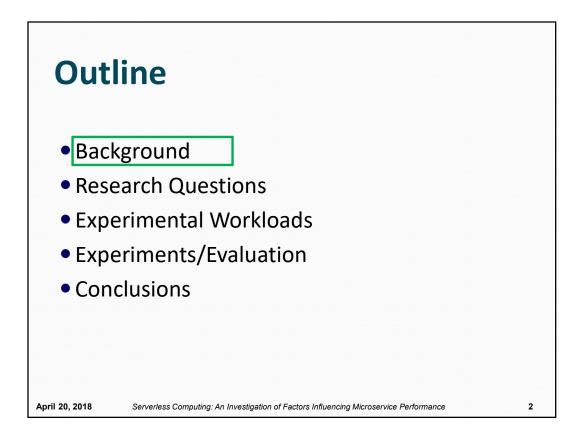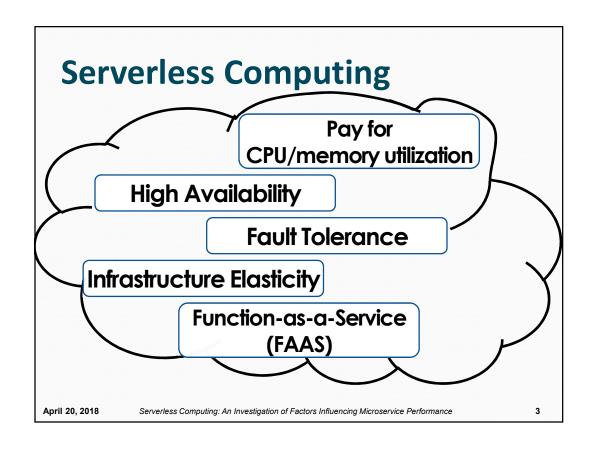# Serverless Computing: An Investigation of Factors Influencing Microservice Performance

Wes Lloyd, Shruti Ramesh, Swetha Chinthalapati,
Lan Ly, Shrideep Pallickara

April 20, 2018

Institute of Technology,
University of Washington, Tacoma, Washington USA
*IC2E 2018*: IEEE International Conference on Cloud Engineering

# Outline

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

April 20, 2018          *Serverless Computing: An Investigation of Factors Influencing Microservice Performance*          2

# Serverless Computing

Pay for
CPU/memory utilization

High Availability

Fault Tolerance

Infrastructure Elasticity

Function-as-a-Service
(FAAS)

# Serverless Computing

**Why Serverless Computing?**

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

*…they are built into the platform*

# Serverless Platforms

AWS Lambda

Azure Functions

IBM Cloud Functions

Google Cloud Functions

*Commercial*

*Open Source* — Apache OpenWhisk

Fn (Oracle)

**April 20, 2018**          *Serverless Computing: An Investigation of Factors Influencing Microservice Performance*          **5**
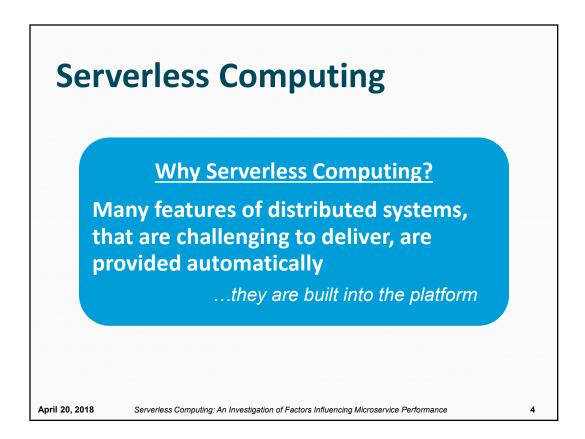
# Serverless Computing
## Research Challenges

