



Tutorial: Enhancing Observability of Serverless Computing with the Serverless Application Analytics Framework (SAAF)

Robert Cordingly, Navid Heydari, Hanfei Yu,
Varik Hoang, Zohreh Sadeghi, Wes Lloyd

School of Engineering and Technology
University of Washington Tacoma

12th ACM/SPEC International Conference on Performance Engineering (ICPE 2021)

SAAF Outline

- Introduction to Serverless Computing
 - Motivation
 - Delivery models and platforms
 - Advantages and challenges
- Serverless Application Analytics Framework (SAAF)
 - Design of SAAF, Supported Languages, Metrics
 - Tools: FaaS Runner, Publish Script
- Analysis Examples with SAAF
 - Programming language comparison, performance modeling
 - Scalability testing, Resource utilization profiling
 - Tracking infrastructure reuse
- Conclusions
- SAAF Demo



How should my app withstand a server failing?

How can I tell if a server has been compromised?

How can I increase utilization of my servers?

Which OS should my servers run?

How much remaining capacity do my servers have?

When should I decide to scale up my servers?

What size servers are right for my budget?

How should I implement dynamic configuration changes on my servers?

How will I keep my server OS patched?

How can I control access from my servers?

Which packages should be baked into my server images?

Servers

(AAHHHHHHHHH!!!)

How will new code be deployed to my servers?

How will the application handle server hardware failure?

How many users create too much load for my servers?

What size server is right for my performance?

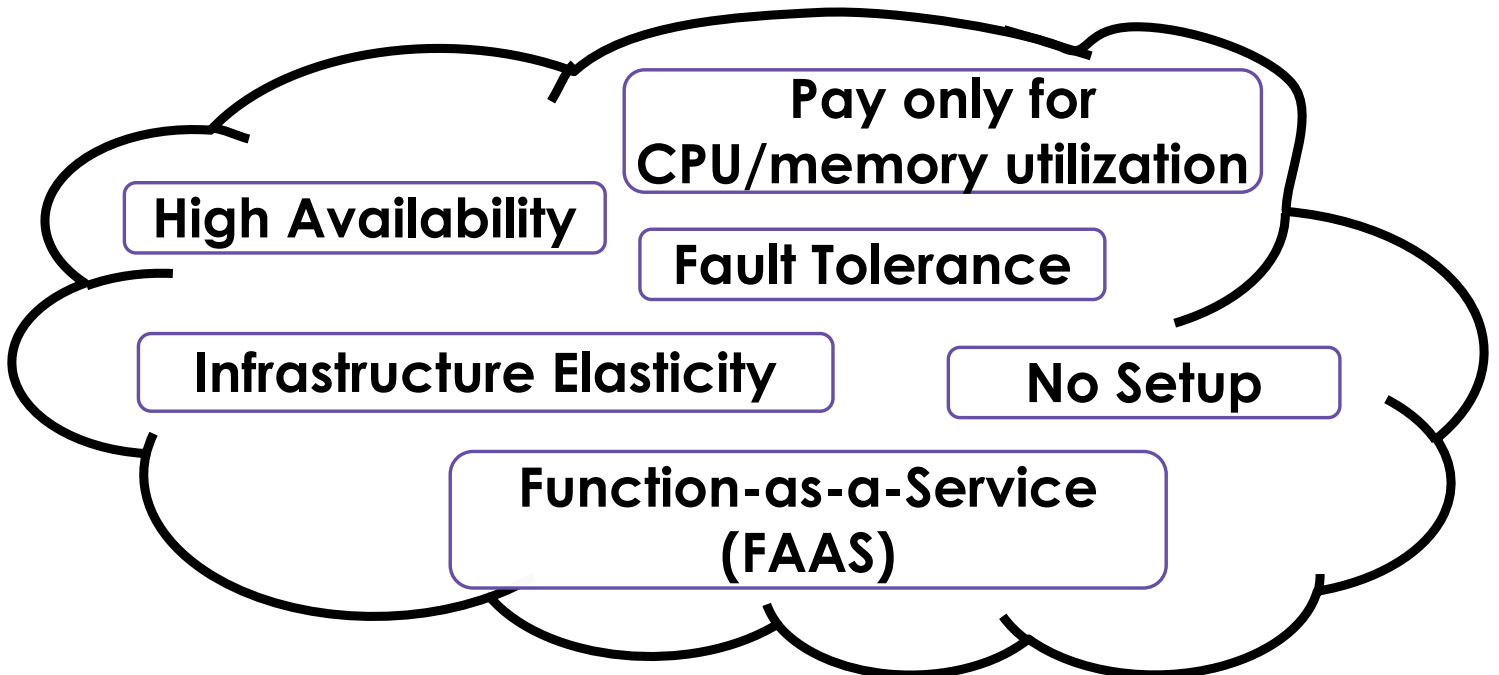
Which users should have access to my servers?

Should I tune OS settings to optimize my application?

When should I decide to scale out my servers?

How many servers should I budget for?

Serverless Computing



Serverless Computing

Why Serverless Computing?

Many features of distributed systems, that are challenging to deliver, are provided automatically

...they are built into the platform

5

SAAF Outline

- Introduction to Serverless Computing
 - Motivation
 - Delivery models and platforms
 - Advantages and challenges
- Serverless Application Analytics Framework (SAAF)
 - Design of SAAF, Supported Languages, Metrics
 - Tools: FaaS Runner, Publish Script
- Analysis Examples with SAAF
 - Programming language comparison, performance modeling
 - Scalability testing, Resource utilization profiling
 - Tracking infrastructure reuse
- Conclusions
- SAAF Demo



6

Serverless Computing Delivery Models

Function-as-a-Service (FaaS)

Container-as-a-Service (CaaS)

Database-as-a-Service (DBaaS)

