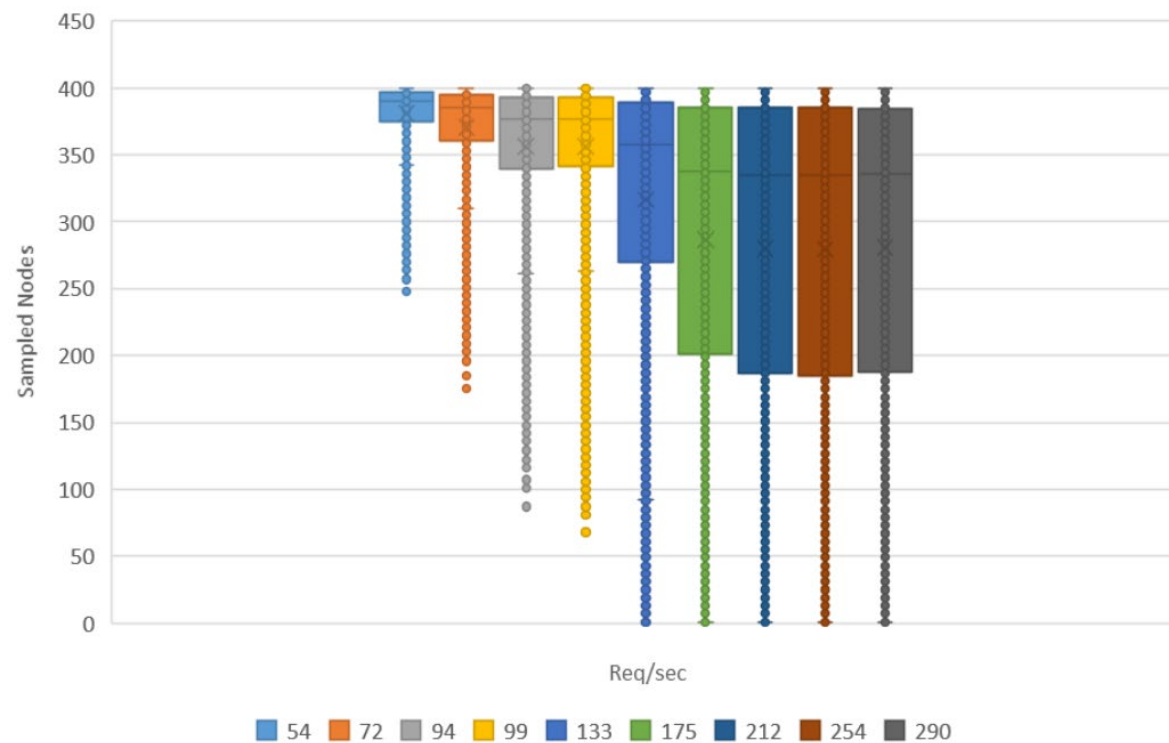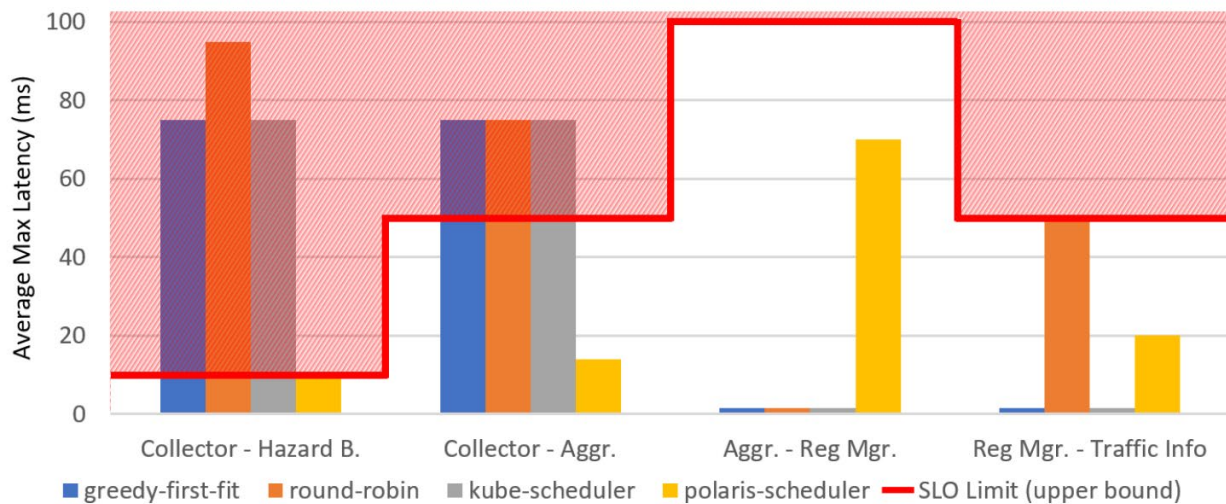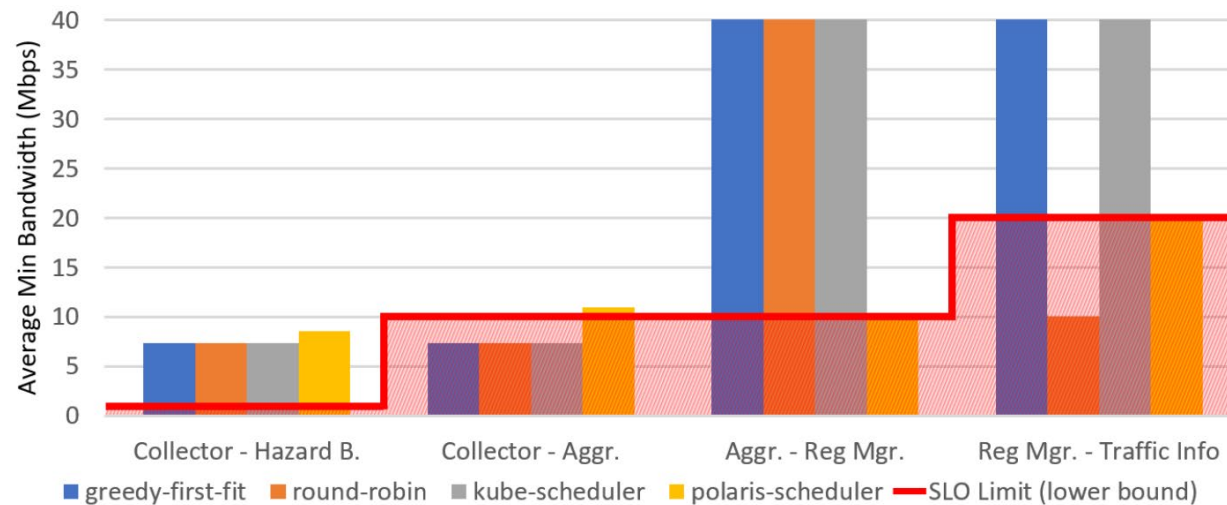Fig. 