



Improving Application Migration to Serverless Computing Platforms: Latency Mitigation with Keep-Alive Workloads

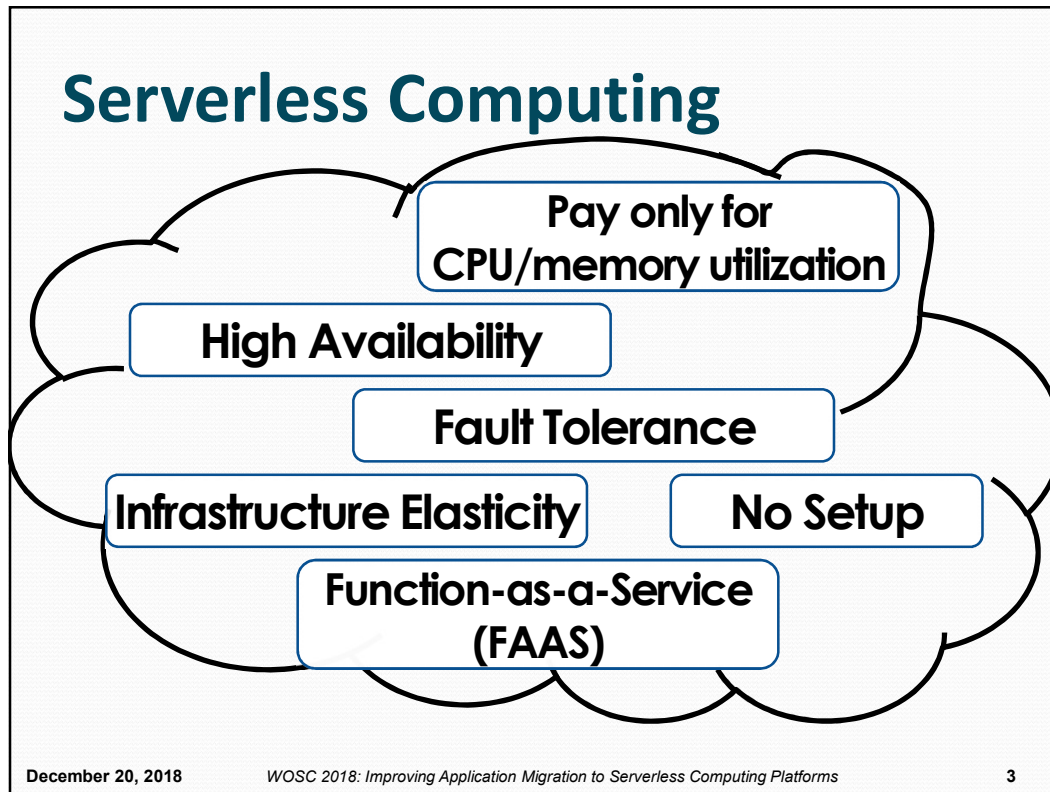
Minh Vu[#], Baojia Zhang[#], Olaf David, George Leavesley,
Wes Lloyd¹

December 20, 2018

School of Engineering and Technology,
University of Washington, Tacoma, Washington USA
WOSC 2018: 4th IEEE Workshop on Serverless Computing (UCC 2018)

Outline

- Background
- Research Questions
- Experimental Workloads
- Experiments/Evaluation
- Conclusions



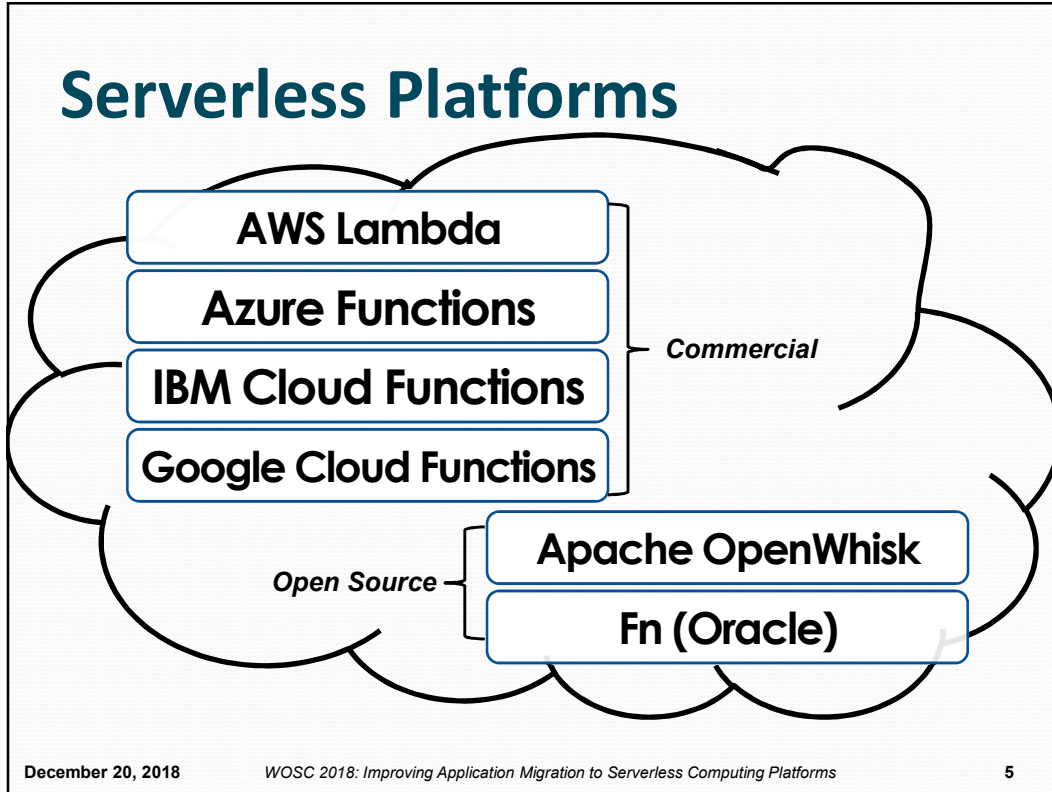
Serverless Computing

Why Serverless Computing?