7

Commercial FaaS Platforms

AWS Lambda

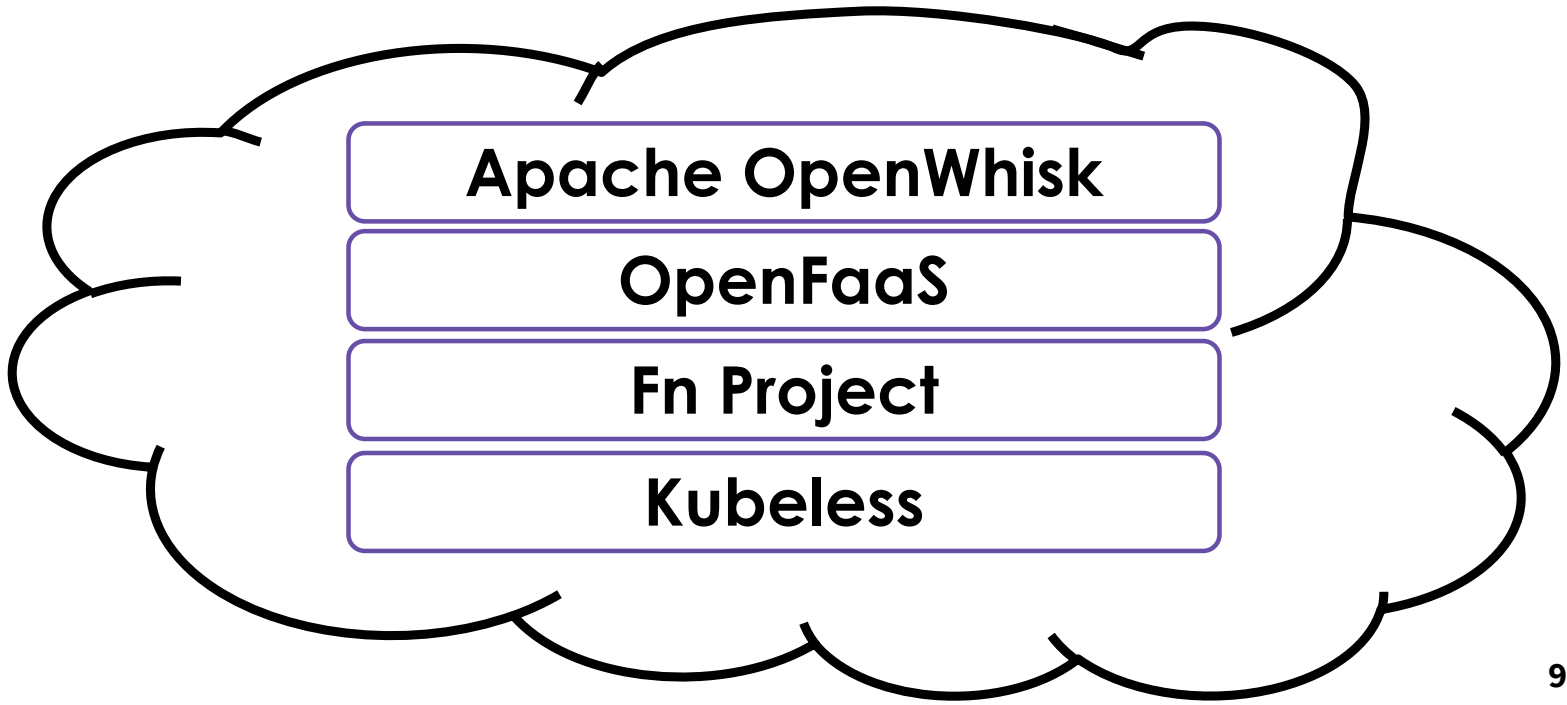
Azure Functions

IBM Cloud Functions

Google Cloud Functions

8

Open Source FaaS Platforms



9

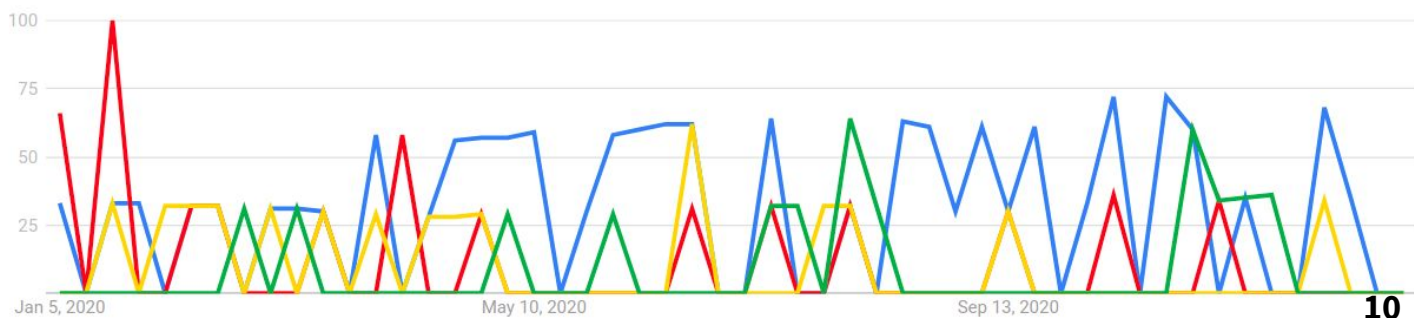
Google Search Trends: Open Source FaaS Platform 2020

OpenFaaS: 31

Apache OpenWhisk: 10

Kubeless: 10

Fn Project: 9



10

FaaS Platform Example: AWS Lambda

Bring your own code

- Languages: Java, Python, Node.js, Go, C#, Ruby, PowerShell, Bash
- Bring you own libraries

Flexible use

- Synchronous or asynchronous
- Integration w/ other AWS services

Simple Resource Model

- Function memory 128 MB to 10 GB
- CPU timeshare and network bandwidth scaled proportional to memory

Flexible Authorization

- Grant access to resources and VPCs
- Fine-grained control for invoking functions

11

SAAF Outline

- Introduction to Serverless Computing
 - Motivation
 - Delivery models and platforms
 - Advantages and challenges
- Serverless Application Analytics Framework (SAAF)
 - Design of SAAF, Supported Languages, Metrics
 - Tools: FaaS Runner, Publish Script
- Analysis Examples with SAAF
 - Programming language comparison, performance modeling
 - Scalability testing, Resource utilization profiling
 - Tracking infrastructure reuse
- Conclusions
- SAAF Demo



12

Function-as-a-Service

Advantages

- No management of servers
- Pay only for actual compute time
- High availability (24/7)
- Scalability (elastic resources)
- Fault tolerance
- Rapid deployment/updates
- Supports vendor workload consolidation

