

Sky Computing for Serverless: Infrastructure Assessment to Support Performance Enhancement

Robert Cordingly, Xinghan Chen, Ling-Hong Hung, Wes Lloyd rcording@uw.edu

School of Engineering and Technology University of Washington Tacoma

IEEE/ACM 18th International Conference on Utility and Cloud Computing (UCC 25) December 1-4, 2025, Nantes France

1

Background and Motivation







What is Serverless?



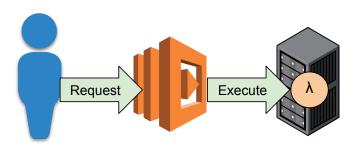


Serverless function-as-a-service (FaaS) platforms offer many desirable features:

- Rapid elastic scaling
- Scale to zero
- No infrastructure management
- Fine grained billing
- Fault tolerance
- No up front cost to deploy an application

3

How FaaS Platforms Work



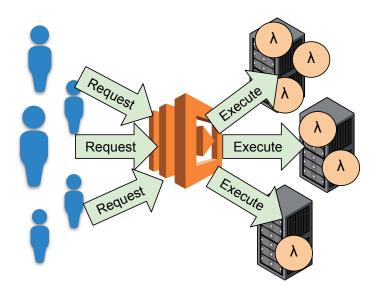
- Clients make requests to the FaaS platform which manage the infrastructure automatically
- FaaS platforms create and execute users code inside environments known as function instances
- Function instances can be hosted on a variety of different kinds of hardware

Client

FaaS Platform

Function Instance

Automatic Scaling



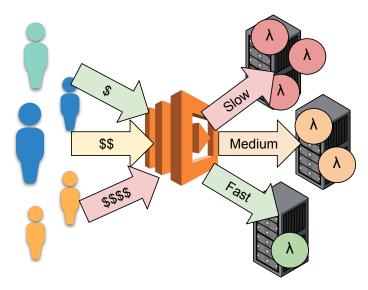
Clients

FaaS Function Function Instances

- As more requests are made, the FaaS platform scales the number of function instances
- Function instances for one function can be spread across many host servers

5

Inconsistent Performance



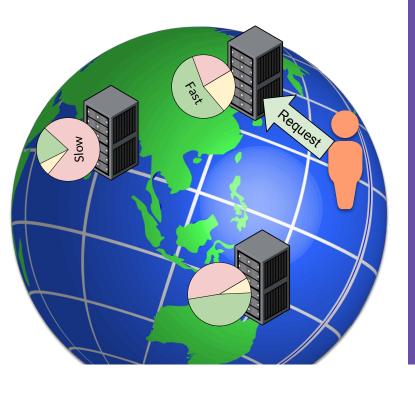
Clients

FaaS Function Function Instances

- Heterogeneous hardware leads to some hosts being faster than others
- Since FaaS platforms bill based off runtime, slower hosts perform worse and are more expensive

6

Request Routing



- If we were able to observe all of the available infrastructure of serverless regions, we could route requests to regions with more favorable infrastructure
- Our previous work built a pserverless sky proxy system to route requests around the world (IC2E 2023)

What is Sky Computing?

Azure

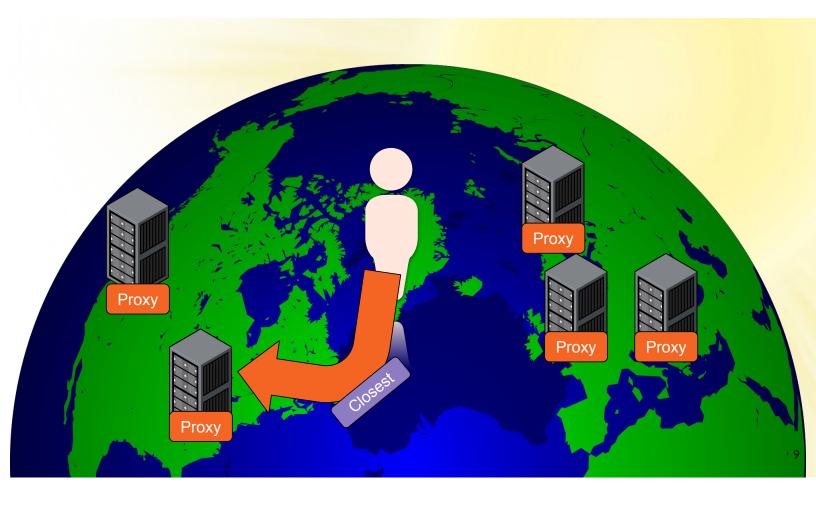
Coogle Cloud

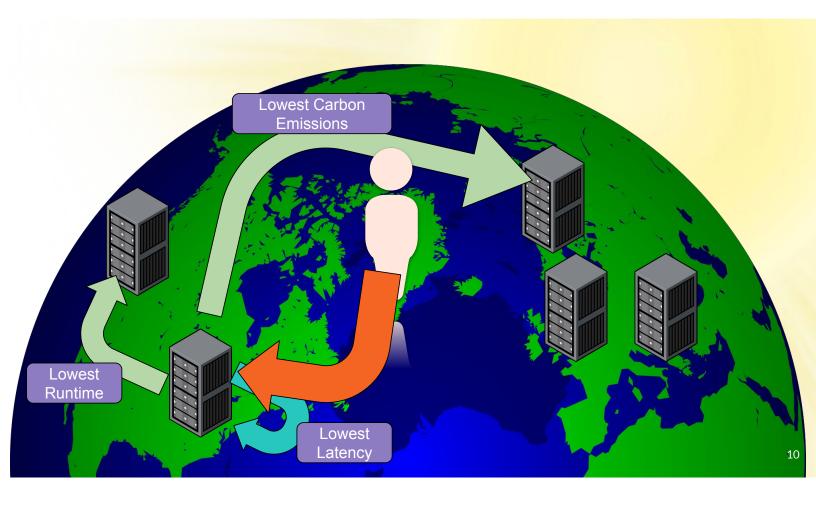
Aws

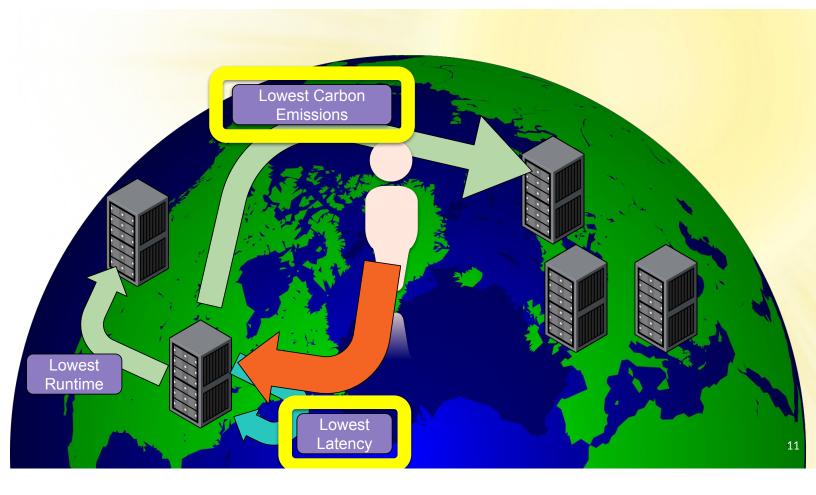
- The Sky sits above the clouds
- Consists of compatibility layers allowing aggregation of resources between cloud regions, availability zones, or cloud providers

