# Tracking Wasted Money

# Tracking Wasted Money

**Job Characteristics on Large-Scale Systems: Long-Term Analysis, Quantification, and Implications***

Tirthak Patel
Northeastern University

Zhengchun Liu, Raj Kettimuthu
Argonne National Laboratory

Paul Rich, William Allcock

Devesh Tiwari

**A Case For Intra-rack Resource Disaggregation in HPC**

GEORGE MICHELOGIANNAKIS, Lawrence Berkeley National Laboratory, USA
BENJAMIN KLENK, NVIDIA, USA
BRANDON COOK, Lawrence Berkeley National Laboratory, USA
MIN YEE TEH and MADELEINE GLICK, Columbia University, USA
LARRY DENNISON, NVIDIA, USA
KEREN BERGMAN, Columbia University, USA
JOHN SHALF, Lawrence Berkeley National Laboratory, USA

TACO, 2022

**FINAL REPORT**

**Quantifying Memory Underutilization in HPC Systems and Using it to Improve Performance via Architecture Support**

Gagandeep Panwar*
Virginia Tech
Blacksburg, USA
gpanwar@vt.edu

Da Zhang*
Virginia Tech
Blacksburg, USA
daz3@vt.edu

Yihan Pang*
Virginia Tech
Blacksburg, USA
pyihan1@vt.edu

Mai Dahshan
Virginia Tech
Blacksburg, USA
mdahshan@vt.edu

Nathan DeBardeleben
Los Alamos National Laboratory
Los Alamos, USA
ndebard@lanl.gov

Binoy Ravindran
Virginia Tech
Blacksburg, USA
binoy@vt.edu

Xun Jian
Virginia Tech
Blacksburg, USA
xunj@vt.edu

v, Jeffrey

for

Enos, and
cations

iv, 2017

MICRO, 2019

**A Holistic View of Memory Utilization on HPC Systems: Current and Future Trends**

Ivy B. Peng*
peng8@llnl.gov
Lawrence Livermore National
Laboratory
USA

Ian Karlin
karlin1@llnl.gov
Lawrence Livermore National
Laboratory
USA

Maya B. Gokhale
gokhale2@llnl.gov
Lawrence Livermore National
Laboratory
USA

Kathleen Shoga
Shoga1@llnl.gov
Lawrence Livermore National
Laboratory
USA

Matthew Legendre
legendre1@llnl.gov
Lawrence Livermore National
Laboratory
USA

Todd Gamblin
gamblin2@llnl.gov
Lawrence Livermore National
Laboratory
USA

and

MEMSYS, 2021

University of Tennessee, Knoxville, TN 37996, USA
{hyou,haozhang}@utk.edu

JSSPP, 2012

2

# Tracking Wasted Money

**Job Characteristics on Large-Scale Systems:**
**Long-Term Analysis, Quantification, and Implications***

Tirthak Patel          Zhengchun Liu, Raj Kettimuthu
Northeastern University     Argonne National Laboratory

Paul Rich, William Allcock          Devesh Tiwari

**A Case For Intra-rack Resource Disaggregation in HPC**

GEORGE MICHELOGIANNAKIS, Lawrence Berkeley National Laboratory, USA
BENJAMIN KLENK, NVIDIA, USA
BRANDON COOK, Lawrence Berkeley National Laboratory, USA
MIN YEE TEH and MADELEINE GLICK, Columbia University, USA
LARRY DENNISON, NVIDIA, USA
KEREN BERGMAN, Columbia University, USA
JOHN SHALF, Lawrence Berkeley National Laboratory, USA

TACO, 2022

FINAL REPORT

**Quantifying Memory Underutilization in HPC Systems and**
**Using it to Improve Performance via Architecture Support**

Gagandeep Panwar*          Da Zhang*          Yihan Pang*
Virginia Tech          Virginia Tech          Virginia Tech
Blacksburg, USA          Blacksburg, USA          Blacksburg, USA
gpanwar@vt.edu          daz3@vt.edu          pyihan1@vt.edu

Mai Dahshan          Nathan DeBardeleben          Binoy Ravindran
Virginia Tech          Los Alamos National Laboratory          Virginia Tech
Blacksburg, USA          Los Alamos, USA          Blacksburg, USA
mdahshan@vt.edu          ndebard@lanl.gov          binoy@vt.edu

Xun Jian
Virginia Tech
Blacksburg, USA
xunj@vt.edu

v, Jeffrey
for

Enos, and
cations

iv, 2017

MICRO, 2019

**A Holistic View of Memory Utilization on HPC Systems:**
**Current and Future Trends**

Ivy B. Peng*          Ian Karlin          Maya B. Gokhale
peng8@llnl.gov          karlin1@llnl.gov          gokhale2@llnl.gov
Lawrence Livermore National          Lawrence Livermore National          Lawrence Livermore National
Laboratory          Laboratory          Laboratory
USA          USA          USA

Kathleen Shoga          Matthew Legendre          Todd Gamblin
Shoga1@llnl.gov          legendre1@llnl.gov          gamblin2@llnl.gov
Lawrence Livermore National          Lawrence Livermore National          Lawrence Livermore National
Laboratory          Laboratory          Laboratory
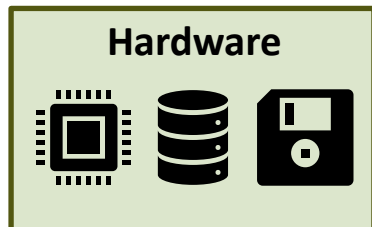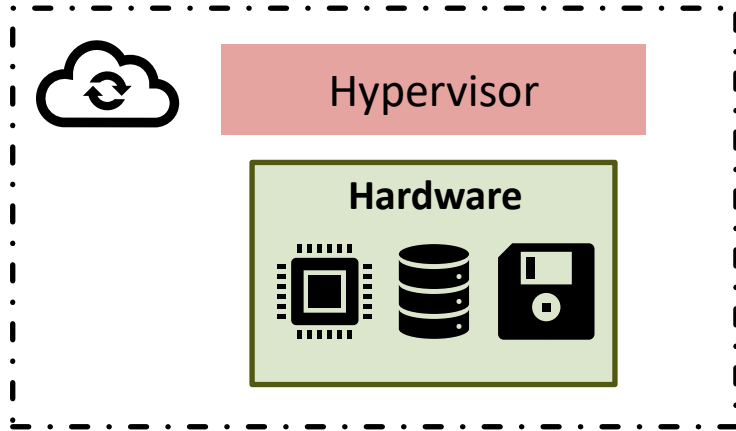USA          USA          USA

and

MEMSYS, 2021

University of Tennessee, Knoxville, TN 37996, USA
{hyou,haozhang}@utk.edu

JSSPP, 2012

**Can we solve underutilization with sharing and fine-grained allocations?**
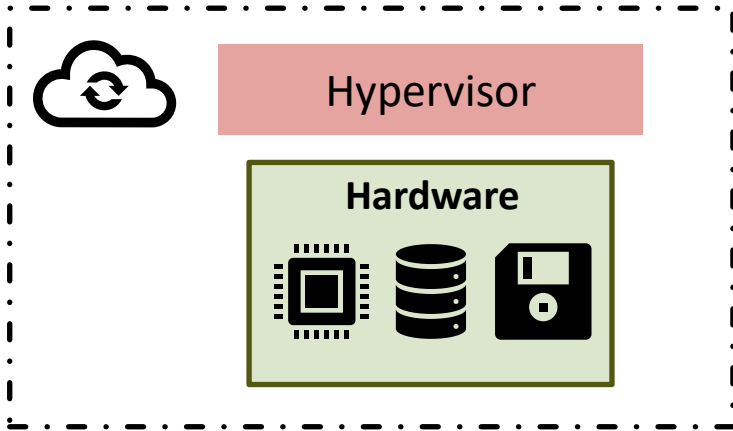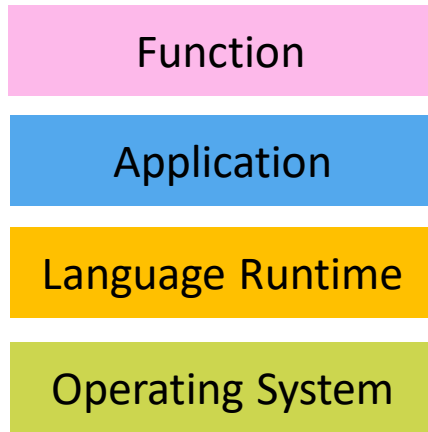
# Cloud and Serverless