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

...they are built into the platform

December 20, 2018 WOSC 2018: Improving Application Migration to Serverless Computing Platforms 4



Serverless Computing

Research Challenges

Serverless Computing
Deploy Applications Without Fiddling With Servers

Image from: <https://mobisoftinfotech.com/resources/blog/serverless-computing-deploy-applications-without-fiddling-with-servers/>

6

Serverless Computing Research Challenges

- Memory reservation
- Infrastructure freeze/thaw cycle
- Vendor architectural lock-in
- Pricing obfuscation
- Service composition

December 20, 2018

WOSC 2018: Improving Application Migration to Serverless Computing Platforms

7

Serverless Computing Research Challenges

- Memory reservation
- Infrastructure freeze/thaw cycle
- Vendor architectural lock-in
- Pricing obfuscation
- Service composition

December 20, 2018

WOSC 2018: Improving Application Migration to Serverless Computing Platforms

8

Memory Reservation Question...



- Lambda memory reserved for functions
- UI provides “slider bar” to set function’s memory allocation
- Resource capacity (CPU, disk, network) coupled to slider bar:
*“every **doubling** of memory, **doubles** CPU...”*
- **But how much memory do model services require?**

▼ Basic settings

Memory (MB) Info
 Your function is allocated CPU proportional to the memory configured.

1536 MB

Timeout Info
 3 min 0 sec

Description

Performance

December 20, 2018

WOSC 2018: Improving Application Migration to Serverless Computing Platforms

9

Infrastructure Freeze/Thaw Cycle



- Unused infrastructure is deprecated
 - *But after how long?*
- AWS Lambda: Bare-metal hosts, firecracker micro-VMs
- Infrastructure states: <https://firecracker-microvm.github.io/>
- **Provider-COLD / Host-COLD**
 - Function package built/transferred to Hosts
- **Container-COLD (firecracker micro-VM)**
 - Image cached on Host
- **Container-WARM (firecracker micro-VM)**
 - “Container” running on Host

Performance



Image from: Denver7 – The Denver Channel News

December 20, 2018

WOSC 2018: Improving Application Migration to Serverless Computing Platforms

10

Outline

- Background
- **Research Questions**
- Experimental Workloads
- Experiments/Evaluation
- Conclusions

December 20, 2018

WOSC 2018: Improving Application Migration to Serverless Computing Platforms

11

Research Questions

- RQ1:** **PERFORMANCE:** What are the performance implications for application migration? How does memory reservation size impact performance when coupled to CPU power?
- RQ2:** **SCALABILITY:** For application migration what performance implications result from scaling the number of concurrent clients? How is scaling affected when infrastructure is allowed to go cold?

December 20, 2018

WOSC 2018: Improving Application Migration to Serverless Computing Platforms

12

Research Questions - 2

- RQ3:** **COST:** For hosting large parallel service workloads, how does memory reservation size, impact hosting costs when coupled to CPU power?
- RQ4:** **PERSISTING INFRASTRUCTURE:** How effective are automatic triggers at retaining serverless infrastructure to reduce performance latency from the serverless freeze/thaw cycle?

December 20, 2018

WOSC 2018: Improving Application Migration to Serverless Computing Platforms

13

Outline

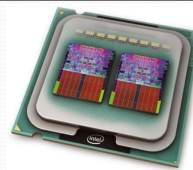
- Background
- Research Questions
- **Experimental Workloads**
- Experiments/Evaluation
- Conclusions

December 20, 2018

WOSC 2018: Improving Application Migration to Serverless Computing Platforms

14

AWS Lambda PRMS Modeling Service



- PRMS: deterministic, distributed-parameter model
- Evaluate impact of combinations of precipitation, climate, and land use on stream flow and general basin hydrology (Leavesley et al., 1983)



- Java based PRMS, Object Modelling System (OMS) 3.0
- Approximately ~11,000 lines of code
- Model service is 18.35 MB compressed as a Java JAR file
- Data files hosted using Amazon S3 (object storage)

Goal: quantify performance and cost implications of memory reservation size and scaling for model service deployment to AWS Lambda