## Serverless Computing
Deploy Applications Without
Fiddling With Servers

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

6

# Vendor architectural lock-in

- Serverless software architecture requires external services/components



**Example:** *Weather Application*

**Client**

*Lambda is triggered*

35° C

S3

API GATEWAY

DYNAMODB

*Front-end code for weather app hosted in S3*

*User clicks on link to get local weather information*

*App makes REST API call to endpoint*

*Lambda runs code to retrieve local weather information and returns data back to user*

Images credit: aws.amazon.com

- Increased dependencies → increased hosting costs

# Serverless Pricing Model

- **EXAMPLE:**        AWS Lambda Pricing

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

- Afterwards:        *obfuscated pricing (AWS Lambda):*
                     $0.0000002 per request
                     $0.000000208 to rent 128MB / 100-ms

# Serverless Computing
## Memory reservation question…

- Lambda memory reserved for functions
- UI provides "slider bar" to set function's memory allocation
- CPU power coupled to slider bar: "*every **doubling** of memory, **doubles** CPU…*"
- **But how much memory does code require?**

**Performance**

# Service Composition

- **How should application code be composed for deployment to serverless computing platforms?**

Monolithic

Client flow control, 4 functions

Server flow control, 3 functions

- Recommended practice: Decompose code into many microservices
- Platform limits: code + libraries ~256MB
- **How does composition impact number of invocations, and memory utilization?**

**Performance**

# Freeze/Thaw Cycle

- Unused infrastructure is deprecated
  - *But after how long?*
- Infrastructure: VMs, "containers"

**Performance**

- **Provider-COLD / VM-COLD**
  - "Container" images - built/transferred to VMs
- **Container-COLD**
  - Image cached on VM
- **Container-WARM**
  - "Container" running on VM



10 MINUTES NON-STOP NEWS | FREEZE-THAW CYCLE CAUSING POTHOLES

Image from: Denver7 – The Denver Channel News

April 20, 2018     *Serverless Computing: An Investigation of Factors Influencing Microservice Performance*     **11**

# Serverless Computing Research Challenges

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

April 20, 2018     *Serverless Computing: An Investigation of Factors Influencing Microservice Performance*     **12**

# Outline

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

April 20, 2018         *Serverless Computing: An Investigation of Factors Influencing Microservice Performance*         **13**

# Research Questions

**RQ1:**     What are the performance implications of infrastructure **_elasticity_** for serverless computing?
*(e.g. COLD vs. WARM performance)*

**RQ2:**     How does **_load balancing_** vary in serverless computing?  How do computational requests impact load balancing, and ultimately performance?
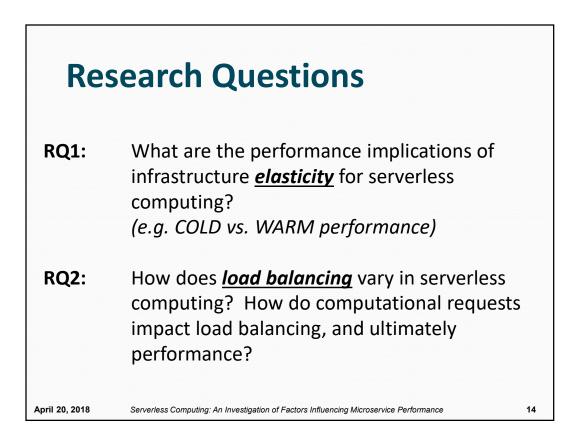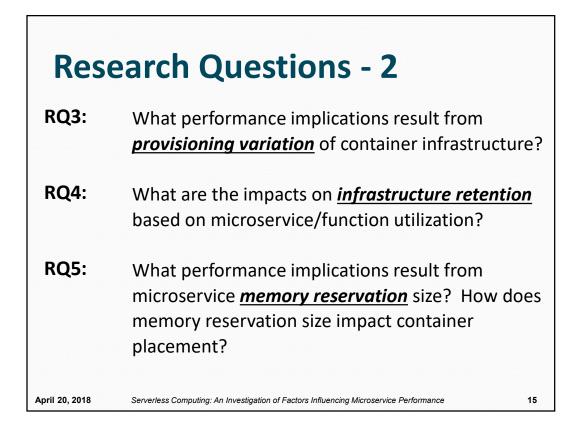
April 20, 2018         *Serverless Computing: An Investigation of Factors Influencing Microservice Performance*         **14**