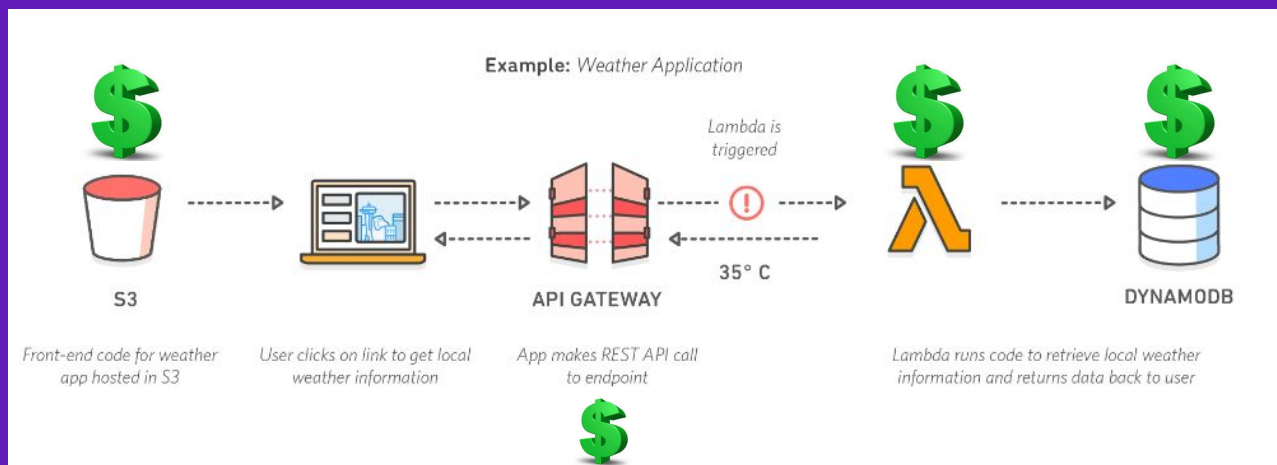
Challenges

- Limited observability of servers
- Multi-dimensional pricing policies
- Vendor lock-In
- Heterogeneous infrastructure
- Performance variation
- Memory reservation size
- Infrastructure freeze/thaw
- Function composition
- Pricing obfuscation

13

Vendor Architectural Lock-In

Cloud native (FaaS) software architecture requires external services/components



Increased dependencies → increased hosting costs

14

Pricing Obfuscation

VM pricing: hourly pricing policy, billed to the nearest second is intuitive to understand...

FaaS pricing:

AWS Lambda Pricing

FREE TIER: first 1,000,000 function calls/month ☐ FREE
first 400,000 GB-sec/month ☐ FREE

Afterwards: \$0.0000002 per request
\$0.000000208 to rent 128MB / 100-ms -or-
\$0.00001667 to rent 1GB / 1-sec

15

IaaS Cloud Pricing Policies

Virtual machines as-a-service at ¢ per hour

No premium to scale:

$$= \frac{1000 \text{ computers}}{1 \text{ computer}} @ \frac{1 \text{ hour}}{1000 \text{ hours}}$$

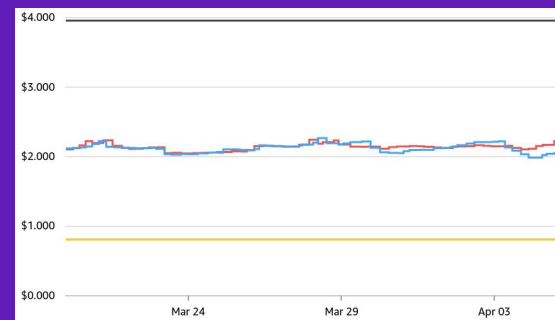
Illusion of infinite scalability to cloud user

As many computers as you can afford

Pricing policies are becoming increasingly granular

- By the minute, second

Preemptive/Spot reduced price instances ☐



16

Price Comparison: IaaS vs. FaaS

Assume 1 month = 30.41667 days (365d / 12)

Workload

Continuous 1-sec service calls: 100% of 2 CPU-cores
100% of 4GB of memory
Workload: 2 continuous client threads
Duration: 1 month (30.41667 days)

ON AWS EC2: Amazon EC2 c5.large 2-vCPU VM x 4GB
c5.large: 8.5¢/hour, 24 hrs/day x 30.41667 days
Hosting cost: \$62.05/month

17

Price Comparison: IaaS vs. FaaS - II

Assume 1 month = 30.41667 days (365d / 12)

Workload

Continuous 1-sec service calls: 100% of 2 CPU-cores
100% of 4GB of memory
Workload: 2 continuous client threads
Duration: 1 month (30.41667 days)

ON AWS Lambda: 2,628,000 function calls, 1-sec @ 4GB
function calls: 2.628 million x 20¢/million
runtime: 2,628,000 sec x 4 GB
Hosting cost: ???

18

Price Comparison: IaaS vs FaaS - III

Workload:	(4GB)	10,512,000 GB-sec	
FREE:	-	<u>400,000 GB-sec</u>	
Runtime/Memory:		10,112,000 GB-sec	
Charge:	x .00001667 GB-sec		<u>\$168.57</u>
Invocations:		2,628,000 calls	
FREE:	-	<u>1,000,000 calls</u>	
Charge:		1,628,000 calls	<u>\$.33</u>
Total:			<u>\$168.90 (272.2%)</u>

BREAK-EVEN POINT = \$62.05 (VM) - \$0.33 (calls) = \$61.72
 $\$61.72 / .00001667 \text{ GB-sec} = \sim 3,702,459 \text{ GB-sec-month} / 4\text{GB/call} = \sim 925,614 \text{ sec}$
 $\sim 10.71 \text{ days}$

Point at which using FaaS costs the same as IaaS

19

Price Comparison: IaaS vs FaaS - III

Workload:	(4GB)	10,512,000 GB-sec	
FREE:	-	<u>400,000 GB-sec</u>	
Runtime/Memory:		10,112,000 GB-sec	
Charge:	x .00001667 GB-sec		<u>\$168.57</u>
Invocations:		2,628,000 calls	
FREE:	-	<u>1,000,000 calls</u>	
Charge:		1,628,000 calls	<u>\$.33</u>
Total:			<u>\$168.90 (272.2%)</u>

BREAK-EVEN POINT = \$62.05 (VM) - \$0.33 (calls) = \$61.72
 $\$61.72 / .00001667 \text{ GB-sec} = \sim 3,702,459 \text{ GB-sec-month} / 4\text{GB/call} = \sim 925,614 \text{ sec}$
 $\sim 10.71 \text{ days}$

Worst-case scenario = ~2.72x !

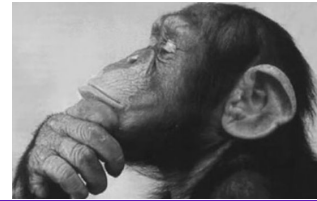
AWS EC2: \$62.05

AWS Lambda: \$168.90

Point at which using FaaS costs the same as IaaS

20

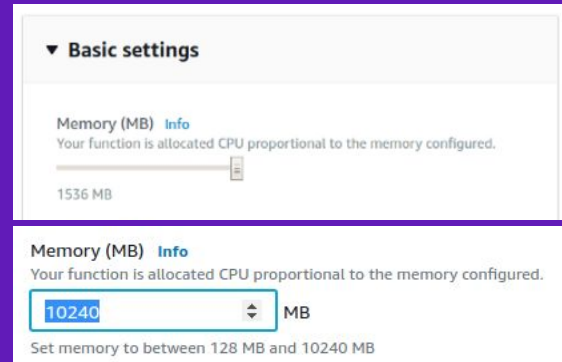
Memory Reservation



Lambda memory reserved for functions

UI provides textbox to set function's memory (*previously a slider*)

Resource capacity (CPU, disk, network) scaled relative to memory



“every doubling of memory, doubles CPU...”