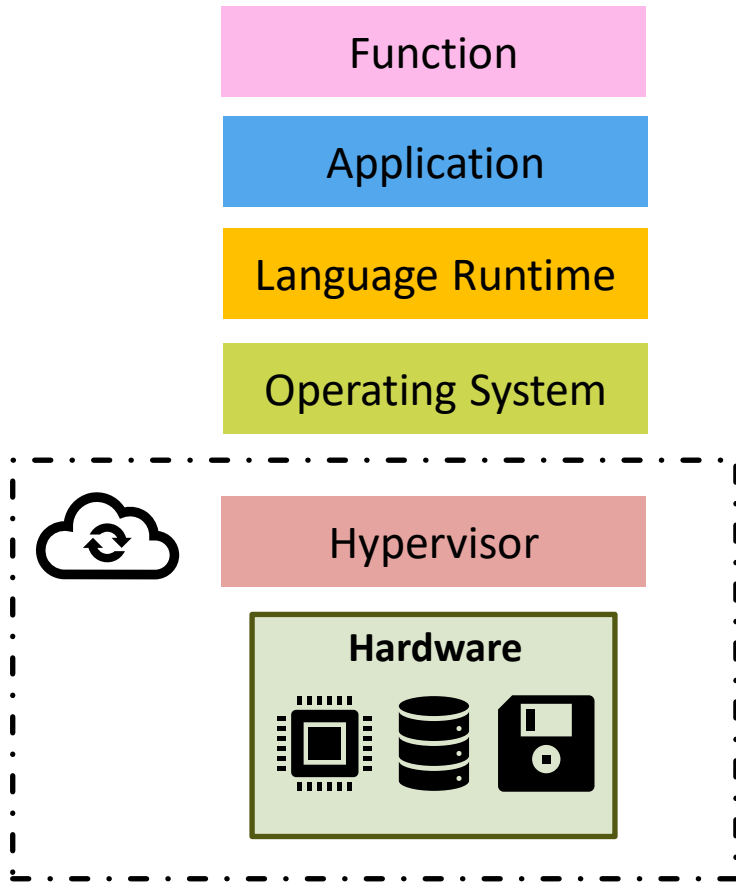# Cloud and Serverless

**Hardware**

# Cloud and Serverless

Hypervisor

**Hardware**

# Cloud and Serverless



Function

Application

Language Runtime

Operating System

Hypervisor

**Hardware**

**Virtual Machine**

# Cloud and Serverless



**Virtual Machine**
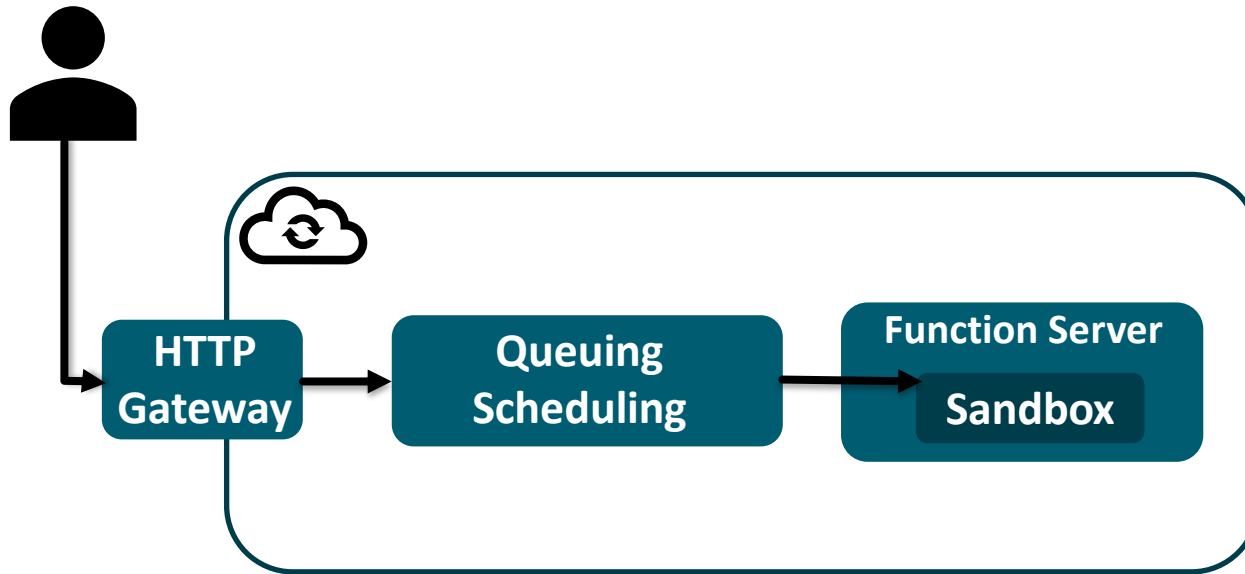
**Containers**

# Cloud and Serverless



**Virtual Machine**
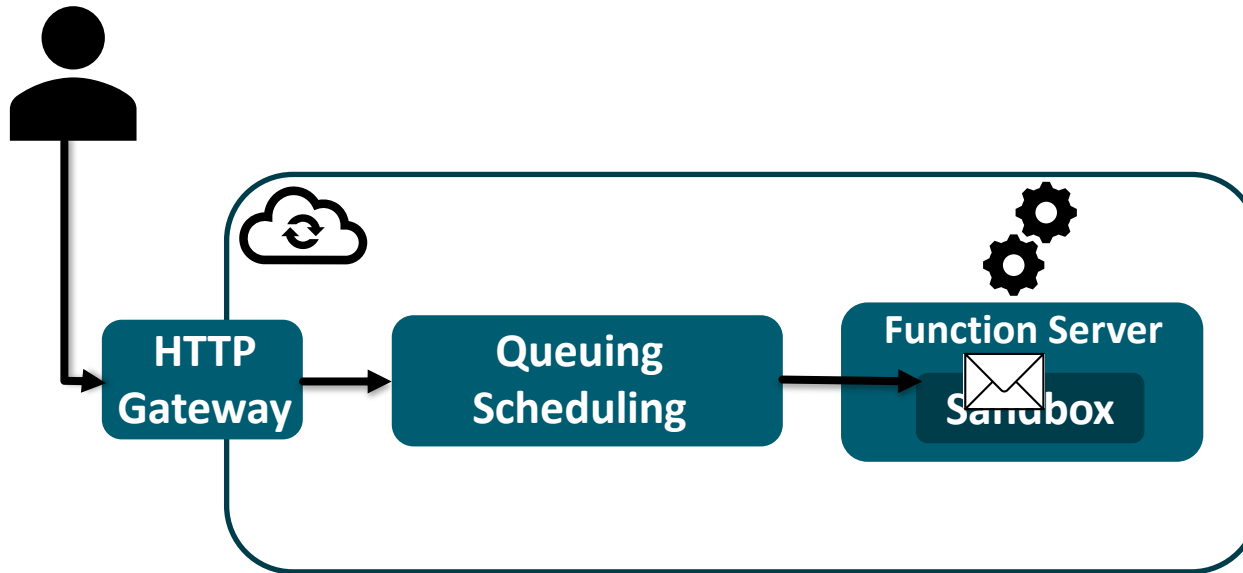
**Containers**

**Functions**

# How does Function-as-a-Service (FaaS) work?

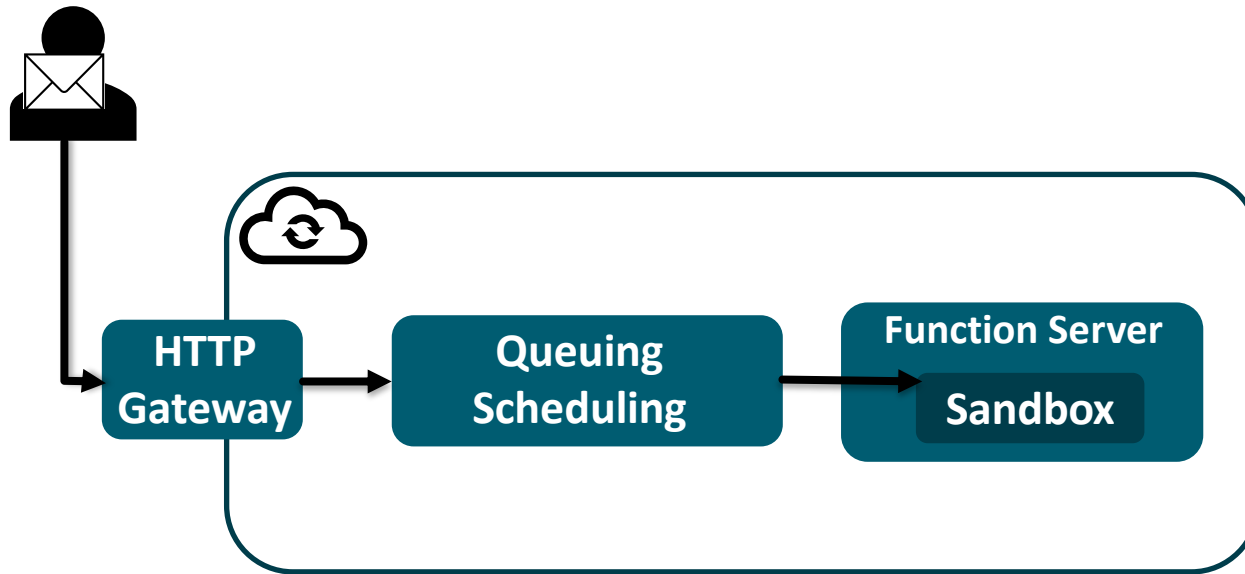# How does Function-as-a-Service (FaaS) work?

# How does Function-as-a-Service (FaaS) work?

# How does Function-as-a-Service (FaaS) work?

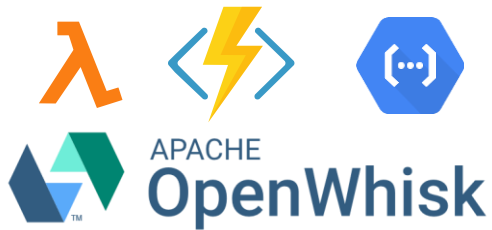# How does Function-as-a-Service (FaaS) work?