# Research Questions - 2

**RQ3:** What performance implications result from
_**provisioning variation**_ of container infrastructure?

**RQ4:** What are the impacts on _**infrastructure retention**_
based on microservice/function utilization?

**RQ5:** What performance implications result from
microservice _**memory reservation**_ size?  How does
memory reservation size impact container
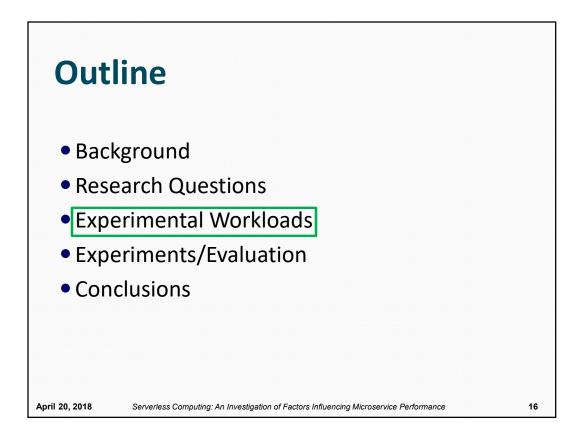placement?

April 20, 2018       Serverless Computing: An Investigation of Factors Influencing Microservice Performance       15

# Outline

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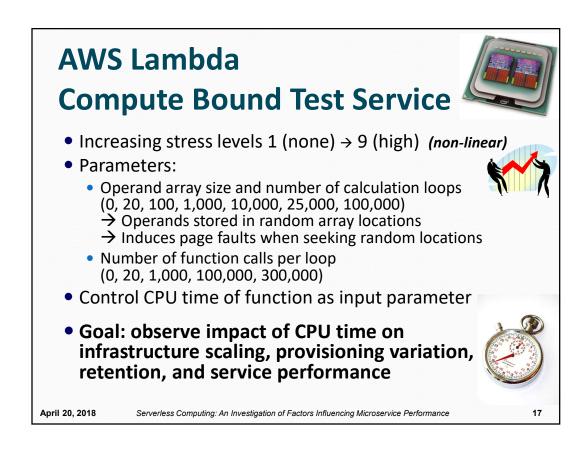- Conclusions

April 20, 2018       Serverless Computing: An Investigation of Factors Influencing Microservice Performance       16

# AWS Lambda
# Compute Bound Test Service

- Increasing stress levels 1 (none) → 9 (high)  *(non-linear)*
- Parameters:
  - Operand array size and number of calculation loops
    (0, 20, 100, 1,000, 10,000, 25,000, 100,000)
    → Operands stored in random array locations
    → Induces page faults when seeking random locations
  - Number of function calls per loop
    (0, 20, 1,000, 100,000, 300,000)
- Control CPU time of function as input parameter

- **Goal: observe impact of CPU time on infrastructure scaling, provisioning variation, retention, and service performance**

April 20, 2018     *Serverless Computing: An Investigation of Factors Influencing Microservice Performance*     17

---

# AWS Lambda Testing

REST/JSON

Images credit: aws.amazon.com



API GATEWAY

Client:
c4.2xlarge

BASH: GNU Parallel
Multi-thread client
**"partest"**

Up to 100 concurrent
synchronous requests

Results of each thread
traced individually

**Fixed-availability zone:
EC2 client / Lambda server
us-east-1e**

CPU-bound
Test Function

Max
service duration:
< 30 seconds

Memory:
128 to 1536MB

April 20, 2018     *Serverless Computing: An Investigation of Factors Influencing Microservice Performance*     18