But how much memory do functions require?

21

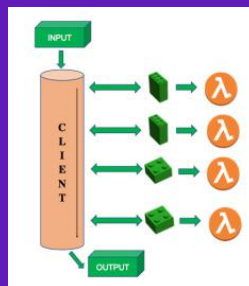
Service Composition

How should applications be decomposed into serverless functions for deployment?

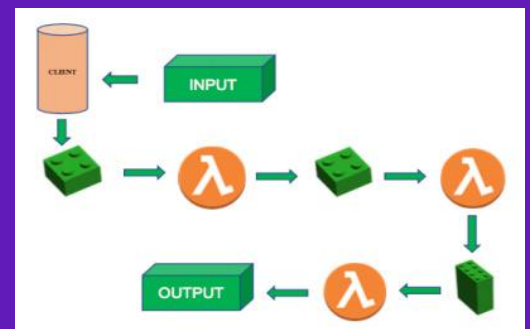
Monolithic Deployment



Client flow control, 4 functions



Server-side flow control, 3 functions



Recommended practice:

Decompose into many microservices

Platform limits: code + libraries ~250MB (uncompressed)

How does composition impact the number of function invocations, and memory utilization?

22

Infrastructure Freeze/Thaw Cycle

Unused infrastructure is deprecated

• But after how long?

FaaS Infrastructure known as “function instances”

Implemented using VMs/microVMs, containers, or software-based isolates

STATES:

Physical Host-COLD

• Function code not yet transferred to any server

Container-COLD

• Function code transferred to server, but infrastructure not created

Container-WARM

• Active container running

23

SAAF Outline

- Introduction to Serverless Computing
 - Motivation
 - Delivery models and platforms
 - Advantages and challenges
- **Serverless Application Analytics Framework (SAAF)**
 - Design of SAAF, Supported Languages, Metrics
 - Tools: FaaS Runner, Publish Script
- Analysis Examples with SAAF
 - Programming language comparison, performance modeling
 - Scalability testing, Resource utilization profiling
 - Tracking infrastructure reuse
- Conclusions
- SAAF Demo



24

SAAF for FAAS

To address these challenges better FaaS profiling & measurement tools are required:

- Limited observability of servers
- Multi-dimensional pricing policies
- Vendor lock-In
- Heterogeneous infrastructure
- Performance variation
- Memory reservation size
- Infrastructure freeze/thaw
- Function composition
- Pricing obfuscation



SAAF: The Serverless Application Analytics Framework

Example Output JSON:

The attributes collect can be customized by changing which functions more detailed descriptions of each variable and the functions that c see the framework documentation for each language.

```
{  "version": 0.2,  "lang": "python",  "cpuType": "Intel(R) Xeon(R) Processor @ 2.50G",  "cpuModel": 63,  "vmuptime": 1551727835,  "uuid": "d241c618-78d8-48e2-9736-997dc1a931d4",  "vmID": "tiUCnA",  "platform": "AWS Lambda",  "newcontainer": 1,  "cpuUsrDelta": "904",  "cpuNiceDelta": "0",  "cpuKrnDelta": "585",  "cpuIdleDelta": "82428",  "cpuIowaitDelta": "226",  "cpuIrqDelta": "0",  "cpuSoftIrqDelta": "7",  "vmcpusteaDelta": "1594",  "frameworkRuntime": 35.72,  "message": "Hello Fred Smith!",  "runtime": 38.94}
```

Attributes Collected by Each Function

The amount of data collected is determined by which functions are called. If some attributes are not needed, then some functions may not need to be called. If you would like to collect every attribute, the inspectAll() method will run all methods.

Core Attributes

Field	Description
version	The version of the SAAF Framework.
lang	The language of the function.
runtime	The server-side runtime from when the function is initialized until Inspector.finish() is called.
startTime	The Unix Epoch that the Inspector was initialized in ms.

inspectContainer()

Field	Description
uuid	A unique identifier assigned to a container if one does not already exist.
newcontainer	Whether a container is new (no assigned uuid) or if it has been used before.
vmuptime	Time when the host booted in seconds since January 1, 1970 (Unix epoch).

inspectCPU()

Field	Description
cpuType	The model name of the CPU.
cpuModel	The model number of the CPU.
cpuUsr	Time spent normally executing...
cpuNice	

Using SAAF in a Function:

Using SAAF in a function is as simple importing the fram of code. Attributes collected by SAAF will be appended asynchronous functions, this data could be stored into retrieved after the function is finished.

Example Function:

```
from Inspector import *  def myFunction(request):  # Initialize the Inspector and collect da  inspector = Inspector()  inspector.inspectAll()  # Add a "Hello World!" message.  inspector.addAttribute("message", "Hello  # Return attributes collected.
```

SAAF Outline

- Introduction to Serverless Computing
 - Motivation
 - Delivery models and platforms
 - Advantages and challenges
- Serverless Application Analytics Framework (SAAF)
 - Design of SAAF, Supported Languages, Metrics
 - Tools: FaaS Runner, Publish Script
- Analysis Examples with SAAF
 - Programming language comparison, performance modeling
 - Scalability testing, Resource utilization profiling
 - Tracking infrastructure reuse
- Conclusions
- SAAF Demo



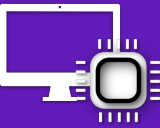
27

Supported Platforms and Languages



28

SAAF Metrics and Design



- Profiles **48** distinct metrics (CPU, memory, I/O utilization), monitors infrastructure state, and observes platform scalability
- Data collection is directed by calling profiling functions
- CPU and Memory metrics are collected from the Linux **procs**
- Cold/Warm infrastructure state is observed by stamping function instances
- Tenancy is determined by introspecting the environment

Example Function:

```
from Inspector import *

def myFunction(request):

    # Initialize the Inspector and collect data.
    inspector = Inspector()
    inspector.inspectAll()

    # Add a "Hello World!" message.
    inspector.addAttribute("message", "Hello " + request['name'])

    # Return attributes collected.
    return inspector.finish()
```

Example Output JSON:

The attributes collect can be customized by changing which functions are called. For more detailed descriptions of each variable and the functions that collect them, please see the framework documentation for each language.

```
{
  "version": 0.2,
  "lang": "python",
  "cpuType": "Intel(R) Xeon(R) Processor @ 2.50GHz",
  "cpuModel": 63,
  "vmuptime": 1551727835,
  "uuid": "d241c618-78d8-48e2-9736-997dc1a931d4",
  "vmID": "tiUCnA",
  "platform": "AWS Lambda",
  "newcontainer": 1,
  "cpuUserDelta": "904",
  "cpuNiceDelta": "0",
  "cpuKrnDelta": "585",
  "cpuIdleDelta": "82428",
  "cpuIowaitDelta": "226",
  "cpuIrqDelta": "0",
  "cpuSoftIrqDelta": "7",
  "vmcpcustealDelta": "1594",
  "frameworkRuntime": 35.72,
  "message": "Hello Fred Smith!",
  "runtime": 38.94
}
```

29

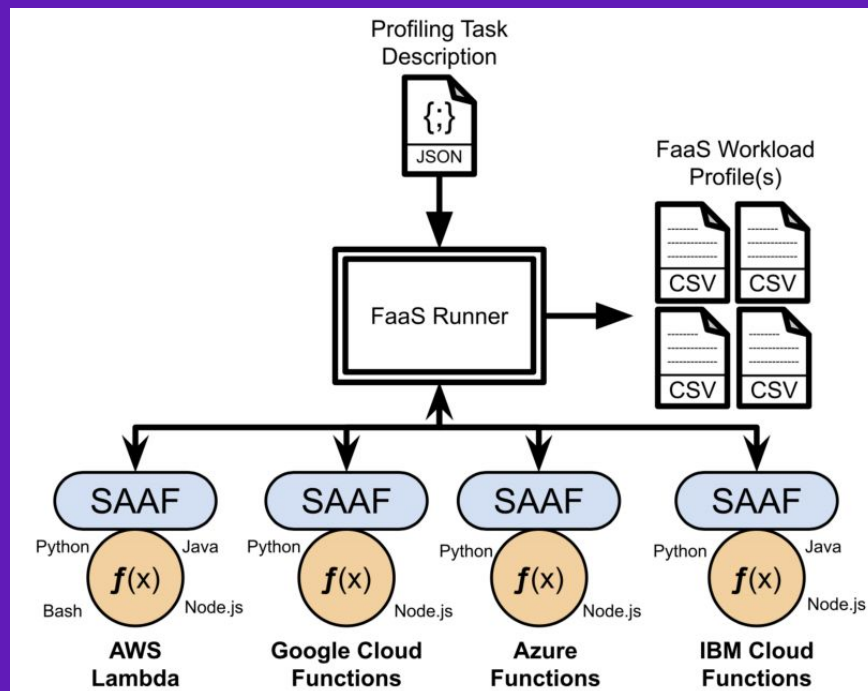
SAAF Outline

- Introduction to Serverless Computing
 - Motivation
 - Delivery models and platforms
 - Advantages and challenges
- Serverless Application Analytics Framework (SAAF)
 - Design of SAAF, Supported Languages, Metrics
 - Tools: FaaS Runner, Publish Script
- Analysis Examples with SAAF
 - Programming language comparison, performance modeling
 - Scalability testing, Resource utilization profiling
 - Tracking infrastructure reuse
- Conclusions
- SAAF Demo



30

Workload Profiling with SAAF and FaaS Runner



31

SAAF Tools: FaaS Runner



- Client for running experiments
- Executes reproducible tests defined by files or command line arguments
 - Automatically change memory settings or redeploy functions
 - Run functions sequentially or concurrently with many threads
 - Run functions synchronously or asynchronously
 - Define payload distribution and creation with inheritance
 - Execute complex pipelines with multiple functions
 - Run multiple iterations of an experiment
- Automatically compile results into a report

32

SAAF + FaaS Runner

- Observations made by FaaS Runner:
 - Network latency
 - Round trip time
 - Runtime concurrency
 - Run/thread IDs to trace pipelines
 - Sum/average/lists of attributes returned by functions
- Combining SAAF and FaaS Runner collects a total of 48 metrics

33

SAAF Tools: Publish Script



Example Hello World Function

```
from Inspector import *  
  
def myFunction(request):  
  
    # Import the module and collect data  
    inspector = Inspector()  
    inspector.inspectAll()  
  
    # Add custom message and finish the function  
    inspector.addAttribute("message", "Hello " + request['name'] + "!")  
  
    inspector.inspectAllDeltas()  
    return inspector.finish()
```



34

SAAF Outline

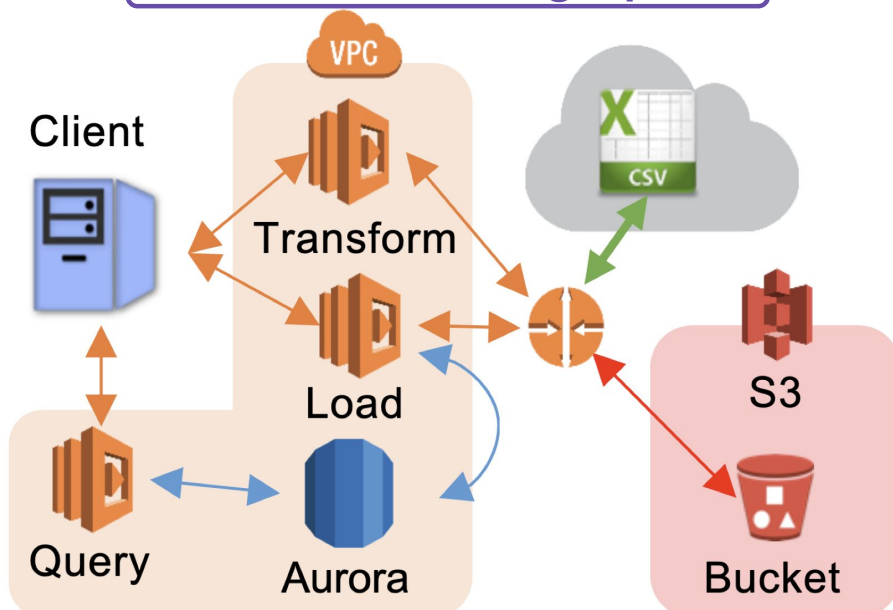
- Introduction to Serverless Computing
 - Motivation
 - Delivery models and platforms
 - Advantages and challenges
- Serverless Application Analytics Framework (SAAF)
 - Design of SAAF, Supported Languages, Metrics
 - Tools: FaaS Runner, Publish Script
- Analysis Examples with SAAF
 - **Programming language comparison** performance modeling
 - Scalability testing, Resource utilization profiling
 - Tracking infrastructure reuse
- Conclusions
- SAAF Demo