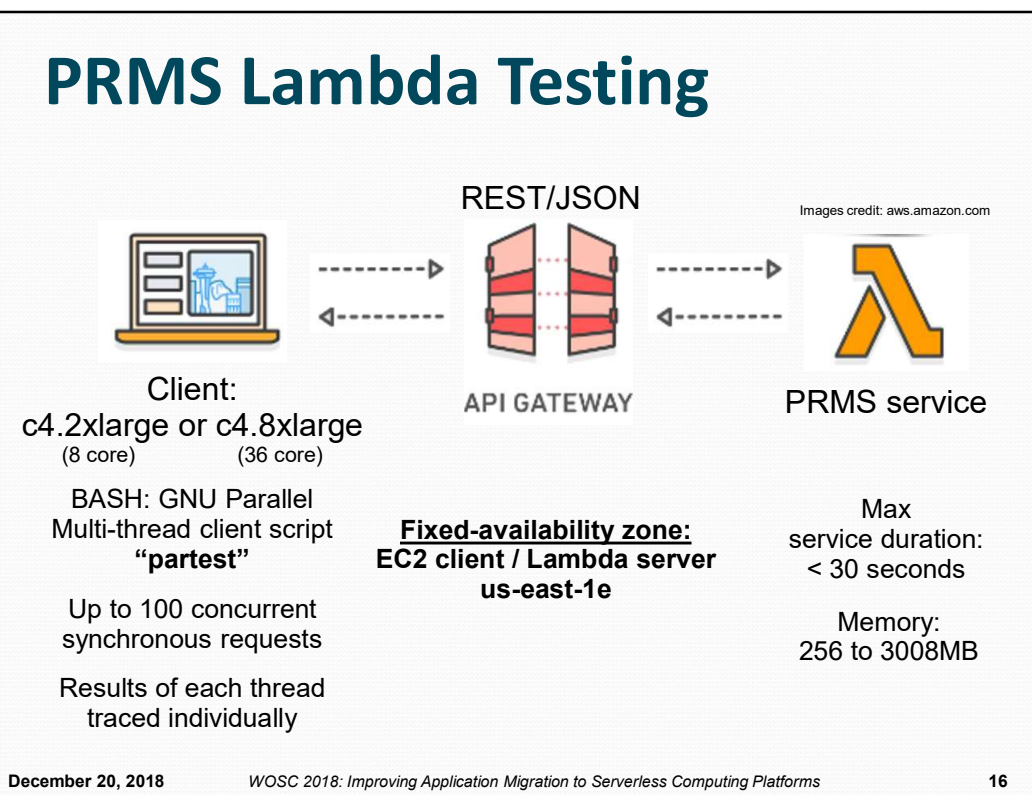


December 20, 2018

WOSC 2018: Improving Application Migration to Serverless Computing Platforms

15

PRMS Lambda Testing

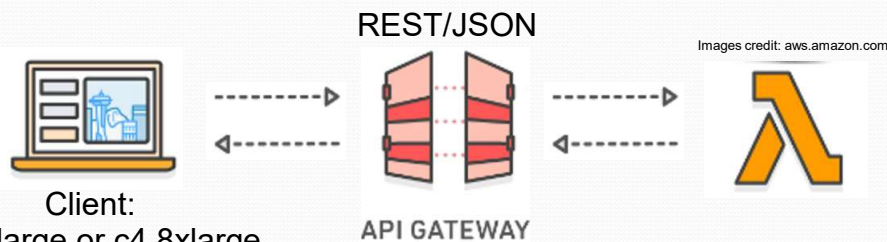


December 20, 2018

WOSC 2018: Improving Application Migration to Serverless Computing Platforms

16

PRMS Lambda Testing - 2



Automatic Metrics Collection⁽¹⁾:

New vs. Recycled Containers/VMs

of requests per container/VM

Avg. performance per container/VM

Avg. performance workload

Standard deviation of requests per container/VM

Container Identification

UUID → /tmp file

VM Identification

btime → /proc/stat

Linux CPU metrics

⁽¹⁾Lloyd, W., Ramesh, S., Chinthalapati, S., Ly, L., & Pallickara, S. (April 2018). Serverless computing: An investigation of factors influencing microservice performance. In Cloud Engineering (IC2E), 2018 IEEE International Conference on (pp. 159-169). IEEE.

Outline

- Background
- Research Questions
- Experimental Workloads
- Experiments/Evaluation
- Conclusions

RQ-1: Performance

Infrastructure

What are the performance implications of memory reservation size ?

RQ-1: AWS Lambda Memory Reservation Size

▼ Basic settings

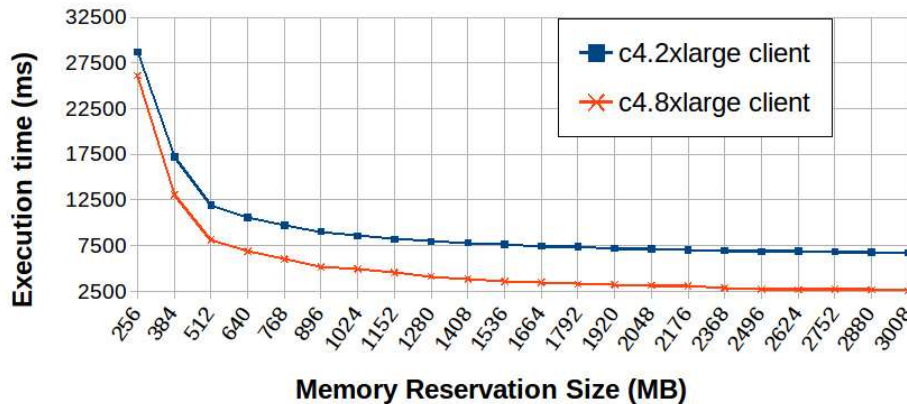
Memory (MB) info
Your function is allocated CPU proportional to the memory configured.
1536 MB