---

# AWS Lambda Testing

REST/JSON

Images credit: aws.amazon.com

API GATEWAY

Client:
c4.2xlarge

CPU-bound
Test Function

**Automatic Metrics Collection:**

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

April 20, 2018        *Serverless Computing: An Investigation of Factors Influencing Microservice Performance*        **19**

---

# Azure Functions Testing

- Http-triggered function app, written in C#
- Logs to Azure Table storage (*similar to Dynamo DB*)
  - Unique app service instance IDs
  - Current worker process ID
- Consumption plan → auto-scaled infrastructure
  - vs. app service plan (*deployment to dedicated VMs*)
- Performance testing:
  Visual Studio Team System (VSTS)

April 20, 2018        *Serverless Computing: An Investigation of Factors Influencing Microservice Performance*        **20**

# Outline

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

April 20, 2018    *Serverless Computing: An Investigation of Factors Influencing Microservice Performance*    21

# CPU-Bound Lambda Test Service WARM Performance



April 20, 2018    *Serverless Computing: An Investigation of Factors Influencing Microservice Performance*    22
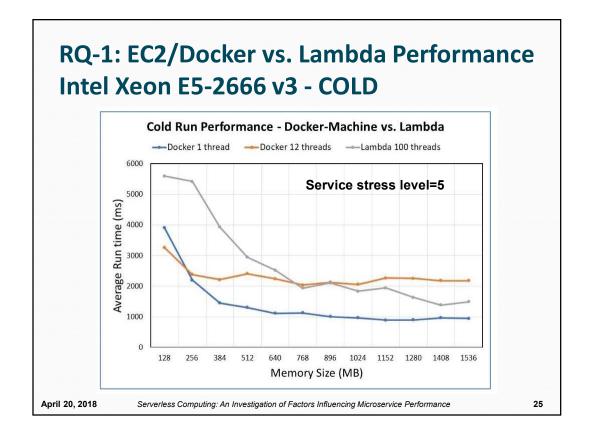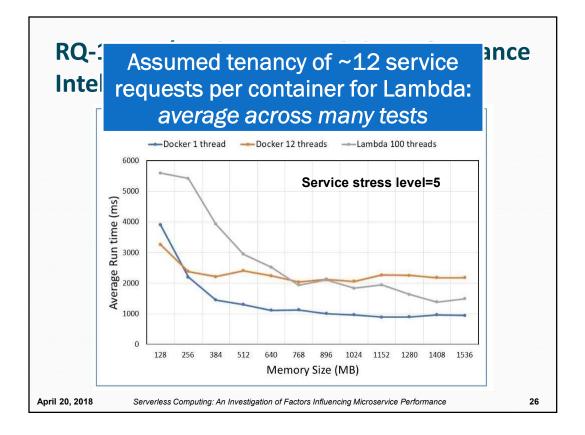
## RQ-1: Elasticity

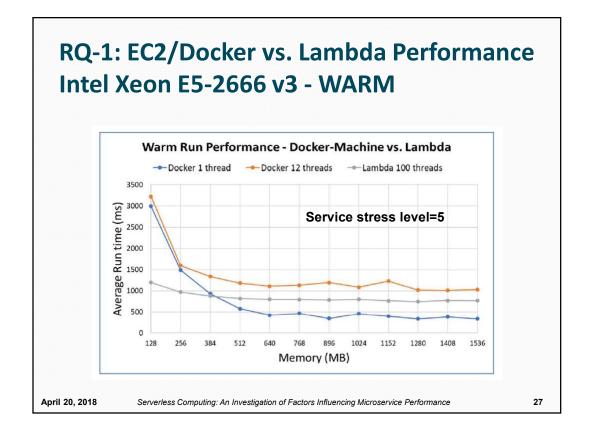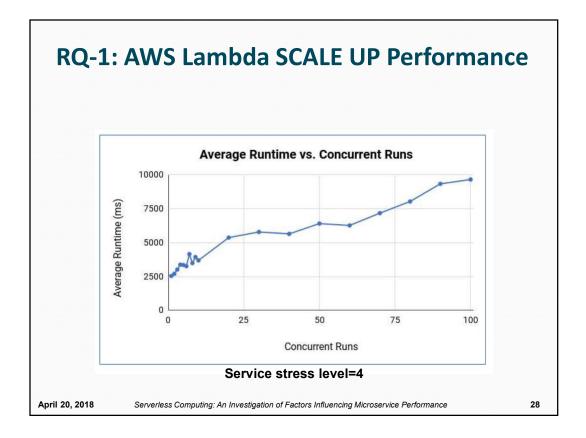What are the performance implications of infrastructure **_elasticity_** for serverless computing?