35

SAAF: FaaS Programming Languages Comparison

ETL: Data Processing Pipeline



36

FaaS Programming Languages Comparison - II



37

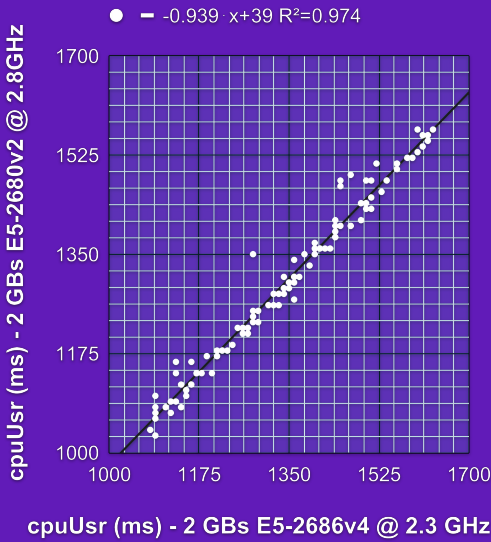
SAAF Outline

- Introduction to Serverless Computing
 - Motivation
 - Delivery models and platforms
 - Advantages and challenges
- Serverless Application Analytics Framework (SAAF)
 - Design of SAAF, Supported Languages, Metrics
 - Tools: FaaS Runner, Publish Script
- Analysis Examples with SAAF
 - Programming language comparison, performance modeling
 - Scalability testing, Resource utilization profiling
 - Tracking infrastructure reuse
- Conclusions
- SAAF Demo



38

Runtime Prediction Scenarios



CPU:

- 256 MBs a1 → a2
- 256 MBs a1 → a3
- 256 MBs a2 → a3
- 512 MBs a1 → a2
- 512 MBs a1 → a3
- 512 MBs a2 → a3
- 1024 MBs a1 → a2
- 1024 MBs a1 → a3
- 1024 MBs a2 → a3
- 2048 MBs a1 → a2
- 2048 MBs a1 → a3
- 2048 MBs a2 → a3

Memory:

- a1 256MBs → 512MBs
- a1 256MBs → 1024MBs
- a1 256MBs → 2048MBs
- a2 256MBs → 512MBs
- a2 256MBs → 1024MBs
- a2 256MBs → 2048MBs
- a3 256MBs → 512MBs
- a3 256MBs → 1024MBs
- a3 256MBs → 2048MBs

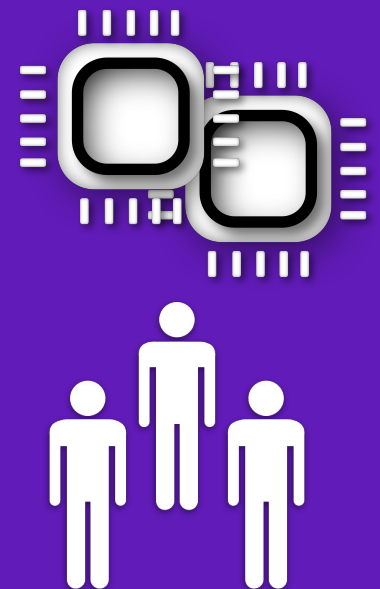
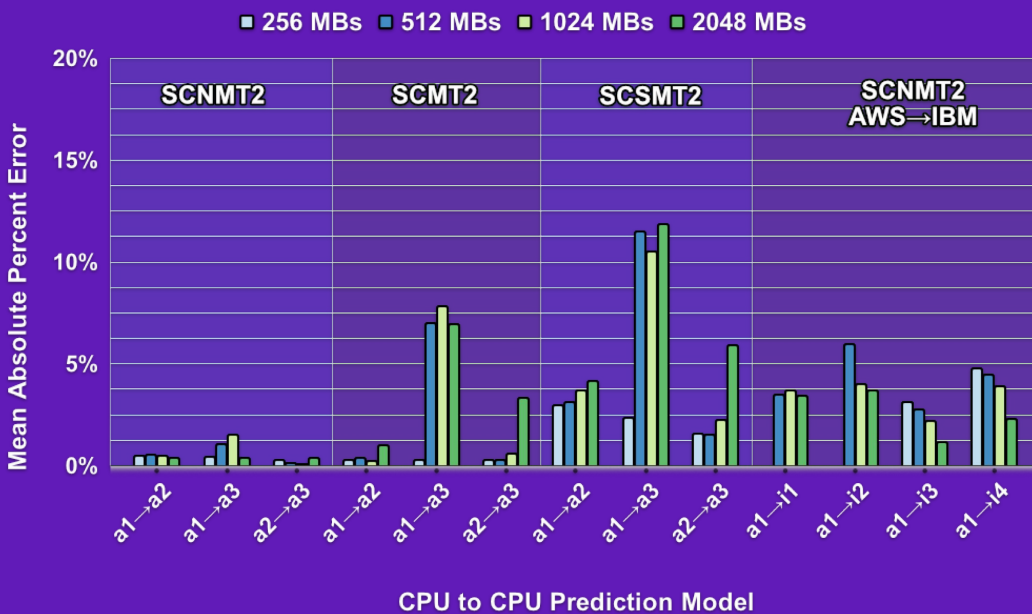
Platform:

- 256MBs a1 → i1
- 256MBs a1 → i2
- 256MBs a1 → i3
- 256MBs a1 → i4
- 512MBs a1 → i1
- 512MBs a1 → i2
- 512MBs a1 → i3
- 512MBs a1 → i4
- 1024MBs a1 → i1
- 1024MBs a1 → i2
- 1024MBs a1 → i3
- 1024MBs a1 → i4
- 2048MBs a1 → i1
- 2048MBs a1 → i2
- 2048MBs a1 → i3
- 2048MBs a1 → i4