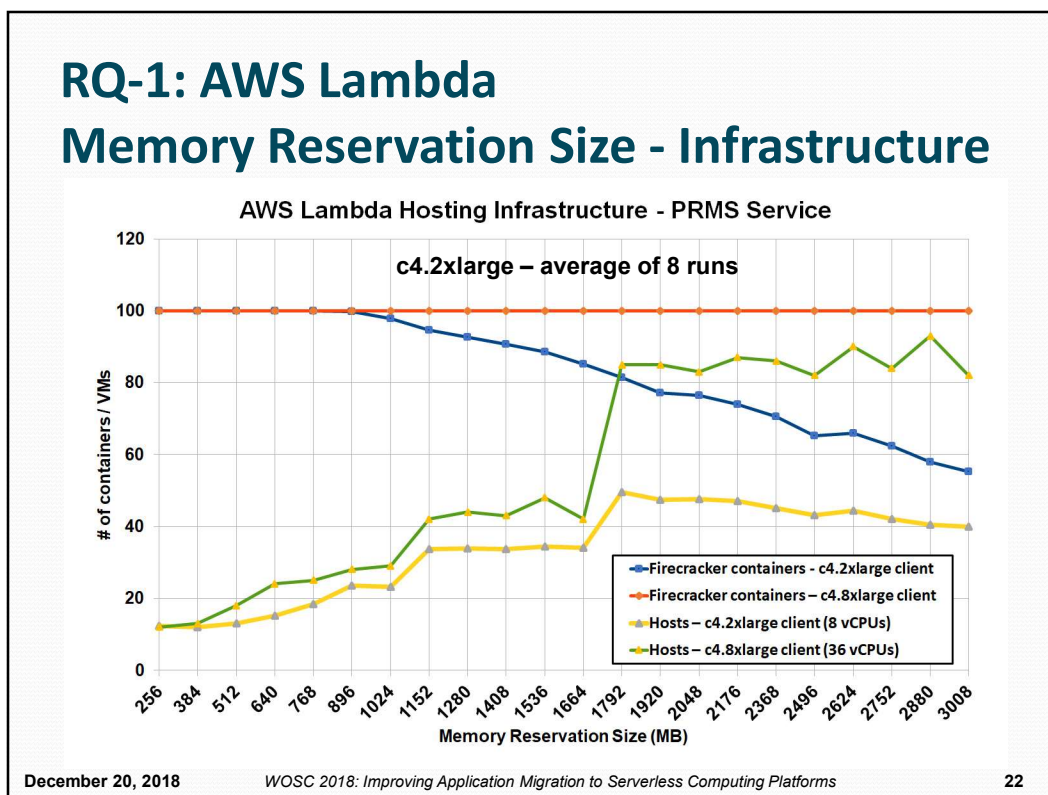
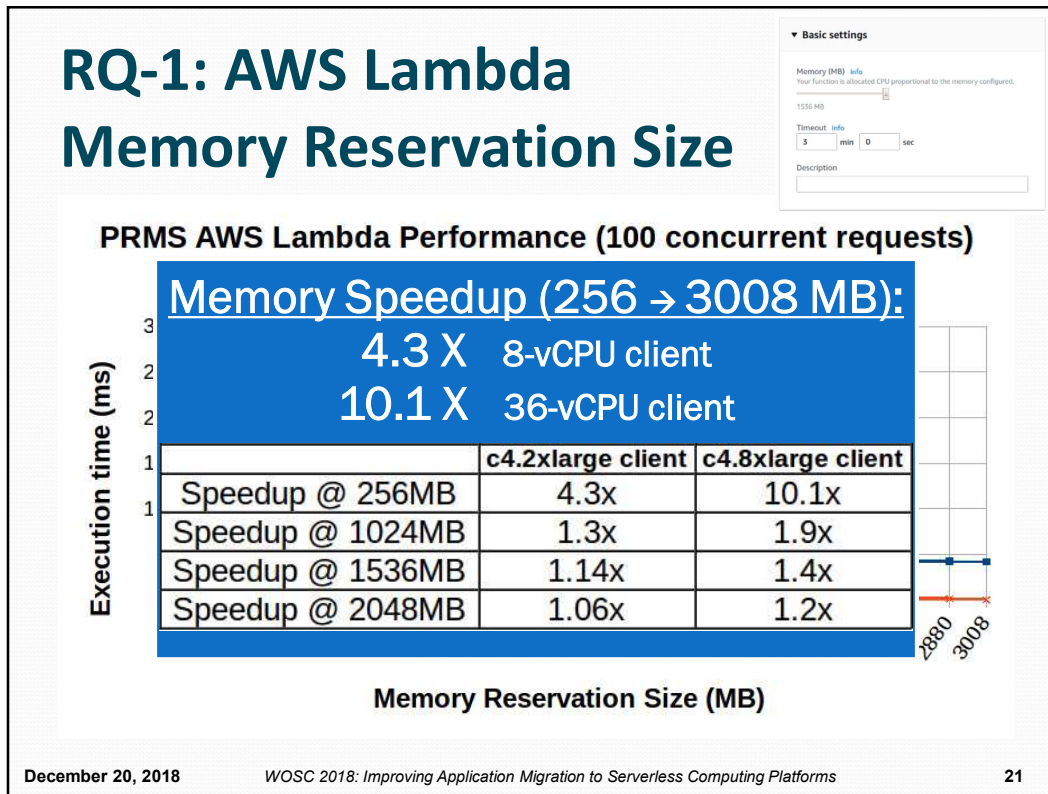
Timeout info
3 min 0 sec