*(e.g. COLD vs. WARM performance)*

23

## RQ-1: AWS Lambda Latency Evaluation

- **AWS Lambda Simulation**
- Harness c4.8xlarge 36 vCPU VM instance
  - Intel Xeon E5-2666v3 CPU – *same as Lambda*
- Lambda JAR file deployed Docker container(s)
  - **Set memory**: `docker run "-m <ram in MB>"`
  - **Set CPUs**: `docker run "—cpus <VCPUs>`
- Compare: 1 and 12 concurrent runs
  - Avg VM tenancy ~12.3 of all tests

- **How does Lambda scale CPU power?**

**Literal Estimates:**

| Memory (MB) | Expected CPU% |
|---|---|
| 128 | 16.6% |
| 256 | 33.3% |
| 384 | 50.0% |
| 512 | 66.6% |
| 640 | 83.3% |
| 768 | 100.0% |
| 896 | 116.7% |
| 1024 | 133.3% |
| 1152 | 150.0% |
| 1280 | 166.60% |
| 1408 | 183.30% |
| 1536 | 200.00% |

April 20, 2018          *Serverless Computing: An Investigation of Factors Influencing Microservice Performance*          24

## RQ-1: EC2/Docker vs. Lambda Performance Intel Xeon E5-2666 v3 - COLD



**Cold Run Performance - Docker-Machine vs. Lambda**

Service stress level=5

## RQ-1: EC2/Docker vs. Lambda Performance Intel...

**Assumed tenancy of ~12 service requests per container for Lambda:** *average across many tests*



Service stress level=5

## RQ-1: EC2/Docker vs. Lambda Performance Intel Xeon E5-2666 v3 - WARM



Warm Run Performance - Docker-Machine vs. Lambda

Service stress level=5

## RQ-1: AWS Lambda SCALE UP Performance



Average Runtime vs. Concurrent Runs

Service stress level=4

## RQ-1: Azure functions COLD Performance
*includes "container" initialization*



**Up to 4 VMs automatically created**

April 20, 2018      *Serverless Computing: An Investigation of Factors Influencing Microservice Performance*      29

## RQ-1 ... ce
*inclu...*

VMs are allocated as opposed
to individual container instances.
Supports better initial performance.



**Up to 4 VMs automatically created**

April 20, 2018      *Serverless Computing: An Investigation of Factors Influencing Microservice Performance*      30

# RQ-2: Load Balancing

How does **_load balancing_** vary in serverless computing?

How do computational requests impact load balancing, and ultimately performance?

31

# RQ-2: COLD Lambda Infrastructure for Scaling

**Infrastructure Elasticity - Increasing Concurrent Requests**



■ containers ■ hosts ■ runs_per_host

concurrent runs

**Service stress level=4**

## RQ-2: COLD Lambda Infrastructure for Scaling

COLD service requests receive separate container instances to amortize startup overhead



**Service stress level=4**

# RQ-2: WARM Lambda Infrastructure for Scaling



**Average for 100 runs**

# RQ-2: WARM Lambda Infr...

**WARM service requests share container instances unless CPU requirements are increased**

## Infrastructure Elasticity - Calculation Stress Levels

containers ■   hosts ●   runs_per_container ●   runs_per_host ▲

**Average for 100 runs**

# RQ-2: COLD Azure Functions Infrastructure for Scaling

## Infrastructure Elasticity - Azure Functions Load Test

Test Duration: ■ 2 min   ■ 5 min   ■ 10 min

App Service Instances Used

Number of Concurrent Service Requests

## RQ-3: Provisioning Variation

What performance implications result from ***provisioning variation*** of container infrastructure?