8: Mean Sampling and E2E-no-queue Times (ms) with Respect to Throughput.
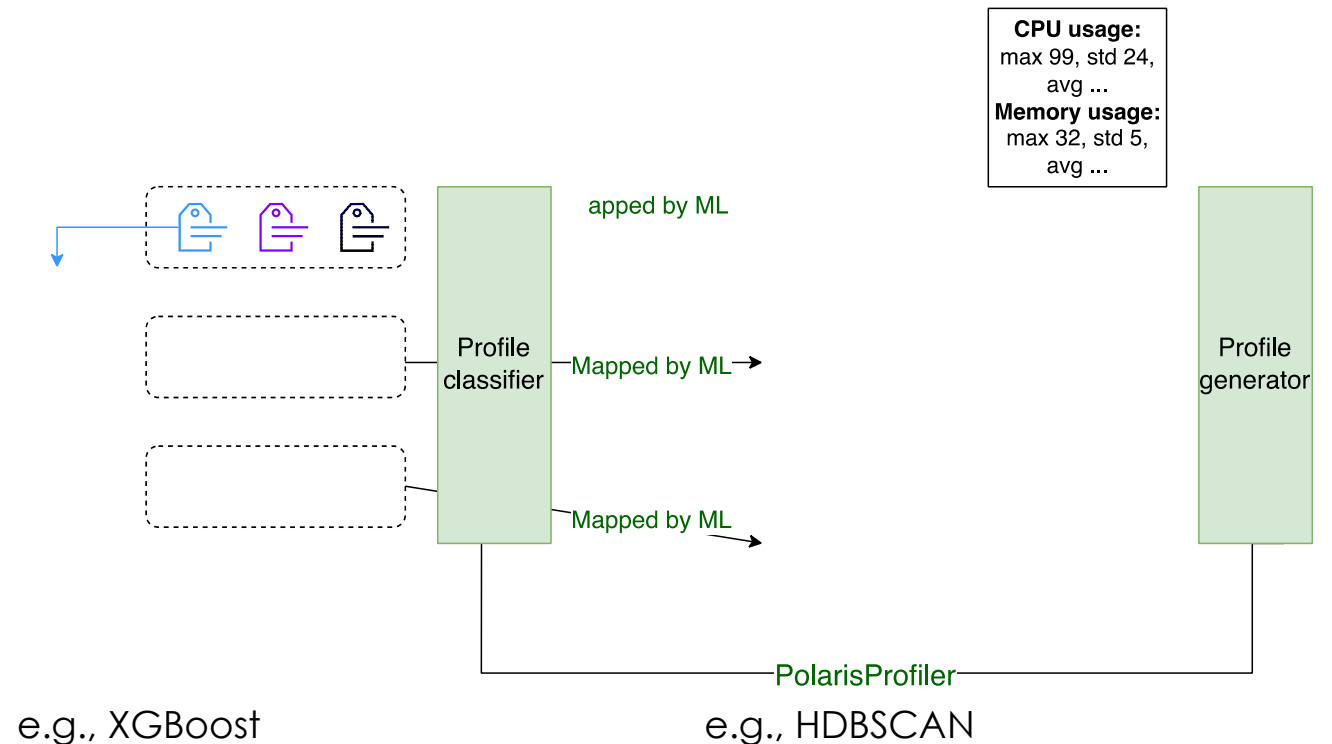