Prediction Scenarios

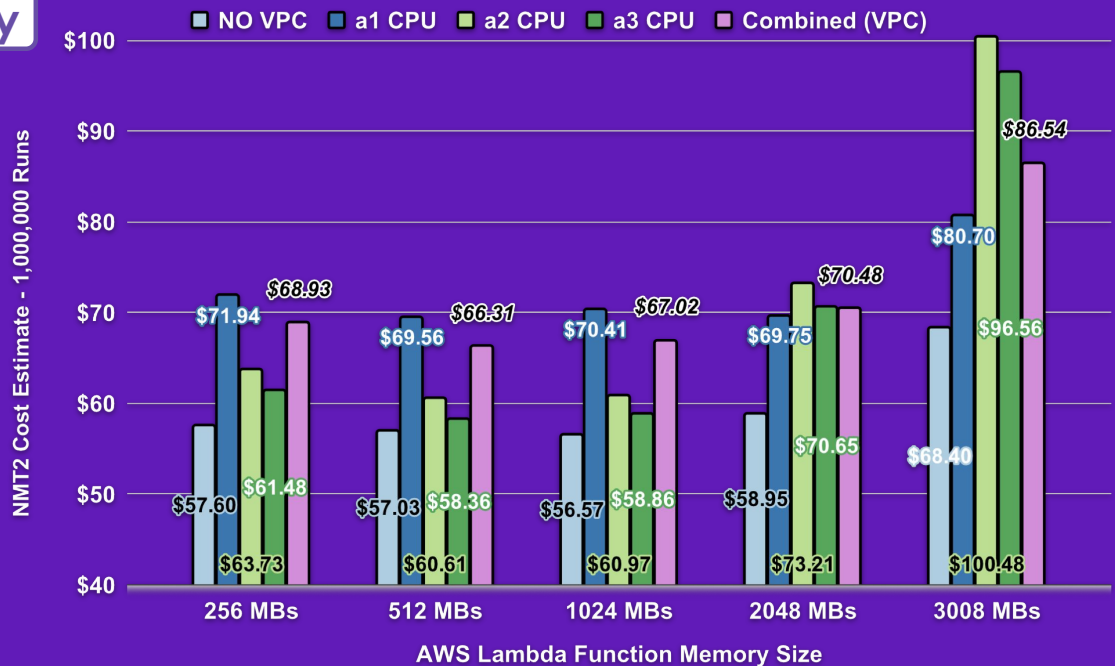
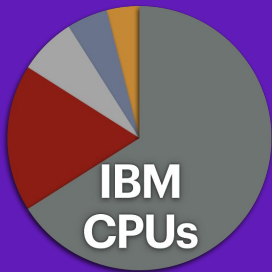
$$\text{Runtime} = \frac{(\text{cpuUsr} + \text{cpuKrn} + \text{cpuIdle} + \text{cpuIOWait} + \text{cpuIntSrcv} + \text{cpuSftIntSrcv})}{(\# \text{ of cores})}$$

Runtime Prediction - Mean Absolute Percentage Error



SAAF: Predicting Hosting Costs

CPU heterogeneity



41

SAAF Outline

- Introduction to Serverless Computing
 - Motivation
 - Delivery models and platforms
 - Advantages and challenges
- Serverless Application Analytics Framework (SAAF)
 - Design of SAAF, Supported Languages, Metrics
 - Tools: FaaS Runner, Publish Script
- Analysis Examples with SAAF
 - Programming language comparison, performance modeling
 - Scalability testing, Resource utilization profiling
 - Tracking infrastructure reuse
- Conclusions
- SAAF Demo



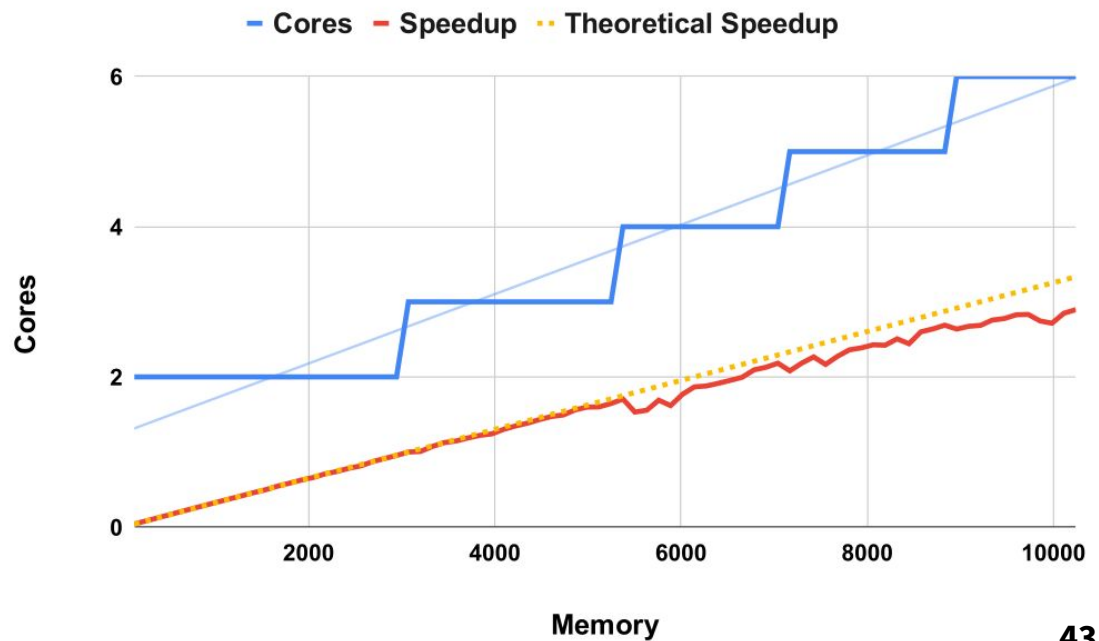
42

AWS Lambda - Scalable Performance Test

Scalability from:
0 to 10 GB

sysbench: prime number
generation w/ 12 threads

Doubling memory
doubles performance
until adding the
4th vCPU:



43

SAAF Outline

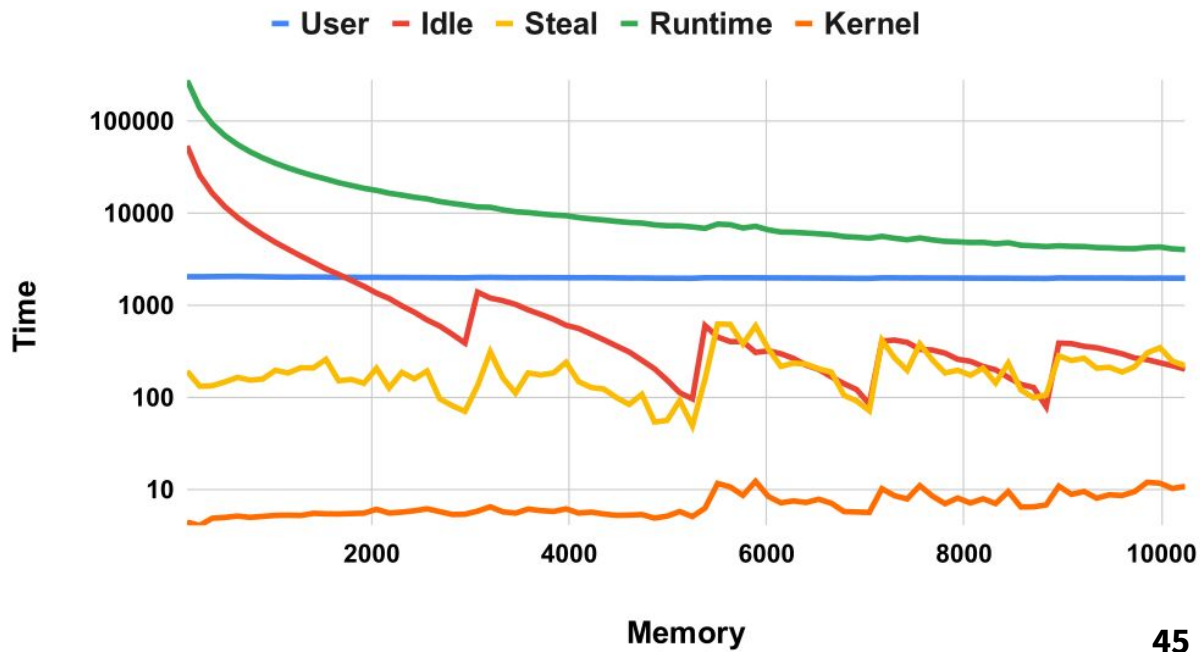
- Introduction to Serverless Computing
 - Motivation
 - Delivery models and platforms
 - Advantages and challenges
- Serverless Application Analytics Framework (SAAF)
 - Design of SAAF, Supported Languages, Metrics
 - Tools: FaaS Runner, Publish Script
- Analysis Examples with SAAF
 - Programming language comparison, performance modeling
 - Scalability testing, Resource utilization profiling
 - Tracking infrastructure reuse
- Conclusions
- SAAF Demo