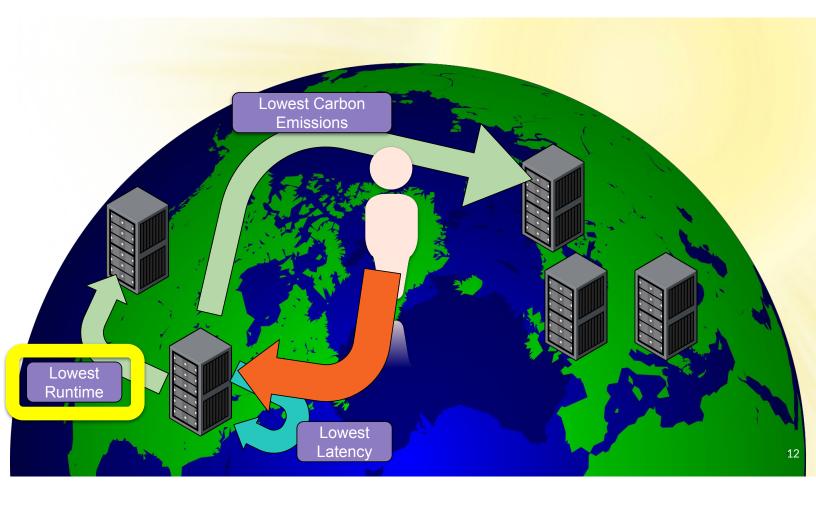
Goals for Serverless Sky Computing:

- Allow applications to take advantage of resources of multiple regions and cloud providers
- Improve performance and reduce costs of FaaS applications









Research Questions

13

Research Questions



• RQ-1 (Infrastructure Variation): What CPU variations can be observed on FaaS platforms across different zones and regions? How does it change over time?



 RQ-2 (Infrastructure Characterization): How many samples are required to accurately infer the hardware pool? How can we balance the cost versus accuracy?



• RQ-3 (**Performance Optimization**): To what extent are runtime and cost improvements possible by targeting specific instructure?

Serverless Infrastructure Observation



Supporting Tools - SAAF

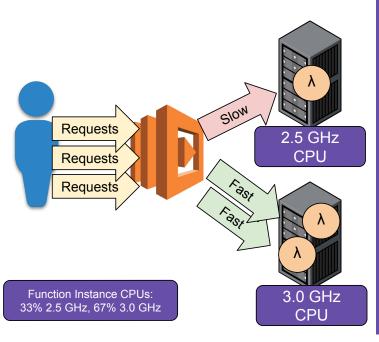
We utilize the **Serverless Application Analytics Framework** to collect various metrics about the infrastructure and platform serverless functions are hosted with.

SAAF collected CPU Timing metrics, latency information, hardware specifications, runtime metrics, and more.

We utilize CPU metrics from SAAF to characterize function instances and observe details about the host infrastructure.

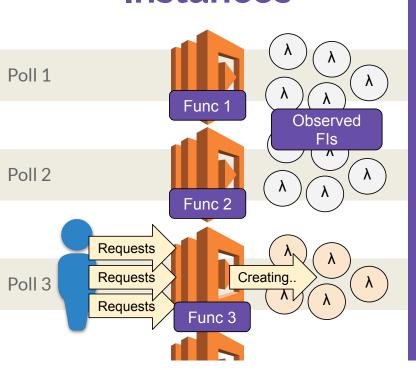
17

Observing Infrastructure



- We deployed "Hello World" sleep functions with SAAF that collect CPU information
- We then make parallel requests to view available infrastructure for a function
- With SAAF's data we build a characterization of the pool of available infrastructure
- But this method is limited...

Creating Function Instances



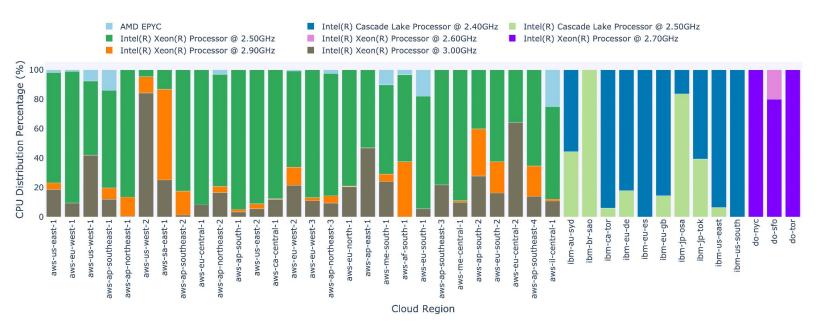
- To observe as much infrastructure as possible we need to create many function instances
- FaaS platforms limit the number of concurrent executions (1000 on AWS Lambda)
- By utilizing many function deployments, we can create more function instances than the concurrency limit
- With this we can build accurate characterizations of available CPUs in availability zones (AZs)

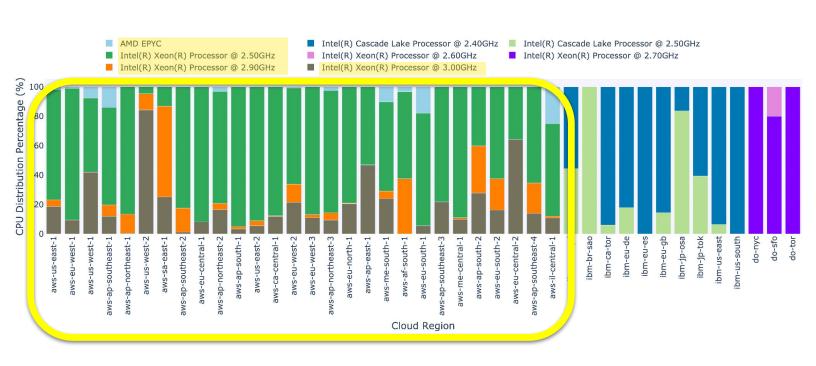
Methodology and Results

19

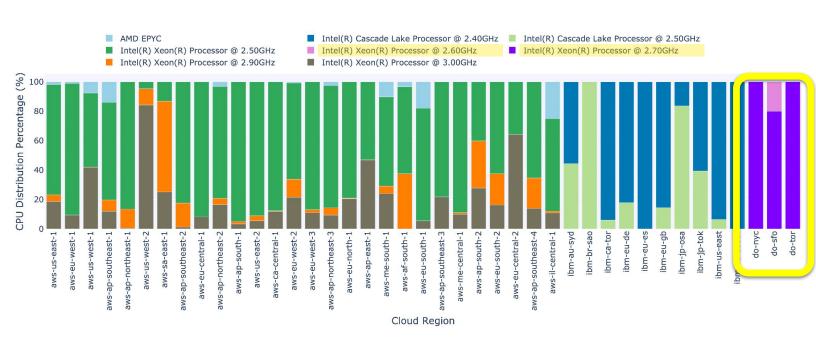
RQ-1 (Part 1): What CPU variations can be observed on FaaS platforms across different zones and regions?