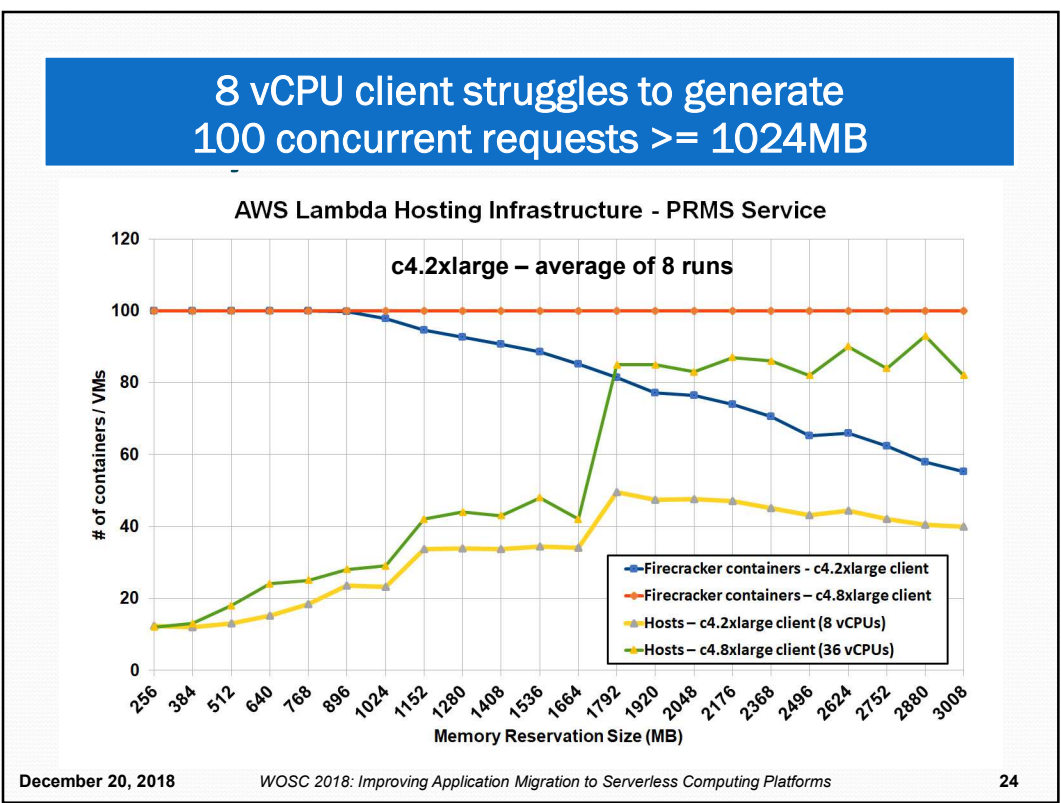
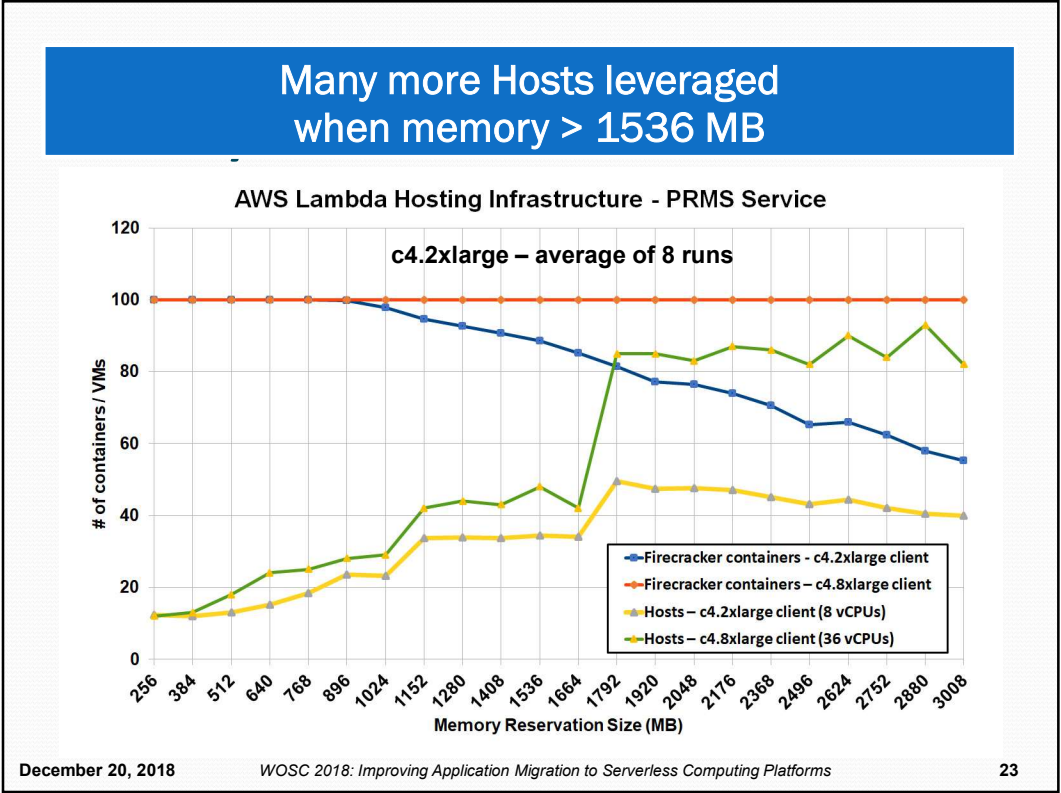
Description

PRMS AWS Lambda Performance (100 concurrent requests)

c4.2xlarge – average of 8 runs







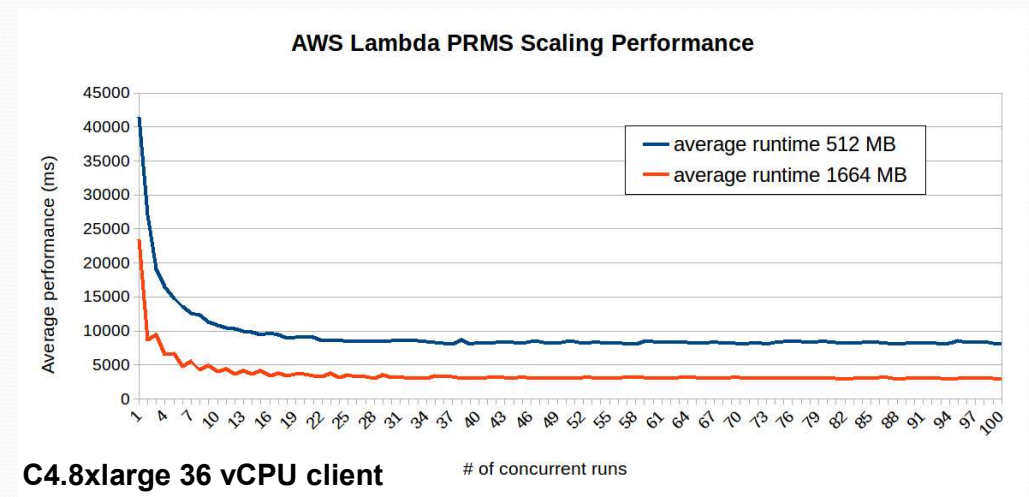
RQ-2: Scalability

How does performance change when increasing the number of concurrent users ?