44

AWS Lambda - Linux CPU time Accounting Metrics

sysbench: prime
number generation
with 12 threads
from 0 to 10 GB
log scale



SAAF Outline

- Introduction to Serverless Computing
 - Motivation
 - Delivery models and platforms
 - Advantages and challenges
- Serverless Application Analytics Framework (SAAF)
 - Design of SAAF, Supported Languages, Metrics
 - Tools: FaaS Runner, Publish Script
- Analysis Examples with SAAF
 - Programming language comparison, performance modeling
 - Scalability testing, Resource utilization profiling
 - Tracking infrastructure reuse
- Conclusions
- SAAF Demo



AWS Lambda - Infrastructure Reuse Testing

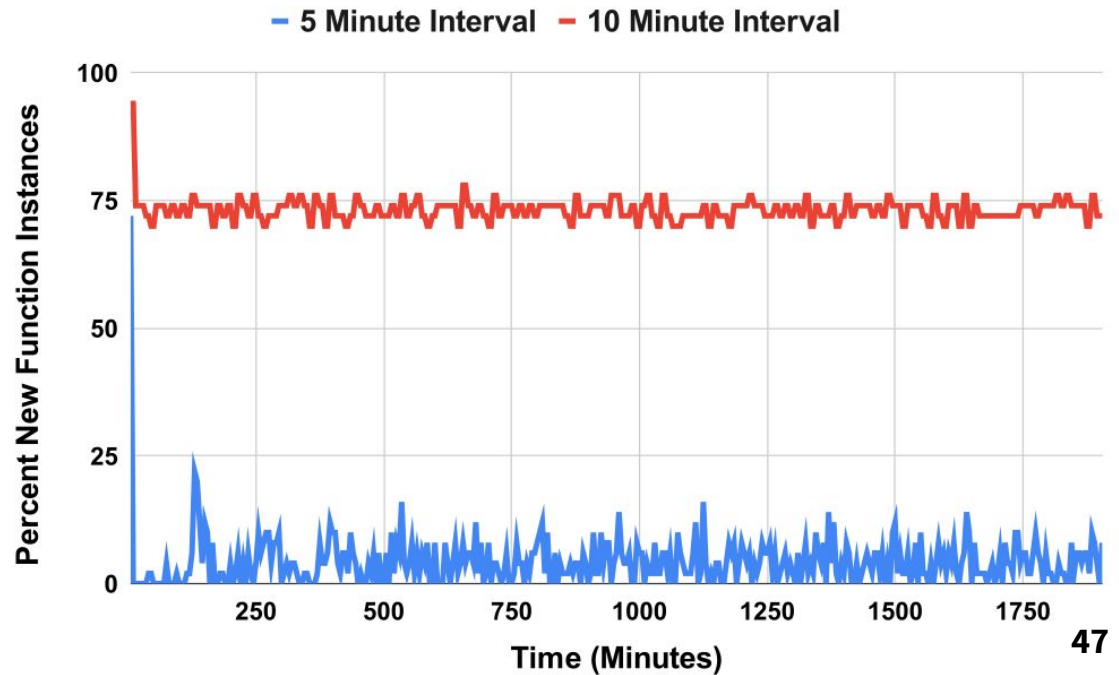
COLD infrastructure is common with the serverless freeze-thaw life cycle

Experiment:

50 concurrent calls

5-min vs 10-min delay

Evaluate % function instances



SAAF Outline

- Introduction to Serverless Computing
 - Motivation
 - Delivery models and platforms
 - Advantages and challenges
- Serverless Application Analytics Framework (SAAF)
 - Design of SAAF, Supported Languages, Metrics
 - Tools: FaaS Runner, Publish Script
- Analysis Examples with SAAF
 - Programming language comparison, performance modeling
 - Scalability testing, Resource utilization profiling
 - Tracking infrastructure reuse
- Conclusions
- SAAF Demo



Conclusions

SAAF's goal is to enable developers and researchers to make educated observations into the factors that impact performance on FaaS platforms

- Design goals:
 - Easy to implement and deploy
 - Low overhead and minimal dependencies
 - Cross platform/language support
 - A complete development workflow with SAAF + FaaS Runner:
 - Development -> Deployment -> Testing -> Data Analysis
 - Available for anyone

49

SAAF Outline

- Introduction to Serverless Computing
 - Motivation
 - Delivery models and platforms
 - Advantages and challenges
- Serverless Application Analytics Framework (SAAF)
 - Design of SAAF, Supported Languages, Metrics
 - Tools: FaaS Runner, Publish Script
- Analysis Examples with SAAF
 - Programming language comparison, performance modeling
 - Scalability testing, Resource utilization profiling
 - Tracking infrastructure reuse
- Conclusions
- SAAF Demo



50

SAAF Demonstration

- SAAF Overview
- Writing a Function with SAAF
- Deploying Functions to all Platforms
- Running Experiments with FaaS Runner
- Generating Reports
- Working with Results in R
- Interactive FaaS with Jupyter

51

Thank You!

Questions or comments?

Please email: rcording@uw.edu or wllloyd@uw.edu

Download the Serverless Application Analytics Framework:
github.com/wllloyd/saaf

SAAF Online Tutorial:
<https://github.com/wllloyd/SAAF/blob/master/tutorial>

Paper Link:
http://faculty.washington.edu/wllloyd/papers/ICPE_SAAF_proof.pdf

This research is supported by NSF Advanced Cyberinfrastructure Research Program (OAC-1849970), NIH grant R01GM126019, and the AWS Cloud Credits for Research program. 52

text