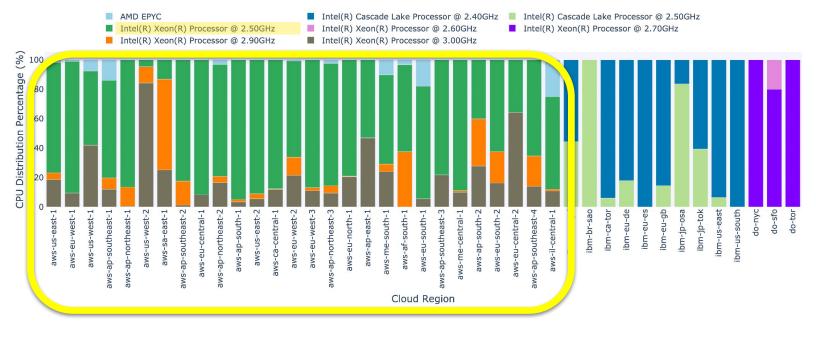


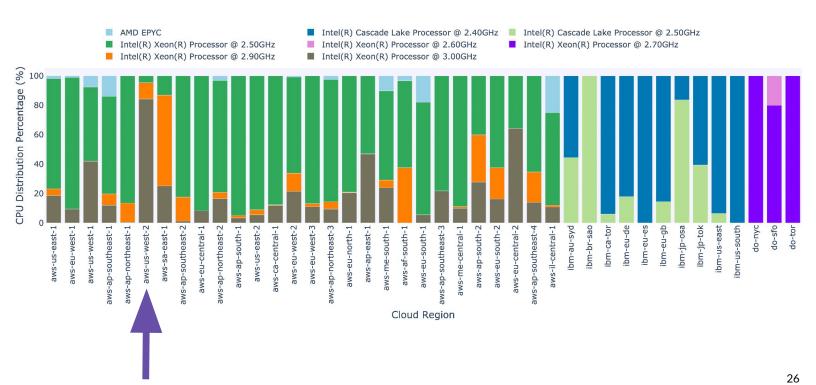


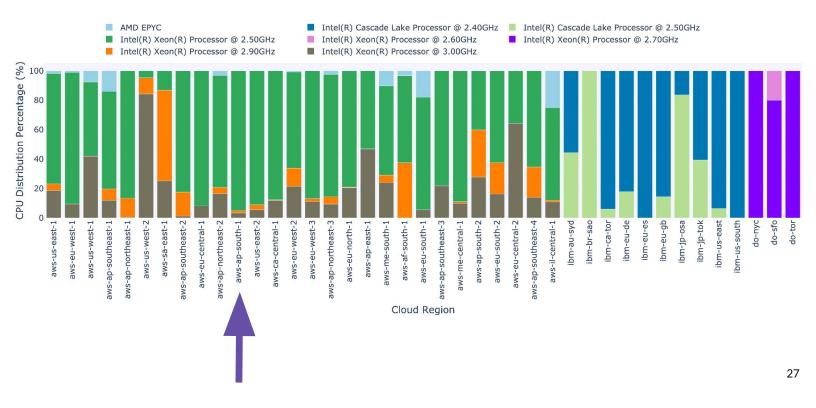


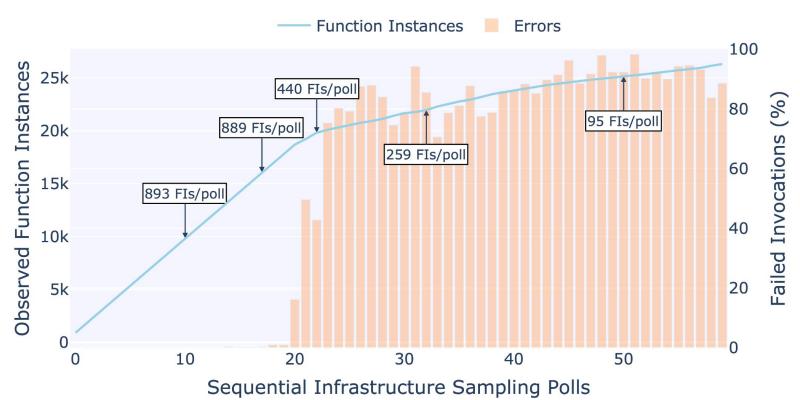


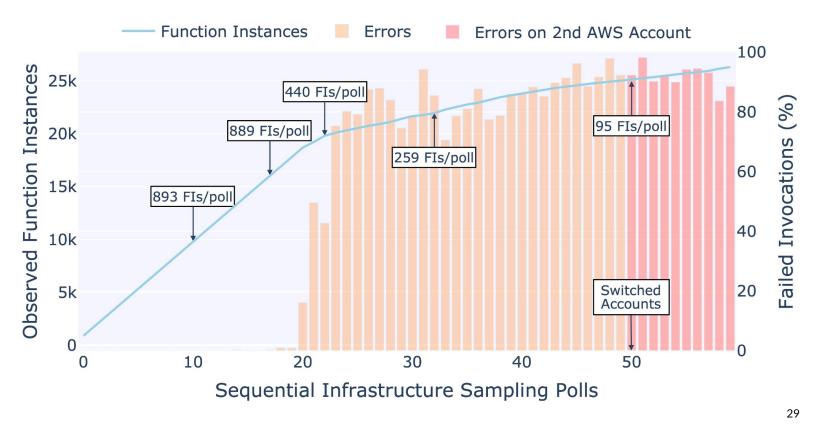


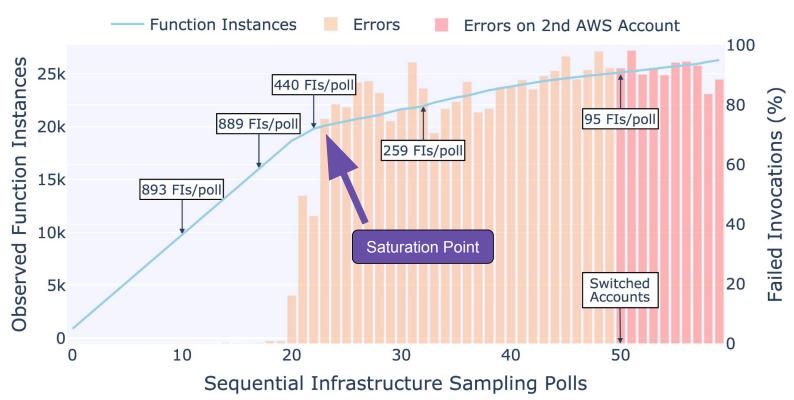






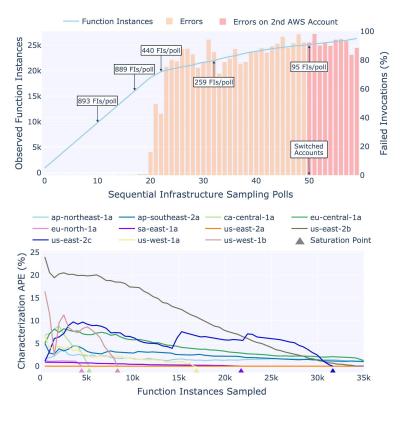




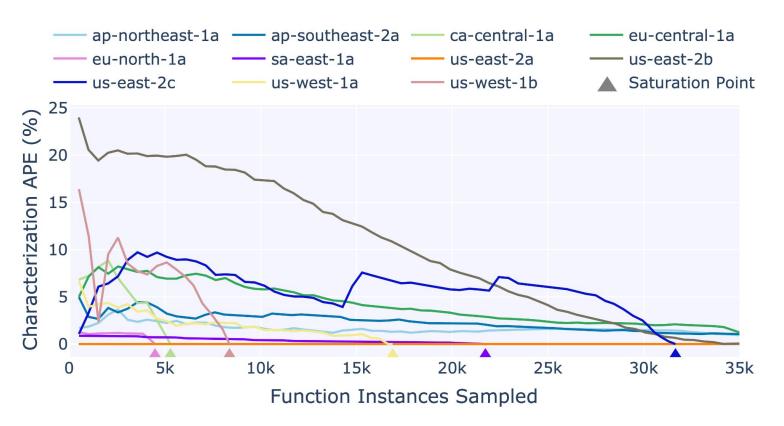


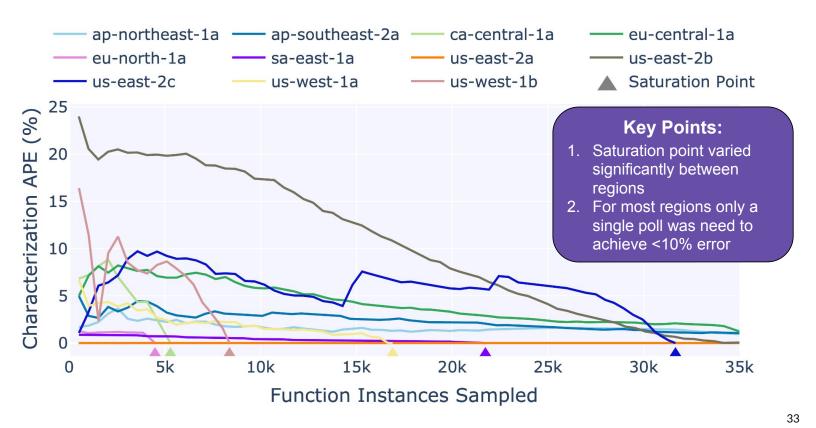
21

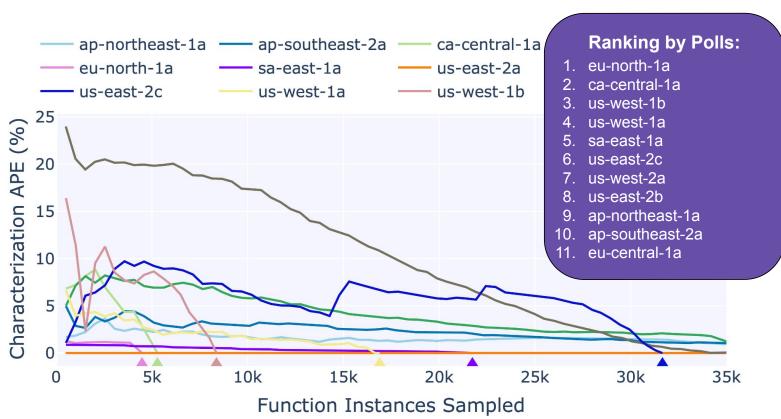
Sanity Checks



- To verify this was not a simple rate limit we:
 - o used separate AWS accounts with:
 - different email addresses
 - different payment methods
 - different locations/network
 - different deployed functions
 - different FaaS configurations (checksum/RAM)
 - different invocation methods:
 - Function URLs
 - API Gateway
 - AWS CLI
- The only similarity was functions shared an availability zone
- Even with all of these tests we still observed these errors