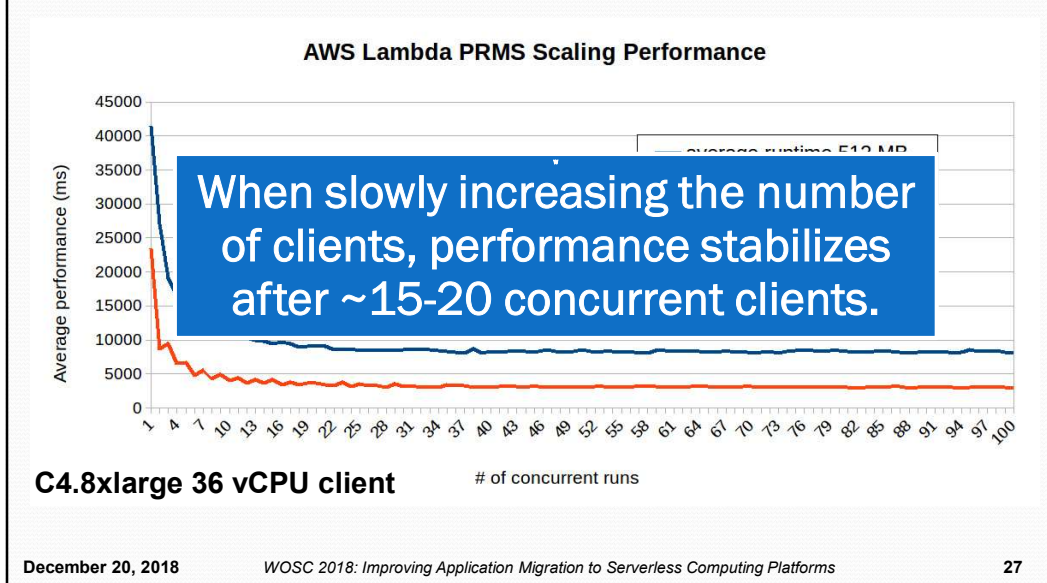
(scaling-up, totally cold, and warm)

25

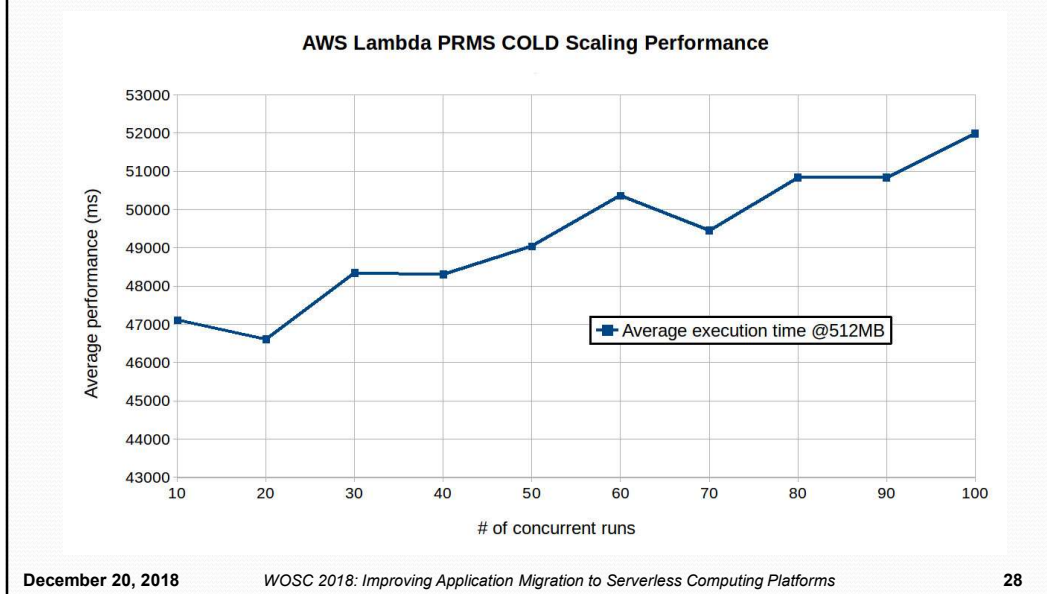
RQ-2: AWS Lambda PRMS Scaling Performance



RQ-2: AWS Lambda PRMS Scaling Performance



RQ-2: AWS Lambda Cold Scaling Performance



RQ-3: Cost

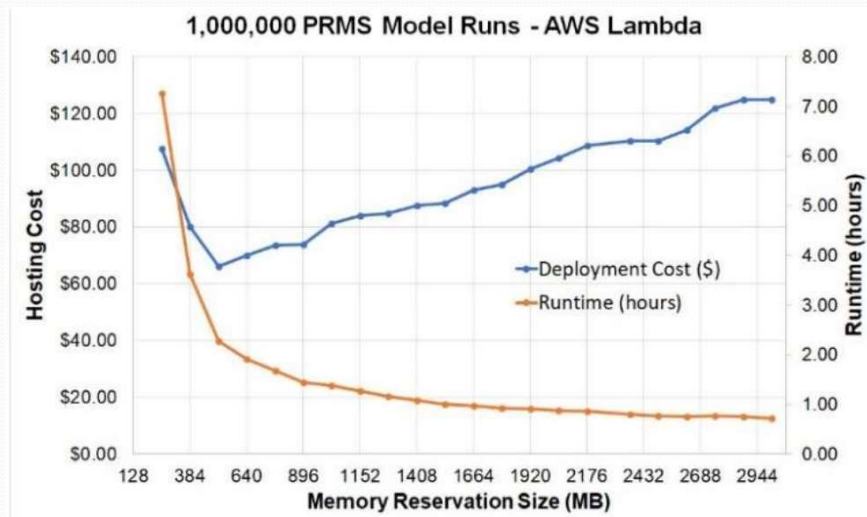
What are the costs of hosting PRMS using a FaaS platform in comparison to IaaS?