# SeBS: The Serverless Benchmark Suite

"SeBS: a Serverless Benchmark Suite for Function-as-a-Service Computing", ACM/IFIP Middleware 2021

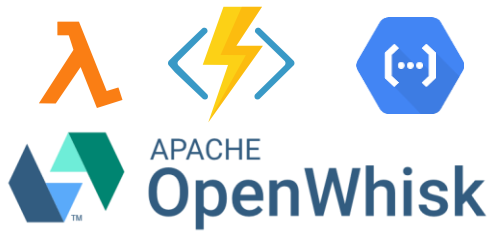# SeBS: The Serverless Benchmark Suite

**Serverless Platforms**

**"SeBS: a Serverless Benchmark Suite for Function-as-a-Service Computing", ACM/IFIP Middleware 2021**

# SeBS: The Serverless Benchmark Suite



**Serverless Platforms**

**Benchmarks**

"SeBS: a Serverless Benchmark Suite for Function-as-a-Service Computing", ACM/IFIP Middleware 2021

# SeBS: The Serverless Benchmark Suite
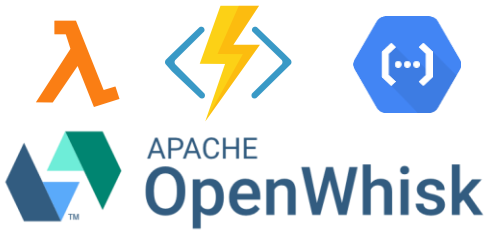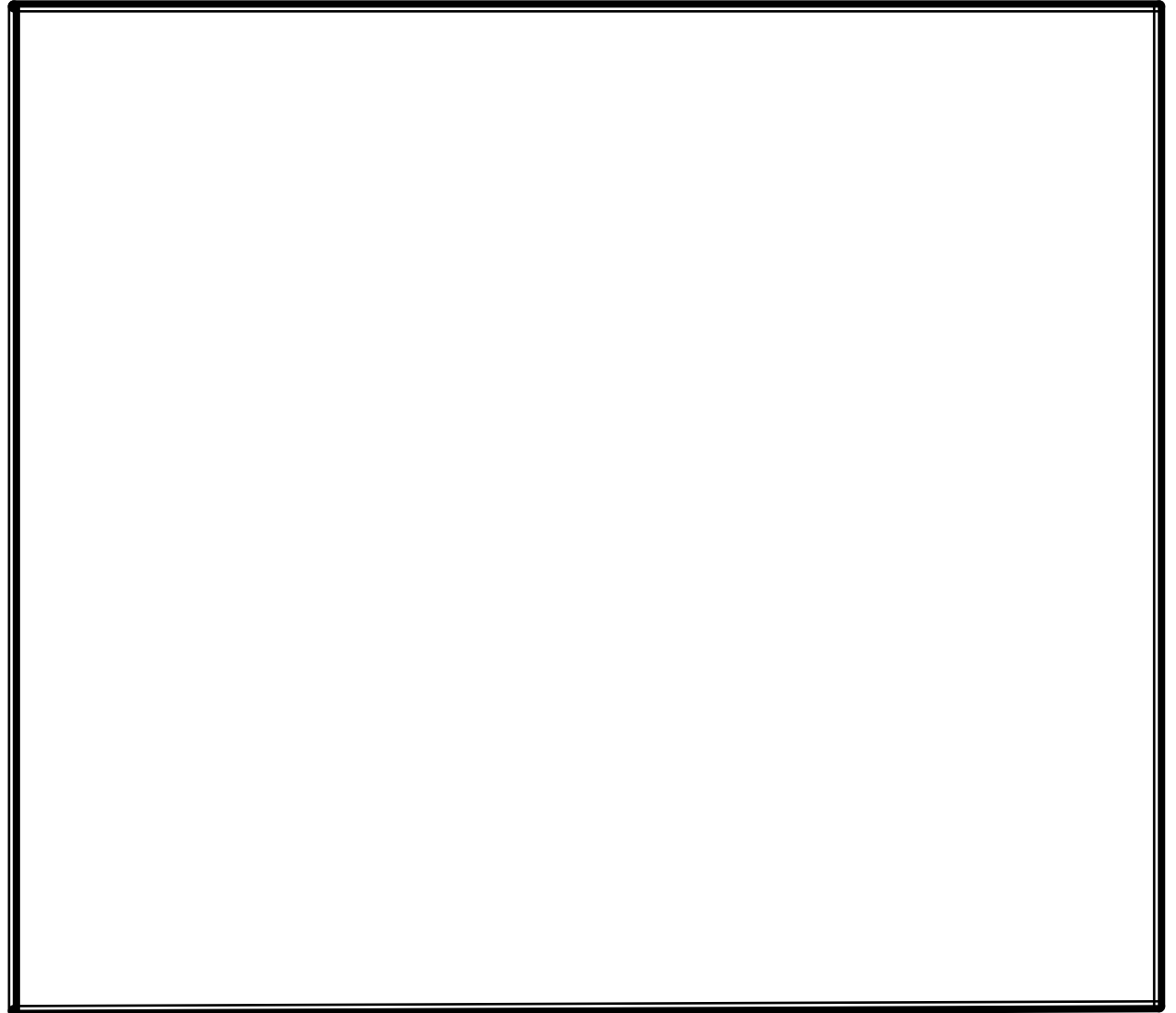
**Serverless Platforms**

**Benchmarks**

**Experiments**

Performance & Cost
Invocation Overhead
Container Eviction
Serverless Communication
Serverless Workflows