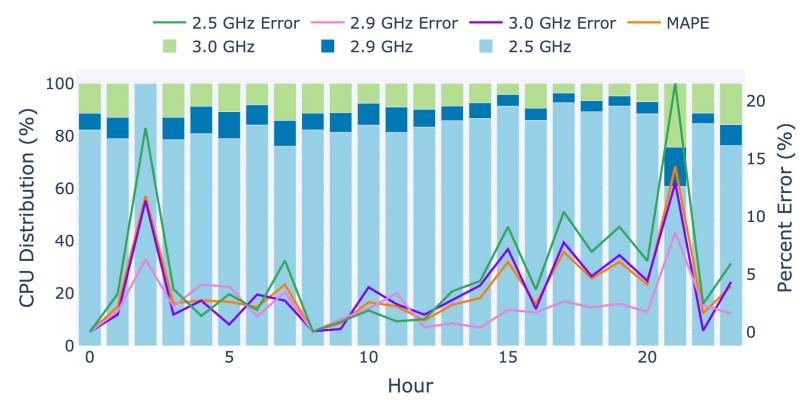




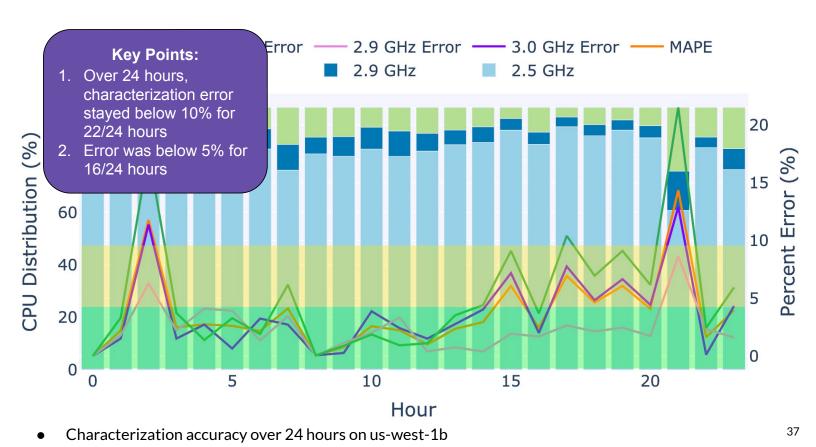


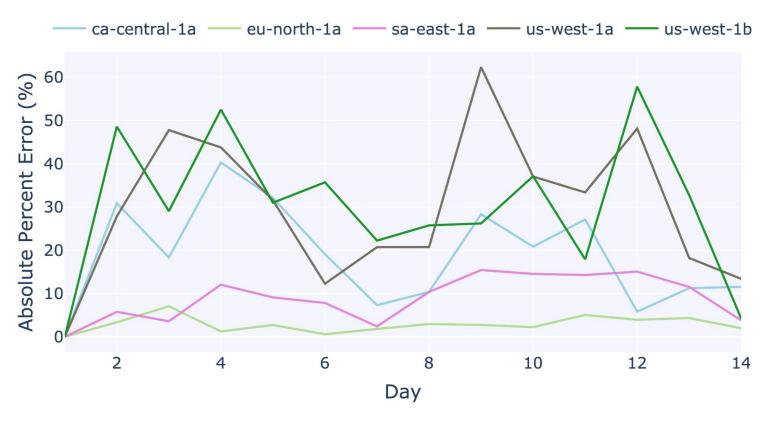
RQ-1 (Part 2): How does characterization accuracy change over time?

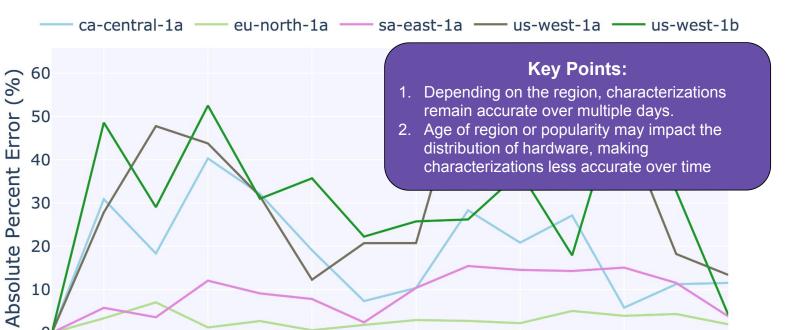
35



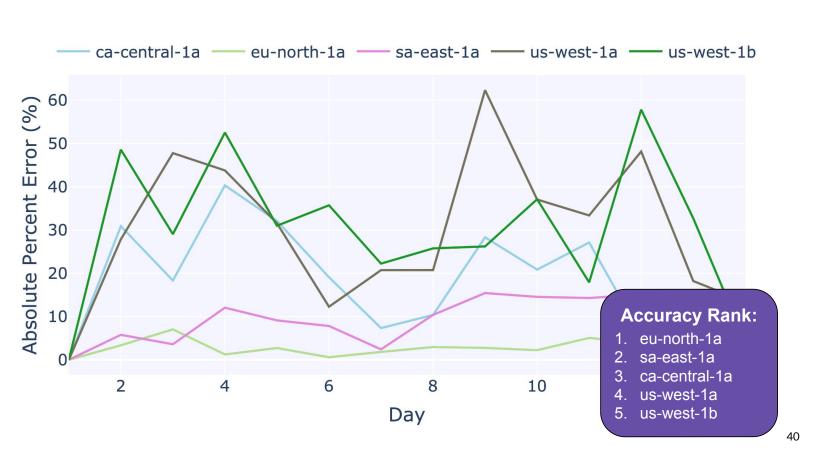
Characterization accuracy over 24 hours on us-west-1b