29

RQ-3: IaaS (EC2) Hosting Cost 1,000,000 PRMS runs

- Using a 2 vCPU c4.large EC2 VM
 - 2 concurrent client calls, no scale-up
- Estimated time: 347.2 hours, **14.46 days**
 - Assume average exe time of 2.5 sec/run
- Hosting cost @ 10¢/hour = **\$34.72**

RQ-3: FaaS Hosting Cost 1,000,000 PRMS runs

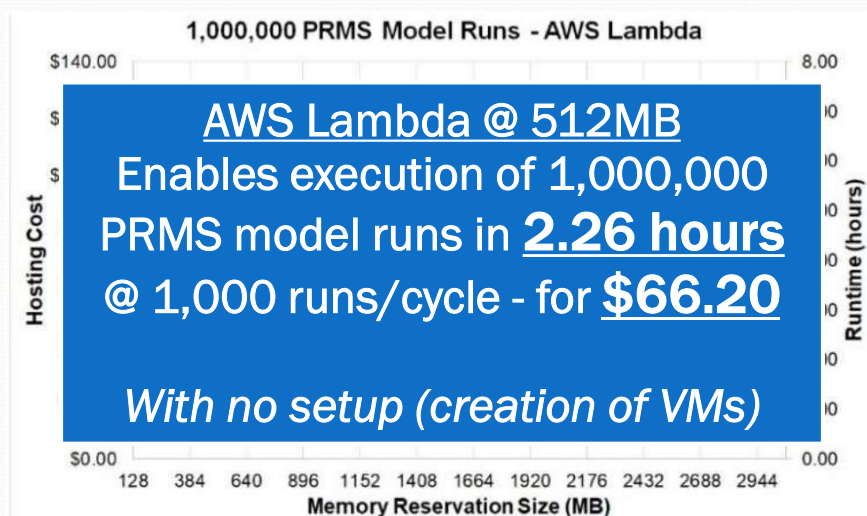


December 20, 2018

WOSC 2018: Improving Application Migration to Serverless Computing Platforms

31

RQ-3: FaaS Hosting Cost 1,000,000 PRMS runs



December 20, 2018

WOSC 2018: Improving Application Migration to Serverless Computing Platforms

32

RQ-4: Persisting Infrastructure

How effective are automatic triggers at retaining serverless infrastructure to reduce performance latency from the serverless freeze/thaw cycle?

33

RQ-4: Persisting Infrastructure

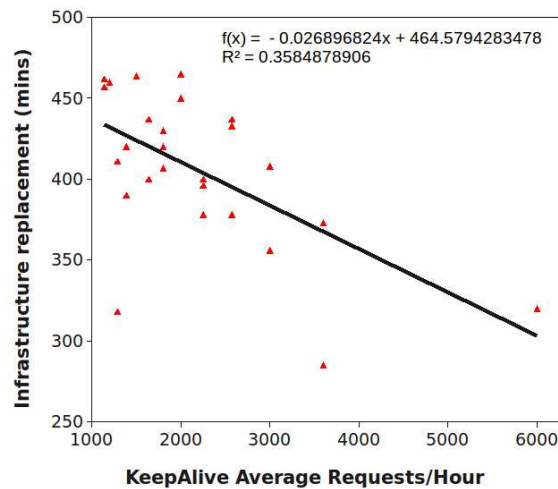
- Goal: preserve 100 firecracker containers for 24hrs
 - Mitigate cold start latency
- Memory: 192, 256, 384, 512 MB
- All initial host infrastructure replaced between ~4.75 – 7.75 hrs
- Replacement cycle (start→finish): ~2 hrs
- Infrastructure generations performance variance observed from: -14.7% to 19.4% (Δ 34%)
- Average performance variance larger for lower memory sizes: 9% (192MB), 3.6% (512MB)

December 20, 2018

WOSC 2018: Improving Application Migration to Serverless Computing Platforms

34

RQ-4: Persisting Infrastructure AWS Lambda: time to infrastructure replacement vs. memory reservation size



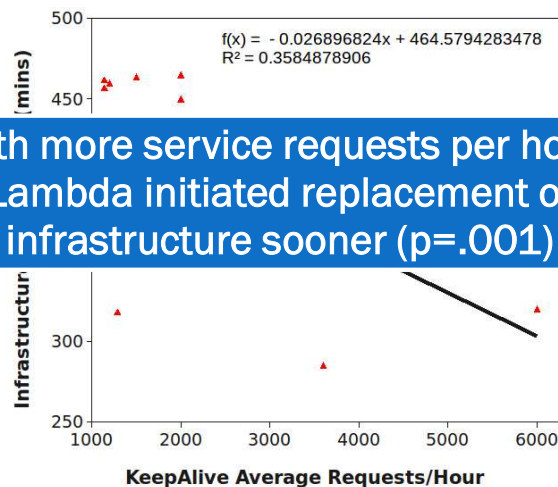
Memory sizes tested: 192, 256, 384, 512 MB

December 20, 2018

WOSC 2018: Improving Application Migration to Serverless Computing Platforms

35

RQ-4: Persisting Infrastructure AWS Lambda: time to infrastructure replacement vs. memory reservation size