Fig. 7: Sampled Nodes.

# Scheduler SLO Compliance Results



Avg. Max Latencies vs SLO Bounds

Avg. Min Bandwidth vs SLO Bounds

Pusztai, T., Nastic, S., Morichetta, A., Casamayor Pujol, V., Raith, P., Dustdar, S., … Zhang, Z. (2022). Polaris Scheduler: SLO- and Topology-aware Microservices Scheduling at the Edge.
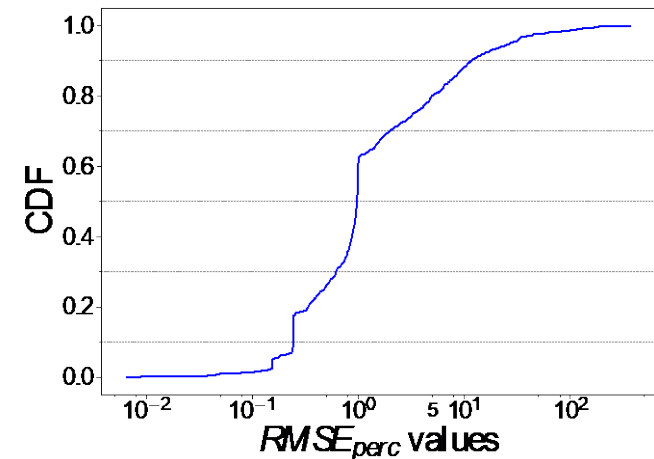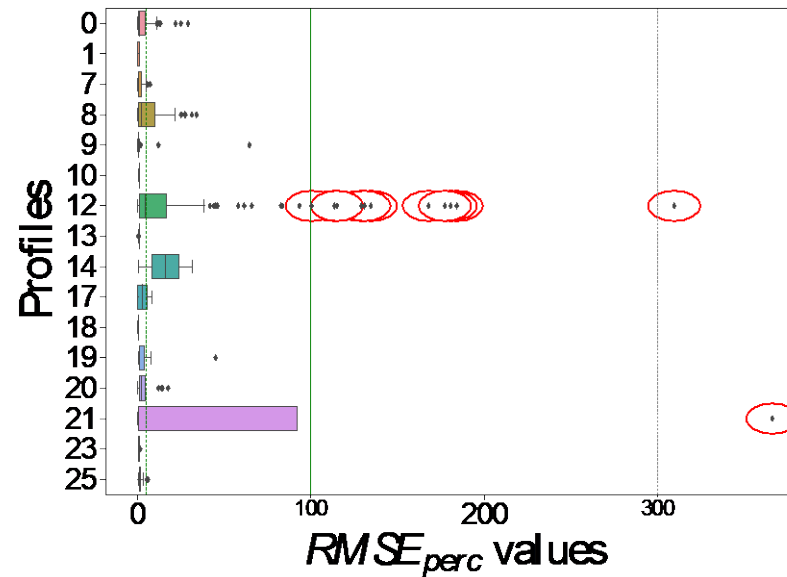
# Metadata-based Workload Profiling