37

---

## RQ-3: Cold Lambda service performance vs. Container Placement

**Service stress level=4**

When more containers were placed on the same VMs for COLD service requests, Lambda Performance suffered up to 5x !

The impact of tenancy vs. performance is quite clear.



$R^2 = 0.9885564979$

Average COLD service execution time (ms)
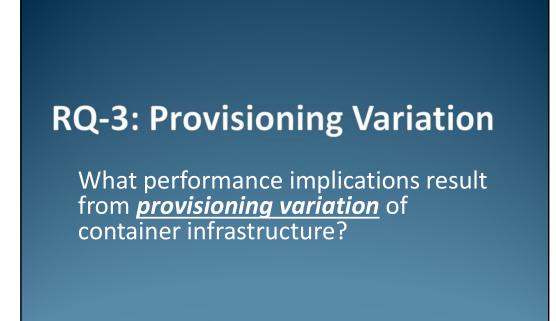
Containers per host (VM)

# RQ-5: Memory Reservation
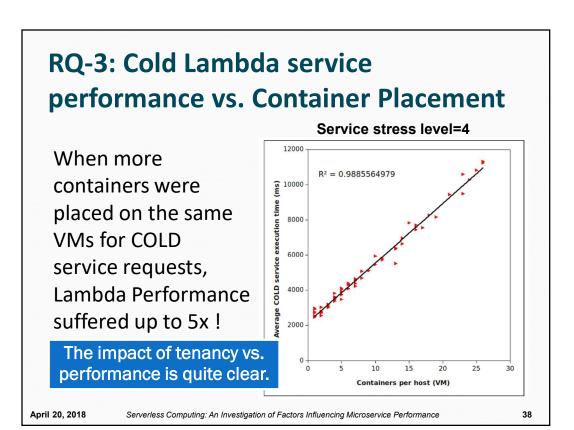
What performance implications result from microservice ***memory reservation*** size?

How does memory reservation size impact container placement?

41

# RQ-5: Slider Bar Test: Memory vs. CPU power

**Service stress level=4**

**Memory Size vs. Average Service Performance**



April 20, 2018    *Serverless Computing: An Investigation of Factors Influencing Microservice Performance*    **42**

## RQ-5: Slider Bar Test II:
## Infrastructure vs. Memory Reservation

**Service stress level=4**

## RQ-5: Slider Bar Test II:
## Infrastructure vs. Memory Reservation

**Increasing the memory reservation size results in more hosting infrastructure**

# Outline

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

# Conclusions

- **RQ-1 Elasticity**: Extra infrastructure is provisioned to compensate for initialization overhead of "container" startup
  - VM COLD: up to ~20x slower than WARM
  - Container COLD: ~5x slower than WARM

- **RQ-2 Load Balancing**: Better when COLD. WARM runs only use all original infrastructure when CPU-bound execution time is similar to container initialization execution time
  - Must increase stress level to harness available infrastructure

# Conclusions - 2

- **RQ-3 Provisioning Variation**: Bad placement can lead to ~4.6x degradation in COLD service performance

- **RQ-4 Infrastructure Retention**:
  3 distinct performance states:
  *VM COLD*, *Container COLD*, *WARM*
  - Containers begin to disappear after 10 minutes
  - VM hosts deprecated after ~40 minutes

- **RQ-5 Memory Reservation**:
- For non memory-bound service, performance improves up to ~512-640MB

April 20, 2018        *Serverless Computing: An Investigation of Factors Influencing Microservice Performance*        **47**

## Questions