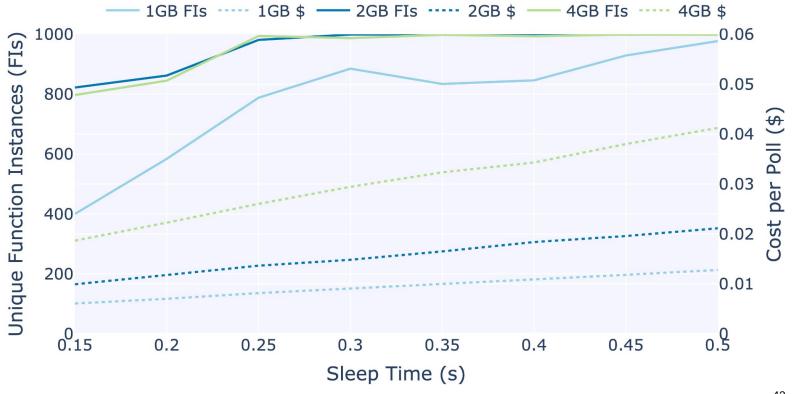


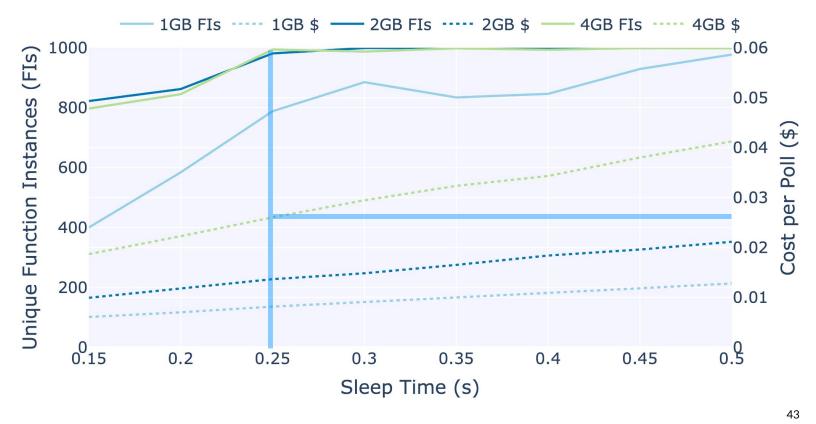


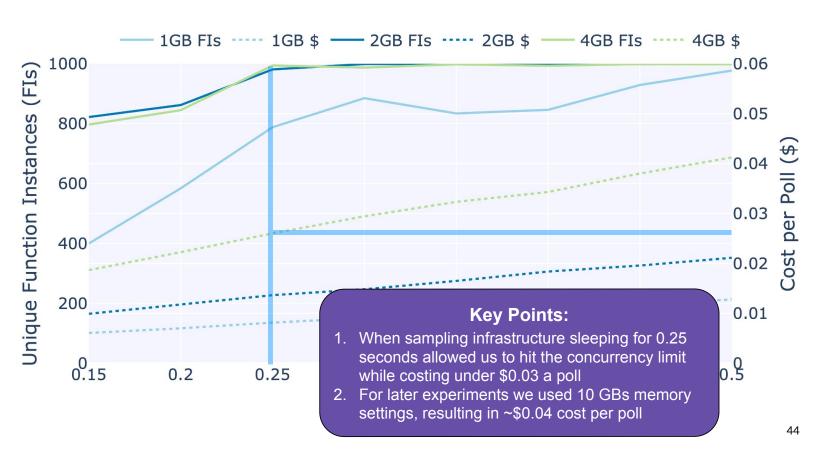
Day

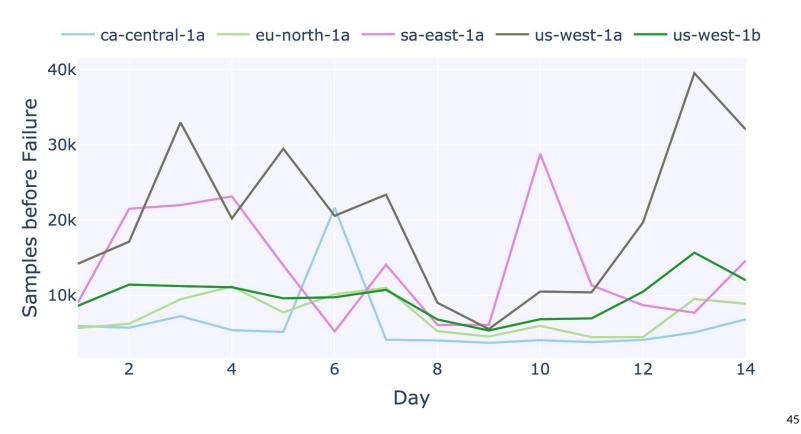


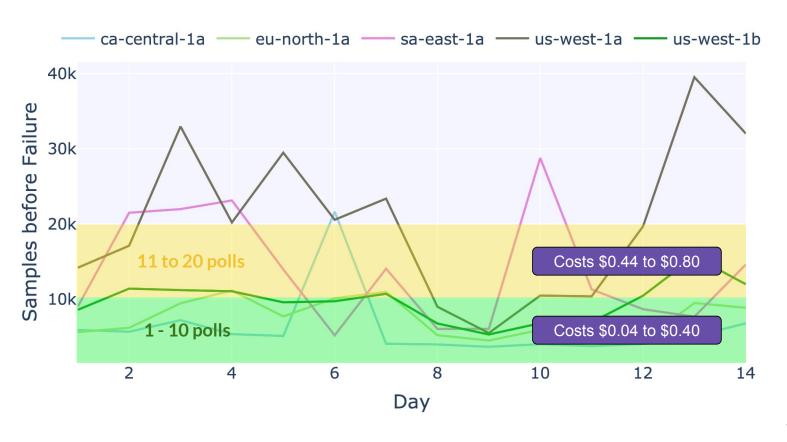
RQ-2: How many samples are required to accurately infer the hardware pool? How can we balance the cost versus accuracy?

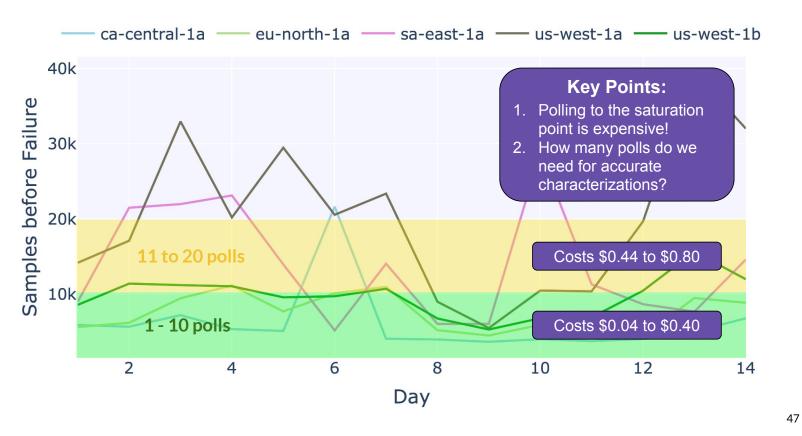


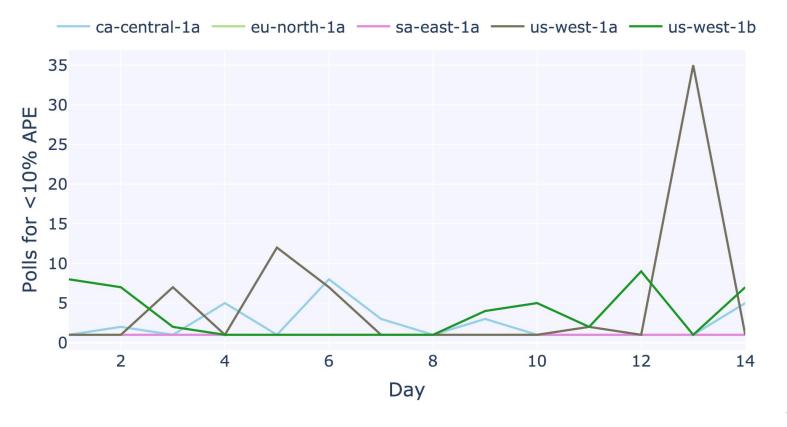


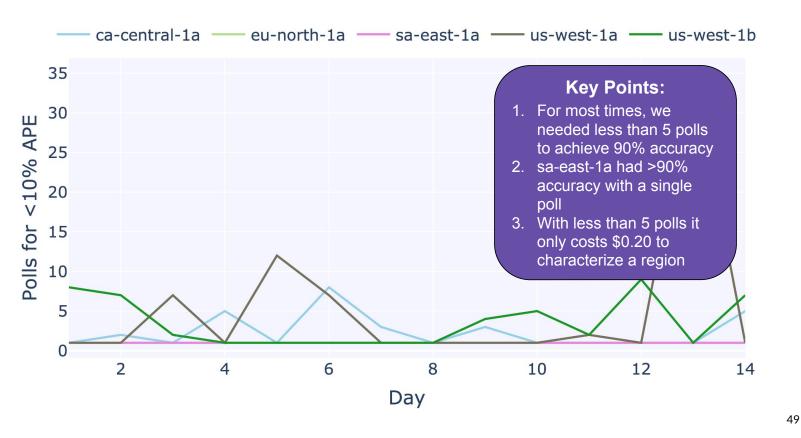




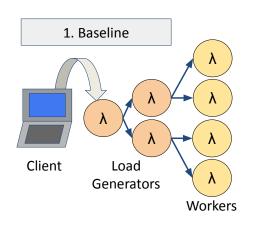


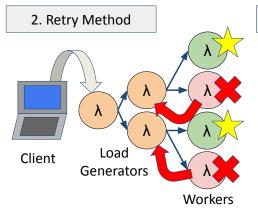


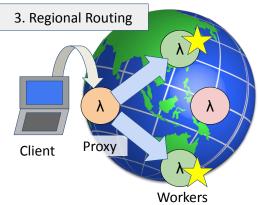




RQ-3: To what extent are runtime and cost improvements possible by targeting specific instructure?