- ▶ Address the SLO and scheduling bootstrapping problem

- ▶ Main idea: use apriori available meta data to profile workloads

  - ▶ Profile generation based on dynamic data

  - ▶ Profile classification based on static metadata

- ▶ Main insight: We can learn a lot about the future workload by looking at its metadata

**CPU usage:**
max 99, std 24, avg ...
**Memory usage:**
max 32, std 5, avg ...

apped by ML

Profile classifier

Mapped by ML →

Mapped by ML →

Profile generator

PolarisProfiler

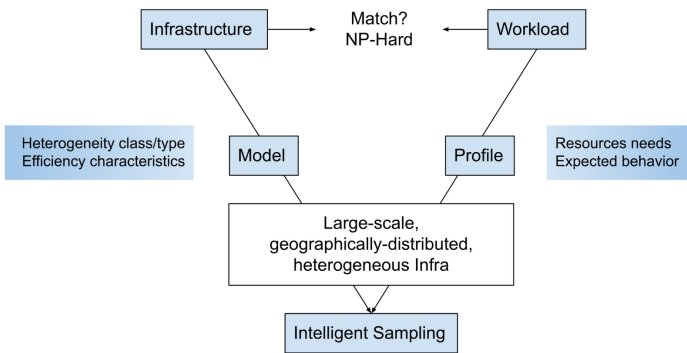e.g., XGBoost                    e.g., HDBSCAN

# Metadata-based Profiler Results

*Predicting job duration*

- 1000 proportionally-sampled jobs not used for training

- Use avg duration in each profile as prediction

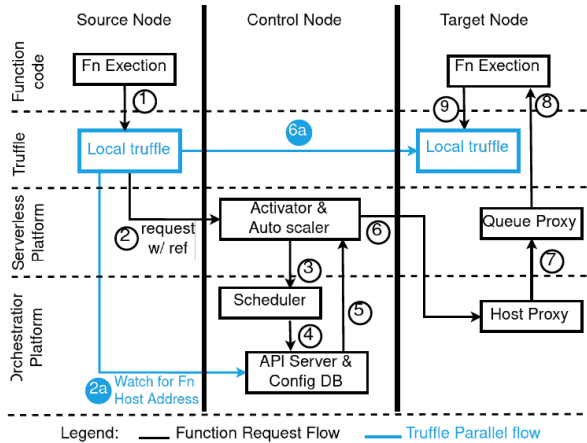- $RMSE_{perc}$ below 5% for more than 80% of the profiled jobs

**Morichetta, A., Casamayor Pujol, V., Nastic, S., Dustdar, S., Vij, D., Xiong, Y., & Zhang, Z. (2023). PolarisProfiler: A Novel Metadata-Based Profiling Approach for Optimizing Resource Management in the Edge-Cloud Continnum. In The 18th IEEE International Symposium on Service-Oriented System Engineering (SOSE 2023).**
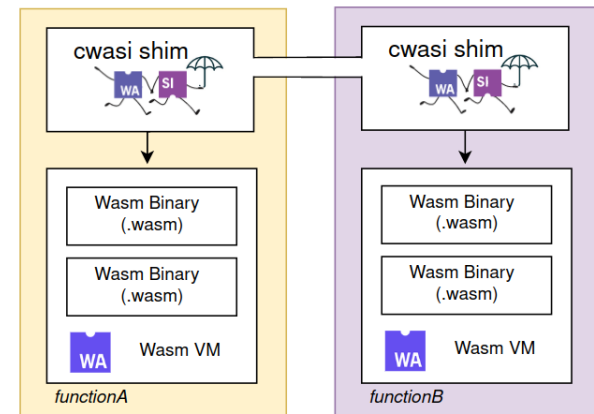
# Work in progress

### Intelligent sampling in a heterogeneous infrastructure



### Using cold starts to improve data transfers between functions



### Optimal function invocation for WASM-based FaaS

# Conclusion

- Main promises of serverless computing
- Main concepts of serverless computing paradigm
  - FaaS – Function as a Service
  - BaaS – Backend as a Service
  - Function virtualization and isolation approaches
- Serverless = FaaS + BaaS
- Project Polaris – towards serverless computing fabric for the edge-cloud continuum
  - SLO managements
  - Scheduling
  - Profiling

# Thank you for your attention!

▶ Find me at: **www.nastic.at**

▶ Drop me an email: **snastic@dsg.tuwien.ac.at**