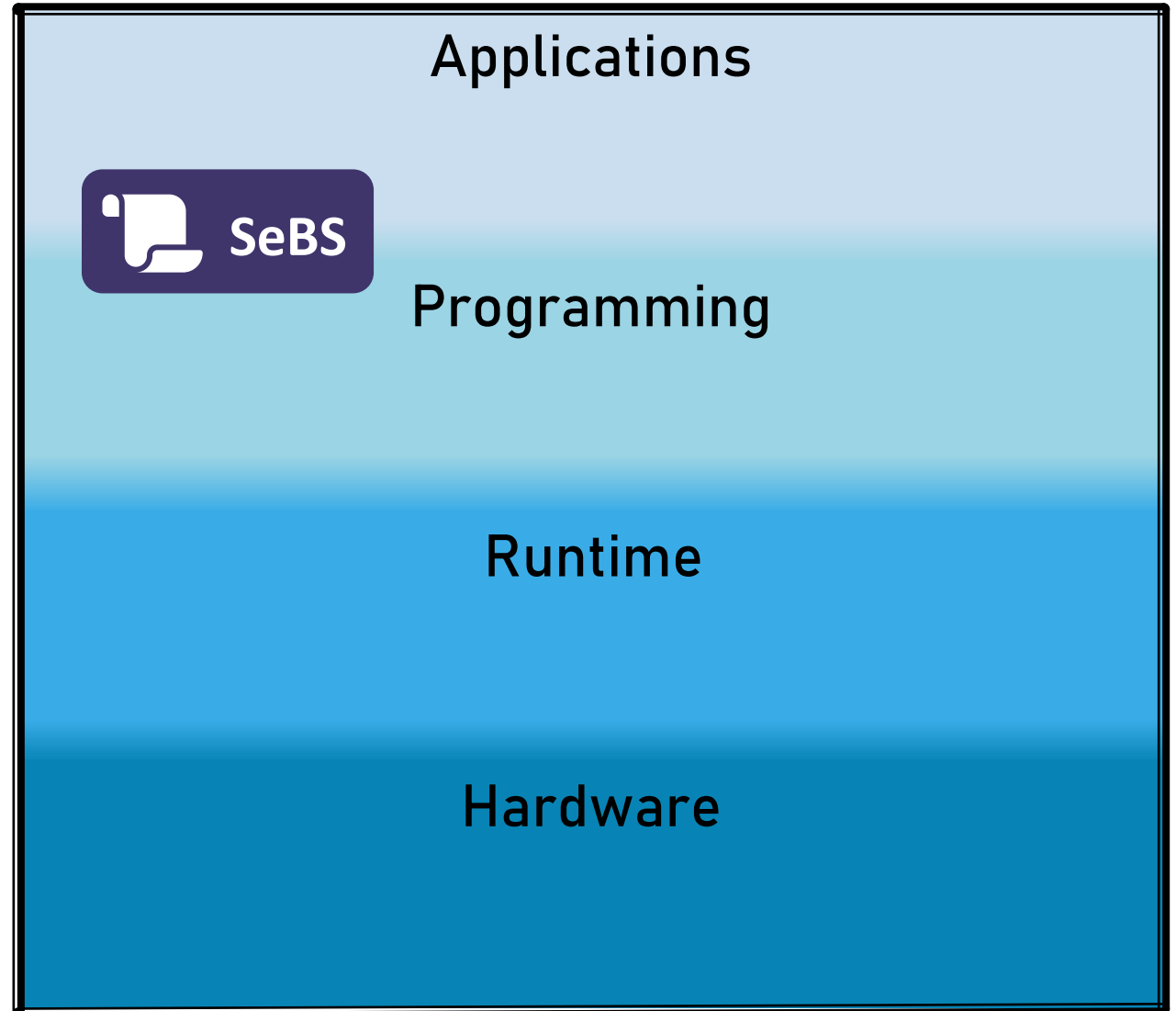**"SeBS: a Serverless Benchmark Suite for Function-as-a-Service Computing", ACM/IFIP Middleware 2021**

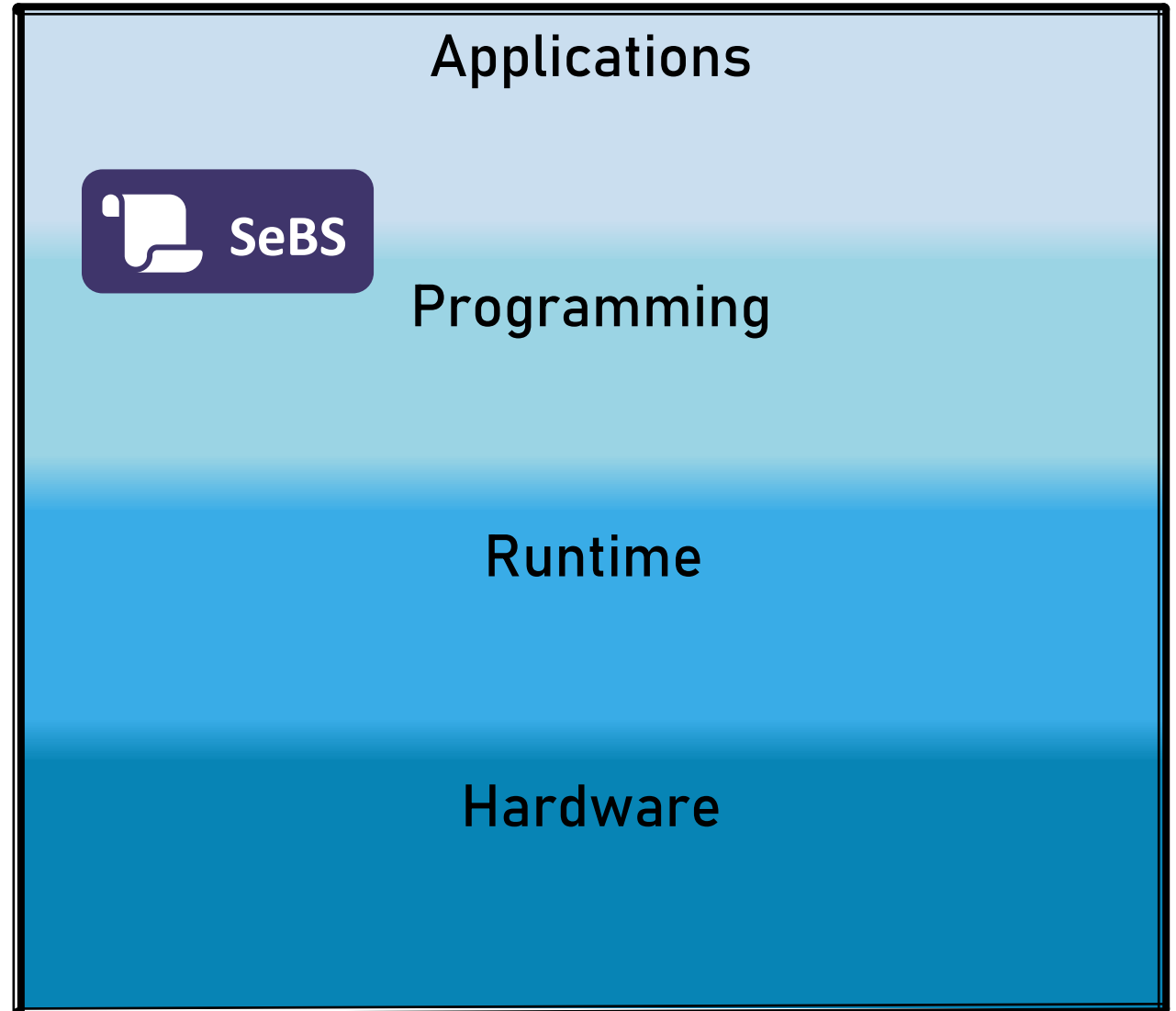# Serverless for High-Performance Applications