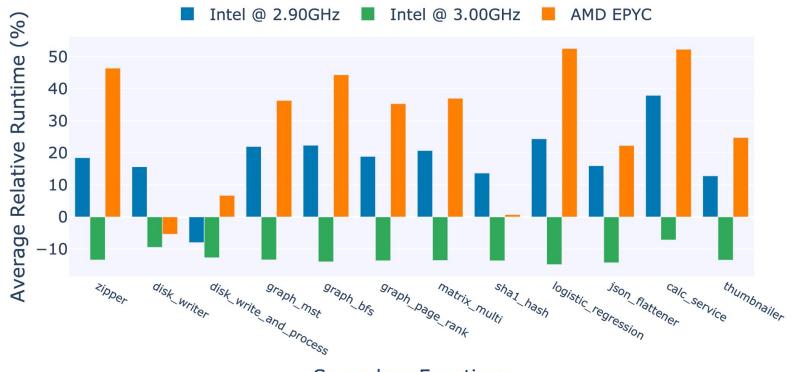


Retry Logic

- 1. **Retry Slow:** Invocations retry on the two slowest CPUs
- Focus Fastest: Invocations retry on everything but the fastest CPU

Hybrid Approach

Combine both Retries and Regional Routing

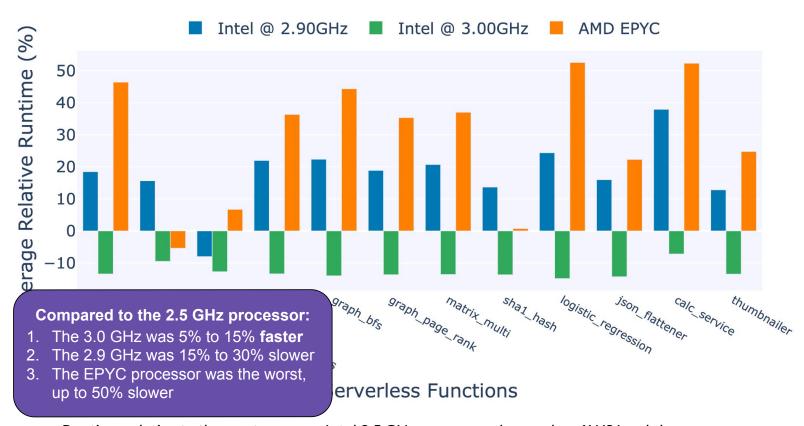


Serverless Functions

Runtime relative to the most common Intel 2.5 GHz processor observed on AWS Lambda



Runtime relative to the most common Intel 2.5 GHz processor observed on AWS Lambda

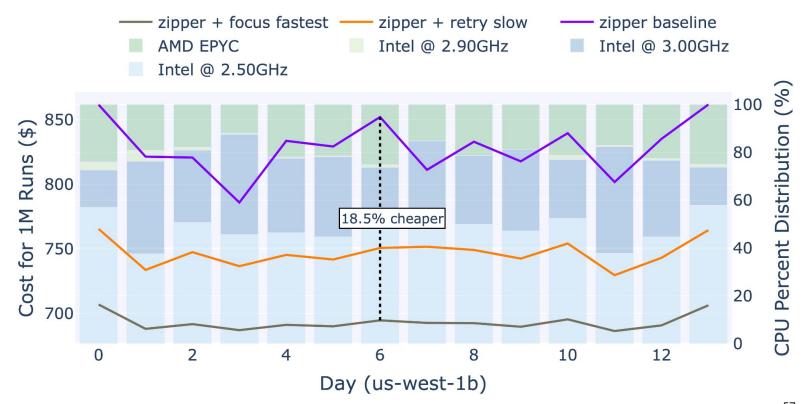


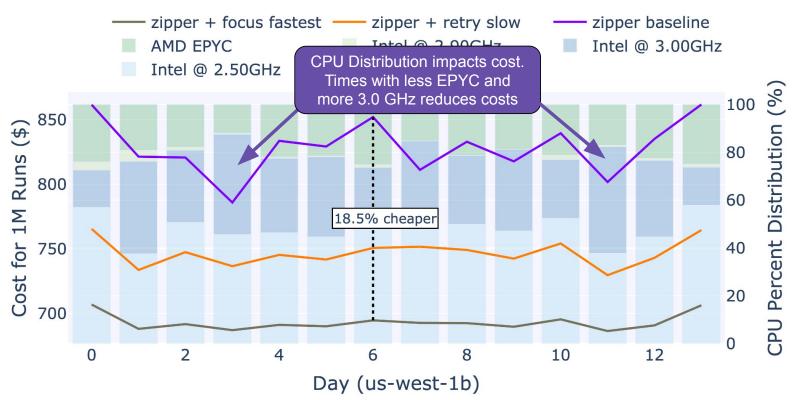
Runtime relative to the most common Intel 2.5 GHz processor observed on AWS Lambda

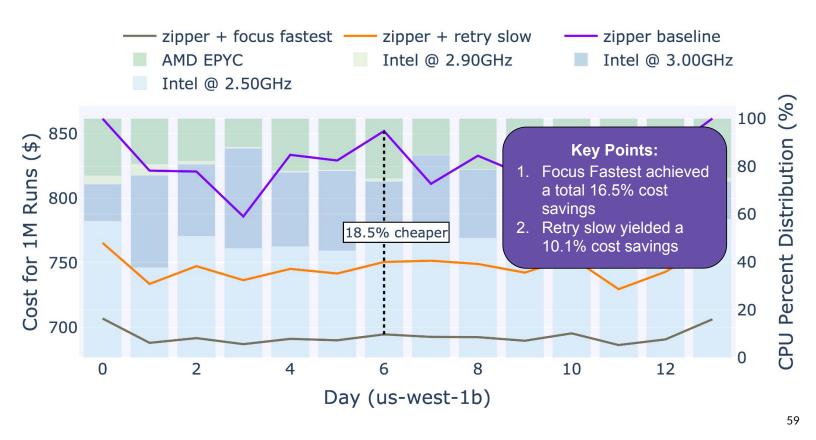
Runtime relative to the most common Intel 2.5 GHz processor observed on AWS Lambda