With more service requests per hour, Lambda initiated replacement of infrastructure sooner (p=.001)

Memory sizes tested: 192, 256, 384, 512 MB

December 20, 2018

WOSC 2018: Improving Application Migration to Serverless Computing Platforms

36

RQ-4: Persisting Infrastructure

Keep-Alive Infrastructure Preservation

- PRMS Service: parameterize for “ping”
 - Perform sleep (idle CPU) – do not run model
 - Provides delay to overlap (n=100) parallel requests to preserve infrastructure
- Ping intervals: tested 3, 4, and 5-minutes
- VM Keep-Alive client:
c4.8xlarge 36 vCPU instance: ~4.5s sleep
- CloudWatch Keep-Alive client:
100 rules x 5 targets: 5-s sleep

December 20, 2018

WOSC 2018: Improving Application Migration to Serverless Computing Platforms

37

RQ-4: Keep-Alive Client Summary

| Client type | c4.8xlarge VM | c4.8xlarge VM | CloudWatch | CloudWatch |
|------------------------------------------------|---------------|---------------|------------|------------|
| Ping interval | 5 min | 3 min | 5 min | 4min |
| Keep-Alive calls/batch | 100 | 100 | 500 | 500 |
| Slowdown vs. WARM | 13.3% | 0.7% | 11.6% | 35.0% |
| Speedup vs. COLD | 4.03x | 4.53x | 4.1x | 3.4x |
| Test runs | 32 | 32 | 26 | 17 |
| Test duration (hrs) | 24 | 24 | 18 | 12 |
| Average new Lambda firecracker containers/test | 2.41 | 0.38 | 5.42 | 14.71 |
| Keep-Alive runtime avg (ms) | 4492 | 4463 | 5200 | 5200 |
| Memory (GB-sec/hour) | 2695 | 4463 | 15600 | 19500 |
| Keep-Alive cost/year | \$4,484.00 | \$4,494.76 | \$2,278.06 | \$2,847.57 |

Keep-Alive clients can support trading off cost for performance for preserving FaaS infrastructure to mitigate cold start latency

December 20, 2018

WOSC 2018: Improving Application Migration to Serverless Computing Platforms

38

RQ-4: Keep-Alive Client Summary

| Client type | c4.8xlarge VM | c4.8xlarge VM | CloudWatch | CloudWatch |
|------------------------------------------------|---------------|---------------|-------------|-------------|
| Ping interval | 5 min | 3 min | 5 min | 4min |
| Keep-Alive calls/batch | 100 | 100 | 500 | 500 |
| Slowdown vs. WARM | 13.3% | 0.7% | 11.6% | 35.0% |
| Speedup vs. COLD | 4.03x | 4.53x | 4.1x | 3.4x |
| Test runs | 32 | 32 | 26 | 17 |
| Test duration (hrs) | 24 | 24 | 18 | 12 |
| Average new Lambda firecracker containers/test | 2.41 | 0.38 | 5.42 | 14.71 |
| Keep-Alive runtime avg (ms) | 4492 | 4463 | 5200 | 5200 |
| Memory (GB-sec/hour) | 2695 | 4463 | 15600 | 19500 |
| Keep-Alive cost/year | \$4,484.00 | \$4,494.76 | \$2,278.06 | \$2,847.57 |

Keep-Alive clients can support trading off cost for performance for preserving FaaS infrastructure to mitigate cold start latency

December 20, 2018

WOSC 2018: Improving Application Migration to Serverless Computing Platforms

39

RQ-4: Keep-Alive Client Summary

| Client type | c4.8xlarge VM | c4.8xlarge VM | CloudWatch | CloudWatch |
|------------------------------------------------|---------------|---------------|-------------|-------------|
| Ping interval | 5 min | 3 min | 5 min | 4min |
| Keep-Alive calls/batch | 100 | 100 | 500 | 500 |
| Slowdown vs. WARM | 13.3% | 0.7% | 11.6% | 35.0% |
| Speedup vs. COLD | 4.03x | 4.53x | 4.1x | 3.4x |
| Test runs | 32 | 32 | 26 | 17 |
| Test duration (hrs) | 24 | 24 | 18 | 12 |
| Average new Lambda firecracker containers/test | 2.41 | 0.38 | 5.42 | 14.71 |
| Keep-Alive runtime avg (ms) | 4492 | 4463 | 5200 | 5200 |
| Memory (GB-sec/hour) | 2695 | 4463 | 15600 | 19500 |
| Keep-Alive cost/year | \$4,484.00 | \$4,494.76 | \$2,278.06 | \$2,847.57 |

Keep-Alive clients can support trading off cost for performance for preserving FaaS infrastructure to mitigate cold start latency

December 20, 2018

WOSC 2018: Improving Application Migration to Serverless Computing Platforms

40

Outline

- Background
- Research Questions
- Experimental Workloads
- Experiments/Evaluation
- **Conclusions**

December 20, 2018

WOSC 2018: Improving Application Migration to Serverless Computing Platforms

41

Conclusions



- **RQ-1 Memory Reservation Size:**
 - MAX memory: 10x speedup, 7x more hosts
- **RQ-2 Scaling Performance:**
 - 1+ scale-up near warm, COLD scale-up is slow
- **RQ-3 Cost**
 - m4.large \$35 (14d), Lambda \$66 (2.3 hr), \$125 (42 min)
- **RQ-4 Persisting Infrastructure (Keep-Alive)**
 - c4.8xlarge VM \$4,484/yr (13.3% slowdown vs warm, 4x ↑), CloudWatch \$2,278/yr (11.6% slowdown vs warm, 4.1x ↑)

December 20, 2018

WOSC 2018: Improving Application Migration to Serverless Computing Platforms

42

Questions