# Serverless for High-Performance Applications
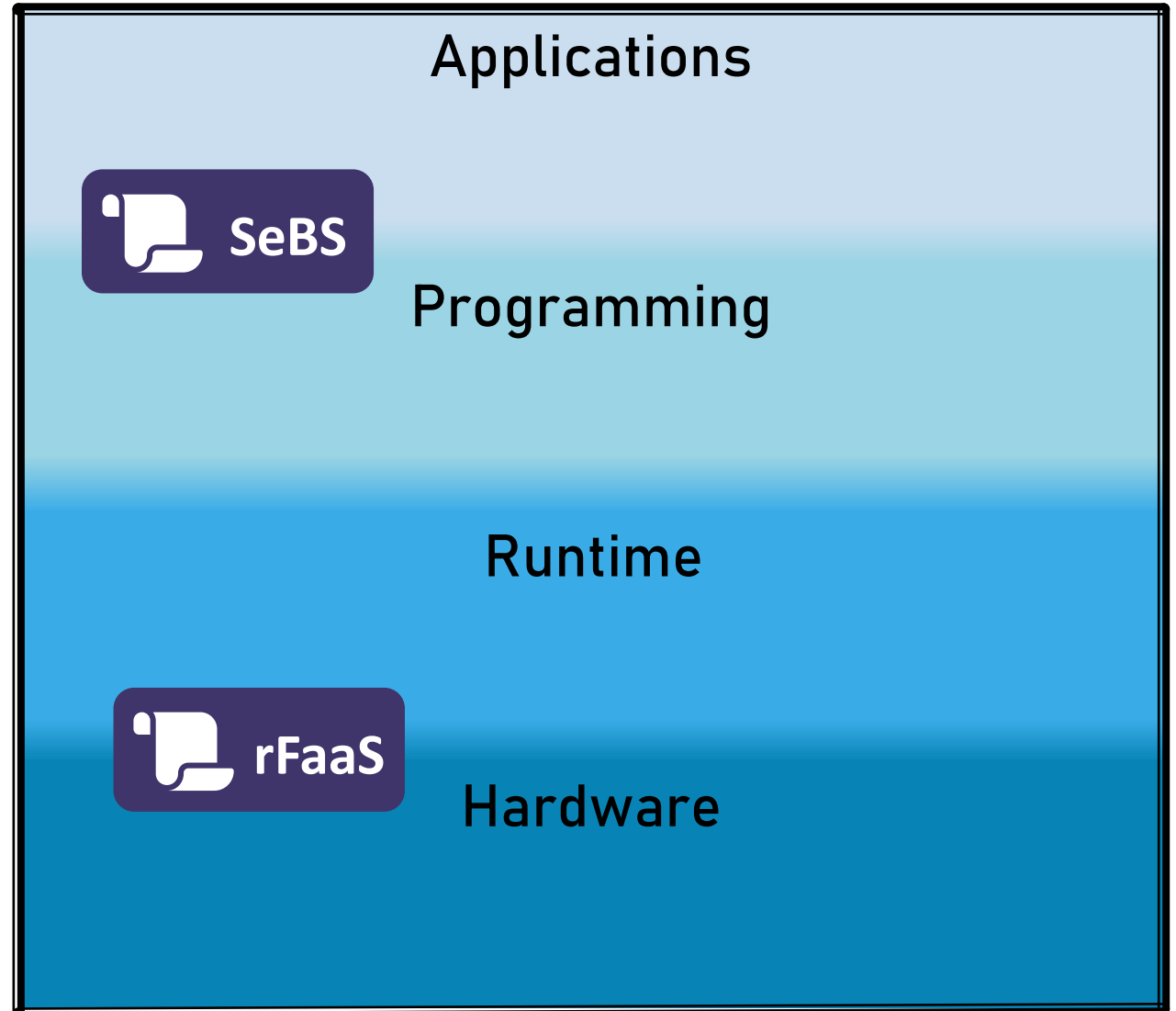
# Serverless for High-Performance Applications



Applications

**SeBS**

Programming

Runtime

Hardware

# Serverless for High-Performance Applications

Functions are expensive to invoke.
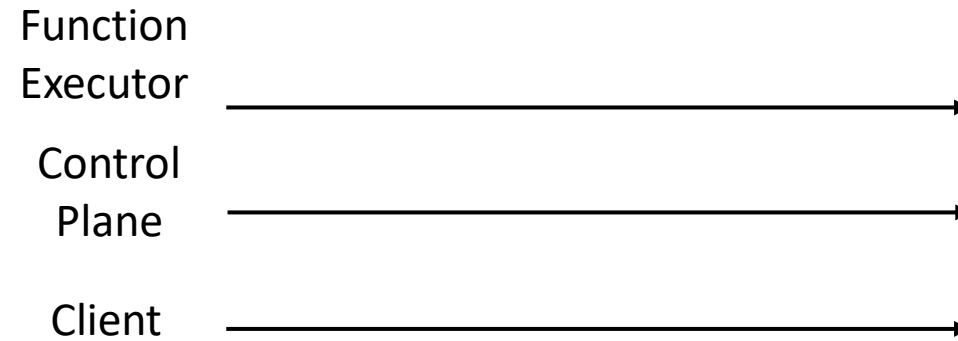
SeBS

Applications

Programming

Runtime

Hardware

# Serverless for High-Performance Applications

Functions are expensive to invoke.

Applications

SeBS

Programming

Runtime

rFaaS

Hardware

# Invocations in FaaS and rFaaS

Function
Executor ————————————————————————→

Control
Plane ————————————————————————→

Client ————————————————————————→

**"rFaaS: Enabling High Performance Serverless with RDMA and Leases", IPDPS'23**

# Invocations in FaaS and rFaaS

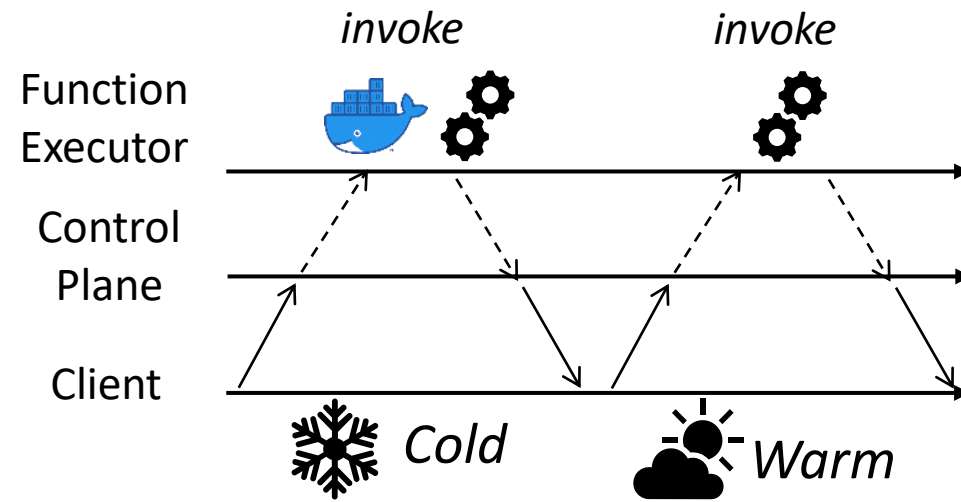**"rFaaS: Enabling High Performance Serverless with RDMA and Leases", IPDPS'23**