Retry Methods



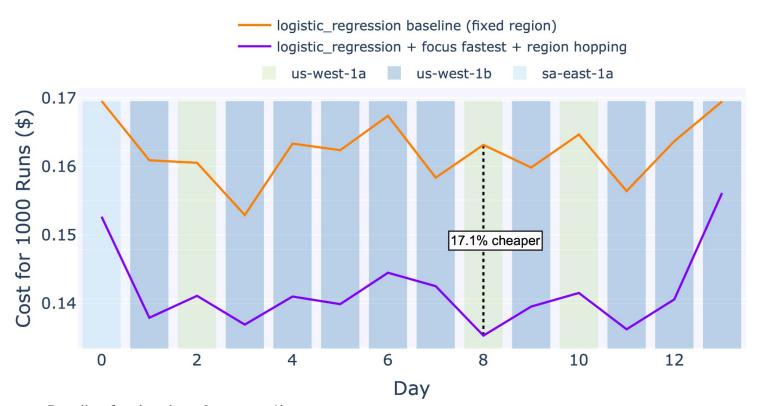


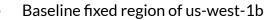


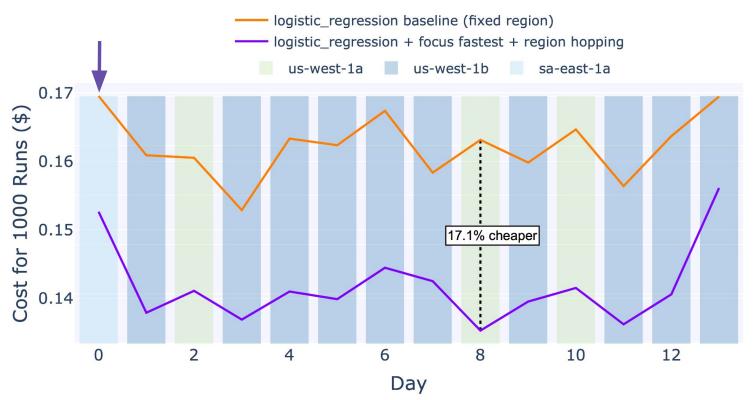


Region Hopping

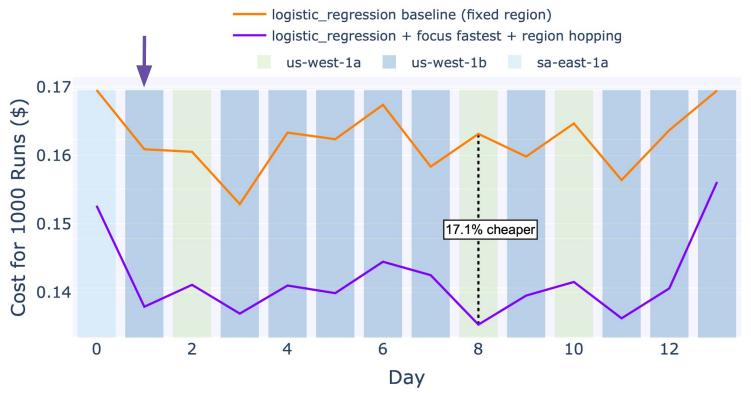




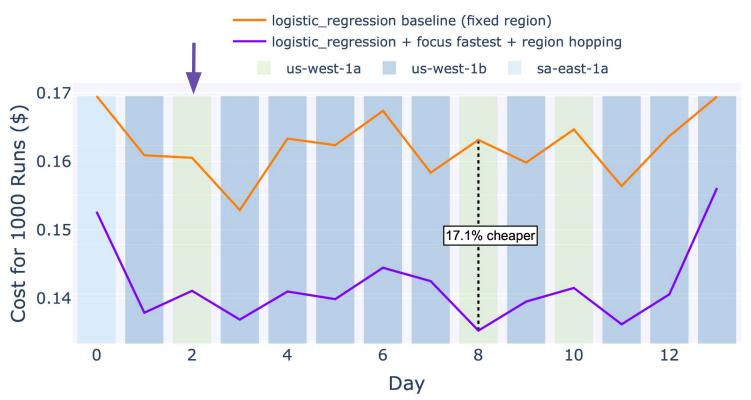




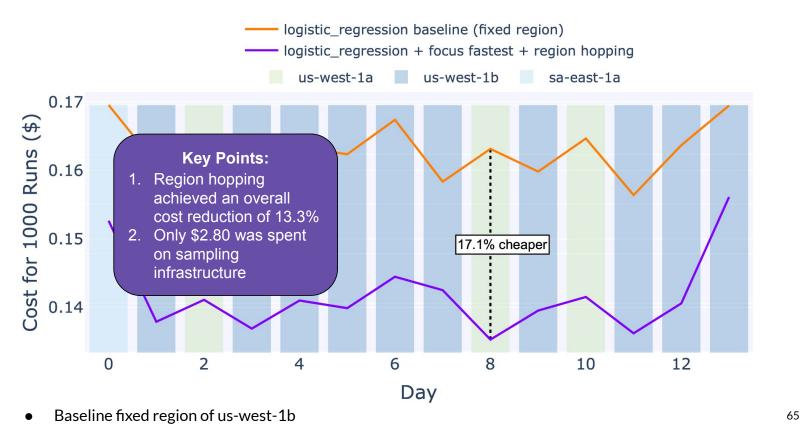
Baseline fixed region of us-west-1b







Baseline fixed region of us-west-1b



Conclusions







Key Results

- **RQ-1:** We were able to observe diverse CPU characterizations across 41 serverless regions
- RQ-1: Our sampling technique appears to be able to saturate entire availability zones creating accurate hardware characterizations
- **RQ-2:** Our sampling technique is able to characterize an availability zone for only \$0.04 with 95% accuracy
- **RQ-3:** The **Retry Slow** method resulted in a consistent 10.1% reduction in runtime and cost
- **RQ-3: Focus Fastest** had up to 18.5% lower costs, averaging 16.5% across all functions
- **RQ-3:** Compared to a baseline of running in us-west-1b, the **Region Hopping** hybrid approach averaged a cost savings of 10.0% up to 18.2%

67

Future Work

Future Work and Limitations

- The retry methods sacrifice responsiveness for cost savings
- Building CPU characterizations for many regions is expensive
- In the future, a complete Serverless Sky Computing platform could build the CPU characterizations as workloads run, eliminating the need for dedicated polling
- Instead of relying on individual characterizations, they could be updated as new functions execute

Thank You!

Any Questions?

This research has been supported by the of Washington Carwein-Andrews Distinguished Fellow Fund and AWS Educate Cloud Credits.