# Invocations in FaaS and rFaaS

# Invocations in FaaS and rFaaS



Function Executor

Control Plane

Client

*invoke*          *invoke*

❄ *Cold*     ⛅ *Warm*

Executor

Control Plane

Client

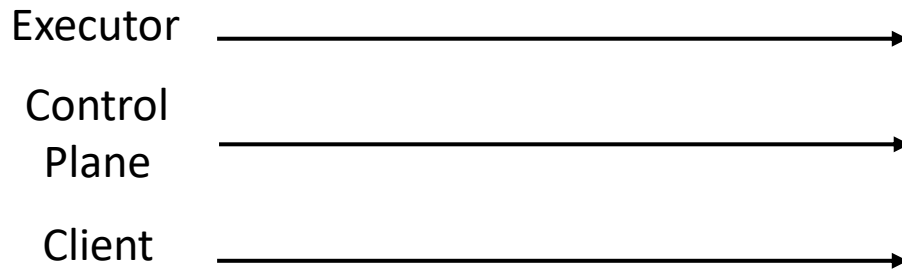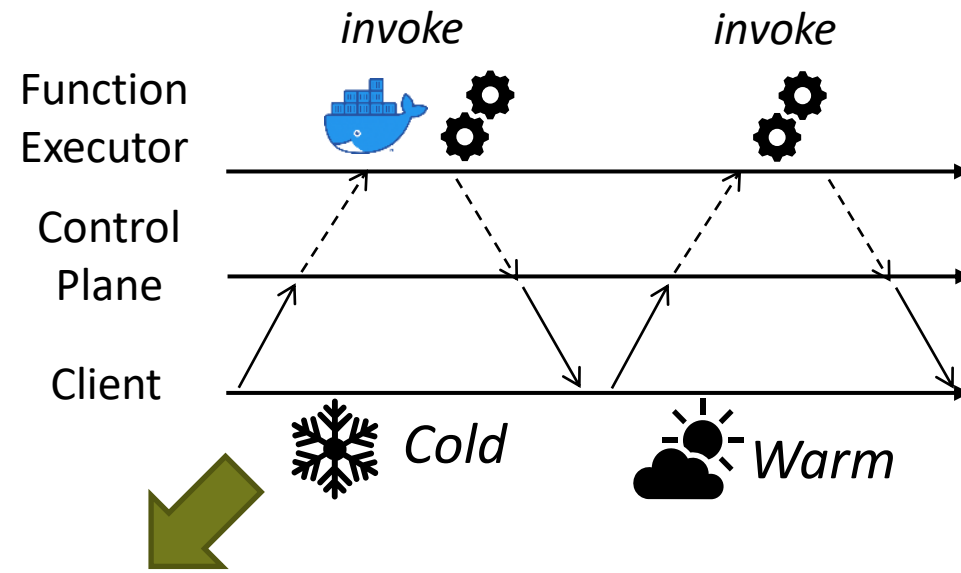**"rFaaS: Enabling High Performance Serverless with RDMA and Leases", IPDPS'23**
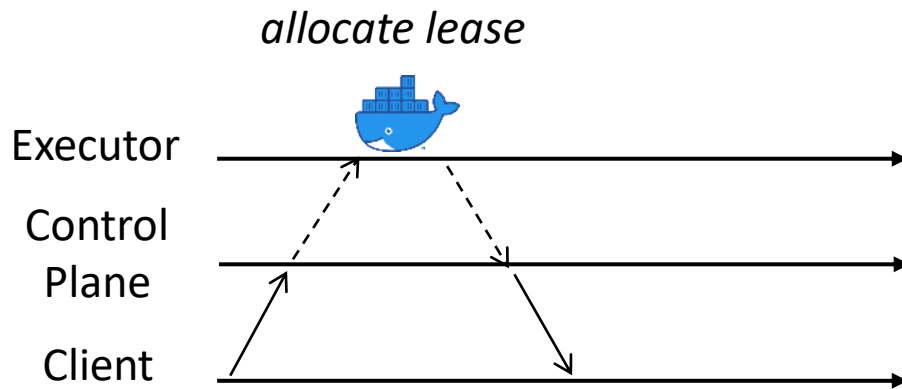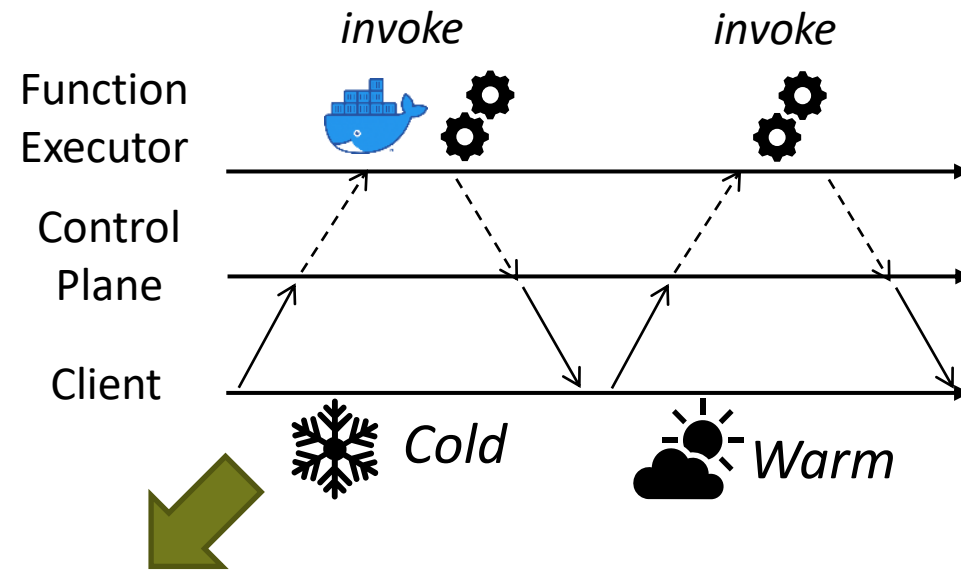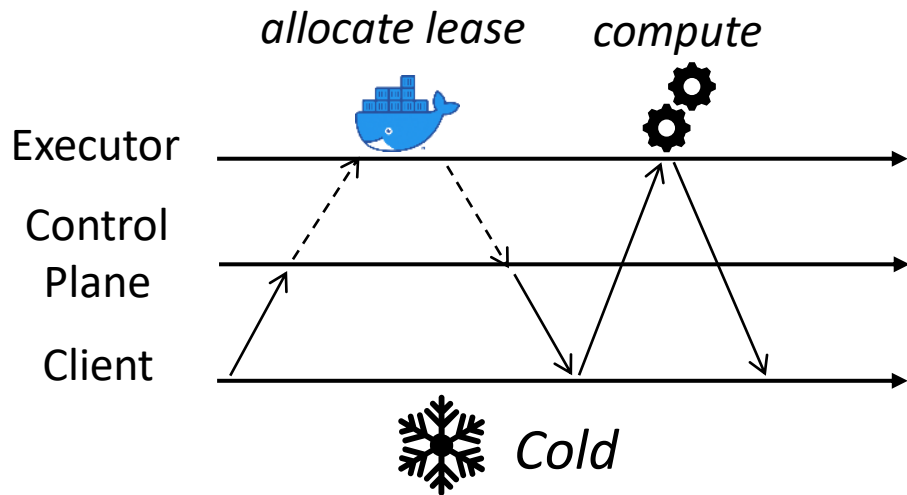
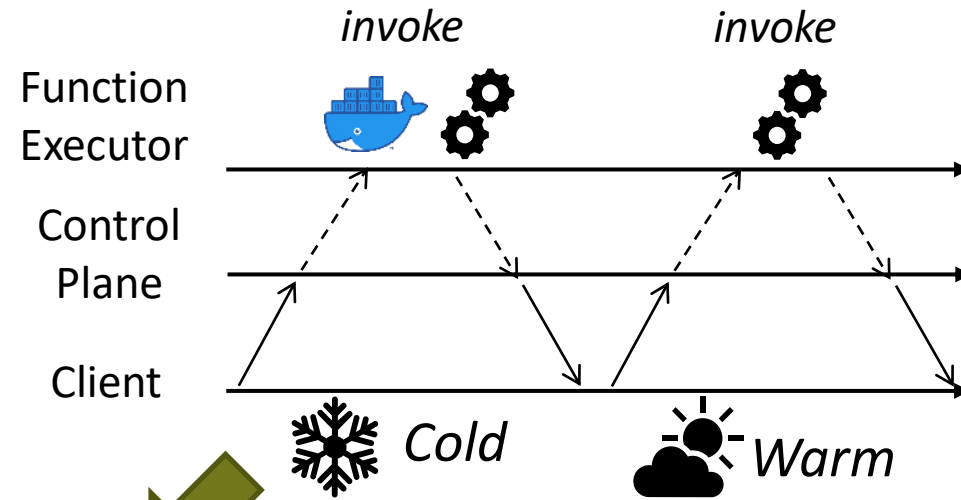# Invocations in FaaS and rFaaS

# Invocations in FaaS and rFaaS

"rFaaS: Enabling High Performance Serverless with RDMA and Leases", IPDPS'23

# Invocations in FaaS and rFaaS



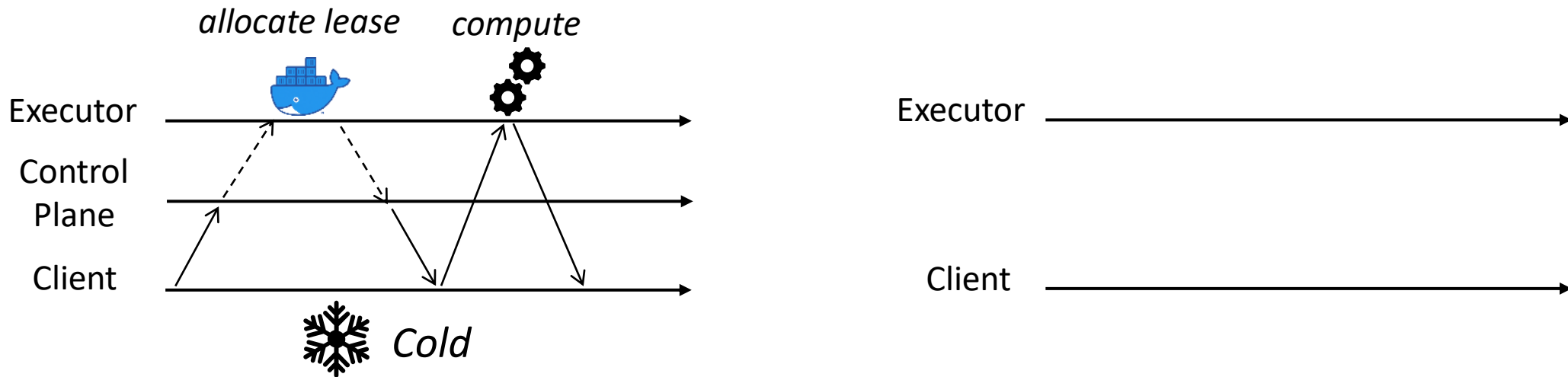"rFaaS: Enabling High Performance Serverless with RDMA and Leases", IPDPS'23

# Invocations in FaaS and rFaaS



"rFaaS: Enabling High Performance Serverless with RDMA and Leases", IPDPS'23

# Invocations in FaaS and rFaaS



"rFaaS: Enabling High Performance Serverless with RDMA and Leases", IPDPS'23

# How fast are invocations in FaaS?

"rFaaS: Enabling High Performance Serverless with RDMA and Leases", IPDPS'23

8

# How fast are invocations in FaaS?

"rFaaS: Enabling High Performance Serverless with RDMA and Leases", IPDPS'23

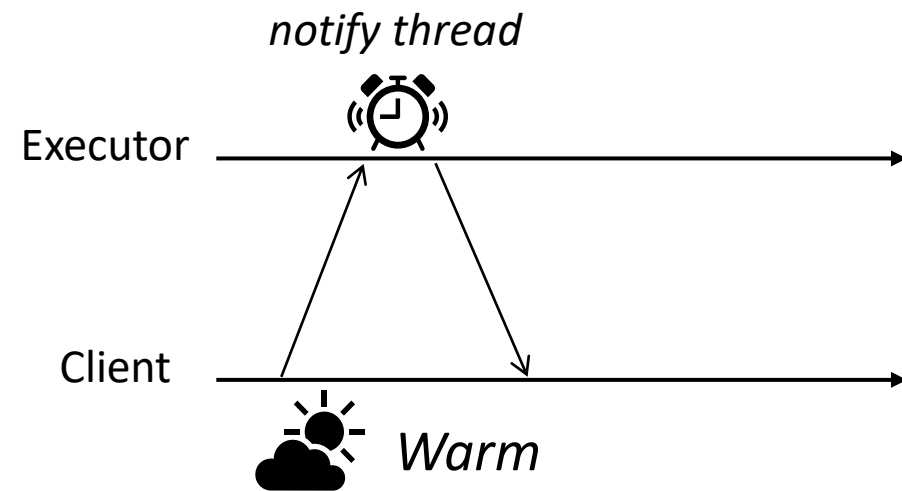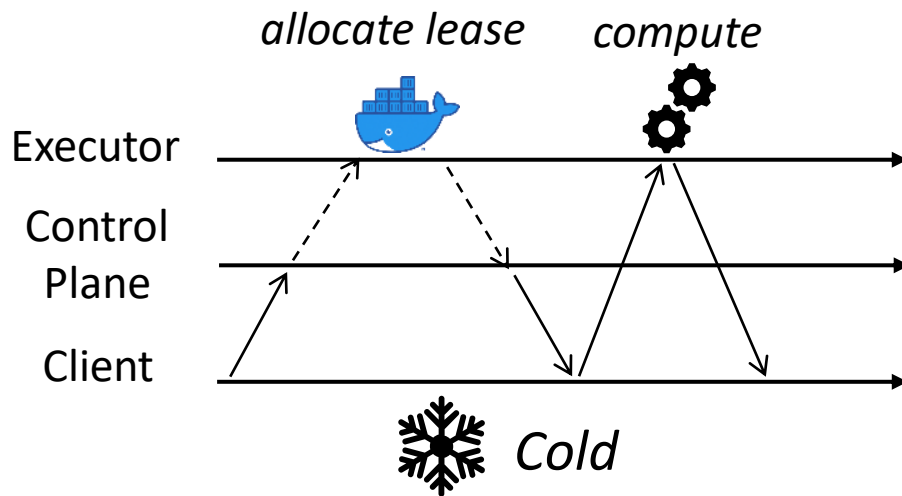# How fast are invocations in FaaS?

36 CPU cores, 377 GB memory.
100 Gbps Ethernet with RoCEv2 support.



OpenWhisk: 119.18 ms

OpenWhisk: 1.79 MB/s, HTTP

AWS: 19.64 ms
nightcore: 209.45 us

AWS: 17.21 MB/s, HTTP

nightcore: 453.72 MB/s, RPC

Round-Trip time [usec] — $10^1$, $10^2$, $10^3$, $10^4$, $10^5$, $10^6$

Message size [kB] — 1, 2, 4, 8, 16, 32, 64, 128, 256, 512, 1024, 2048, 5120

"rFaaS: Enabling High Performance Serverless with RDMA and Leases", IPDPS'23

8

# How fast are invocations in FaaS?

36 CPU cores, 377 GB memory.
100 Gbps Ethernet with RoCEv2 support.

Reduced invocation critical path

Zero-copy RDMA networking

OpenWhisk: 119.18 ms

AWS: 19.64 ms
nightcore: 209.45 us

OpenWhisk: 1.79 MB/s, HTTP

AWS: 17.21 MB/s, HTTP

nightcore: 453.72 MB/s, RPC

Round-Trip time [usec] vs Message size [kB]

"rFaaS: Enabling High Performance Serverless with RDMA and Leases", IPDPS'23
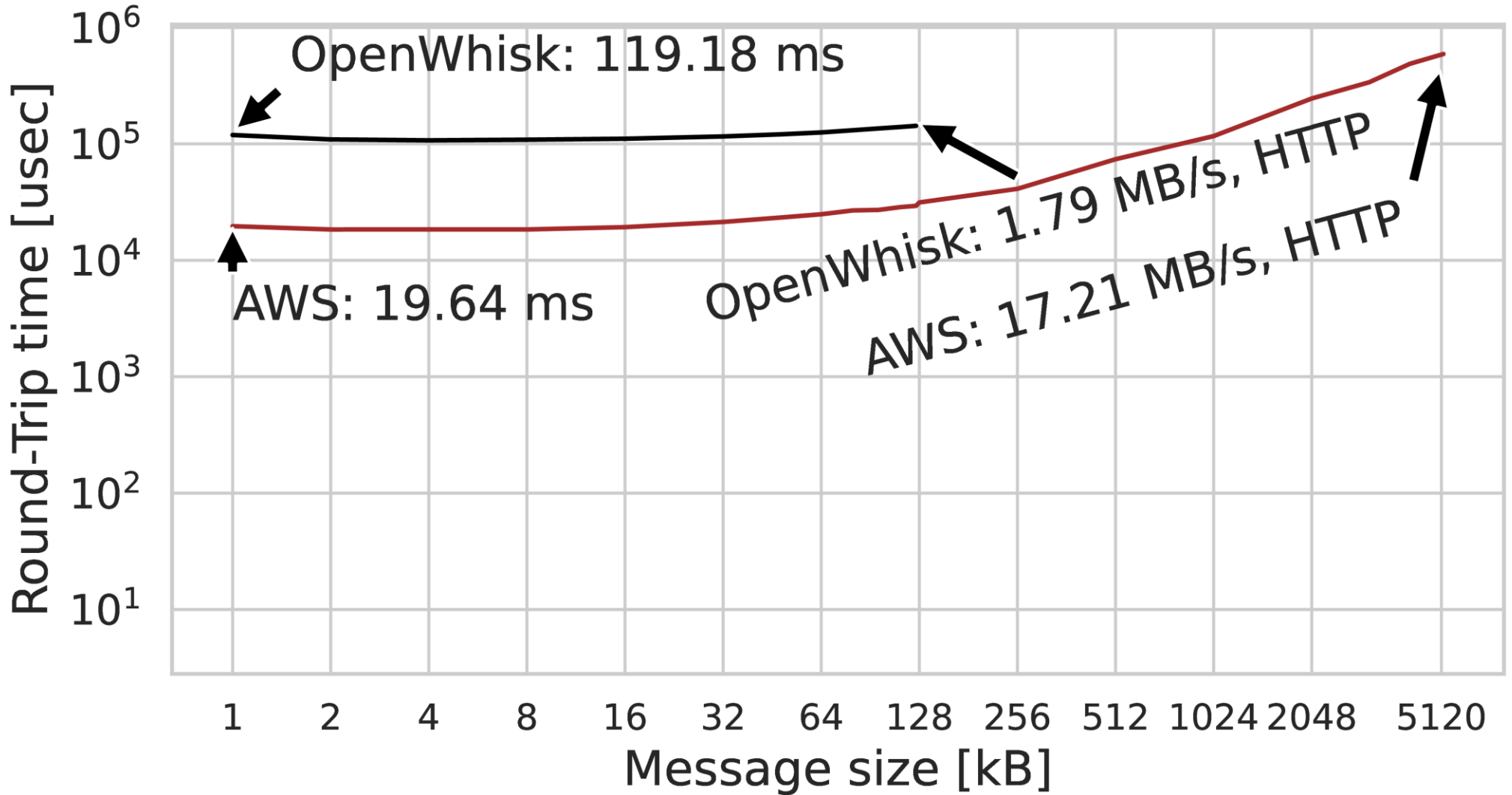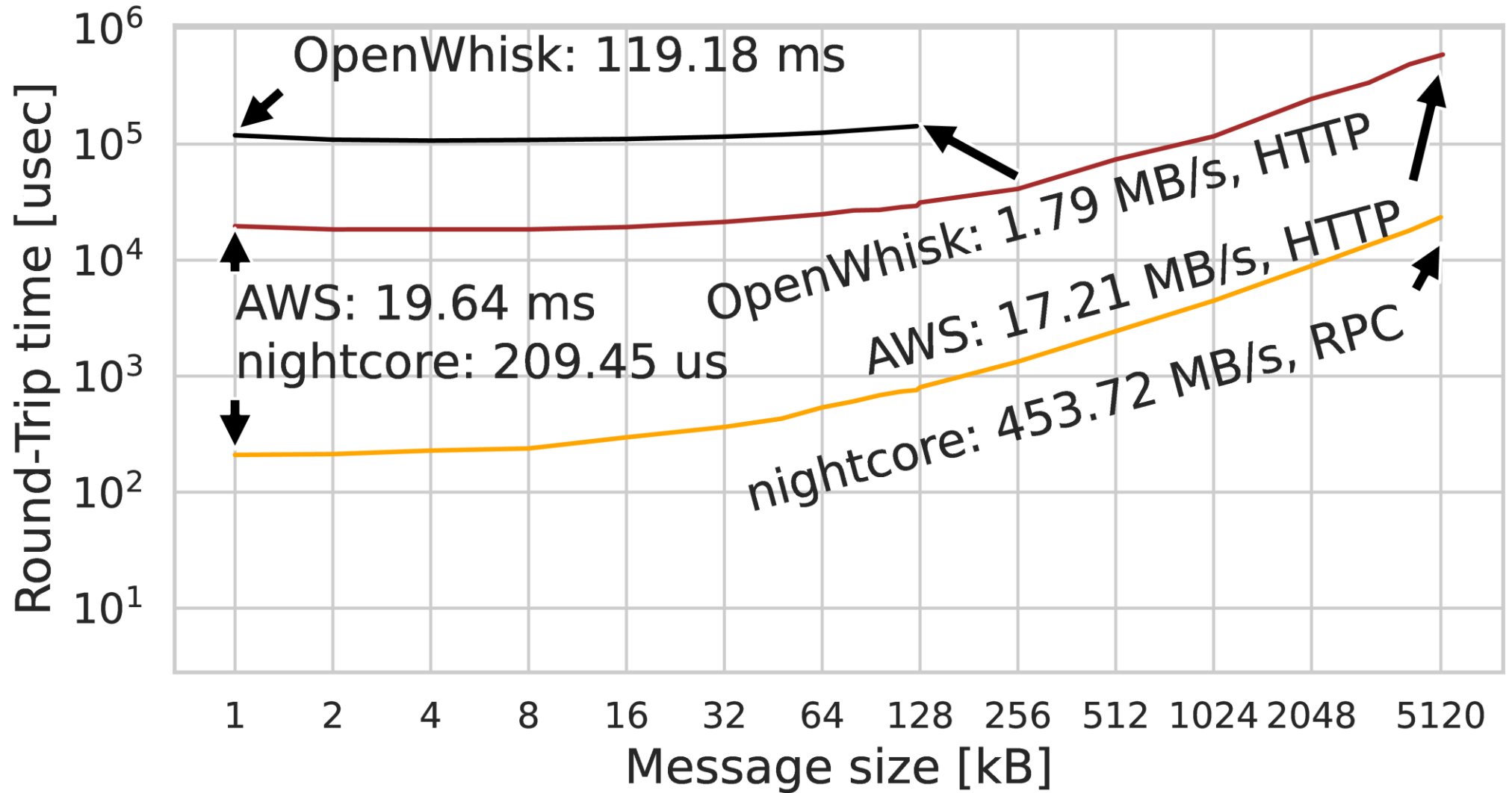
# How fast are invocations in FaaS?

36 CPU cores, 377 GB memory.
100 Gbps Ethernet with RoCEv2 support.

Reduced invocation critical path

Zero-copy RDMA networking



OpenWhisk: 119.18 ms

AWS: 19.64 ms
nightcore: 209.45 us

Warm: 9.3 us
Hot: 5.3 us

OpenWhisk: 1.79 MB/s, HTTP

AWS: 17.21 MB/s, HTTP

nightcore: 453.72 MB/s, RPC

rFaaS: 12 GB/s, RDMA

Round-Trip time [usec]

Message size [kB]

"rFaaS: Enabling High Performance Serverless with RDMA and Leases", IPDPS'23

# Serverless for High-Performance Applications

# Serverless for High-Performance Applications

Functions are expensive to invoke.

Communication is slow and restricted.

Applications

**SeBS**

Programming

Runtime

**rFaaS**

Hardware

# Serverless for High-Performance Applications

**Functions are expensive to invoke.**

**Communication is slow and restricted.**

Applications

SeBS

Programming

FMI

Runtime

rFaaS

Hardware

# Communication in serverless

**"FMI: Fast and Cheap Message Passing for Serverless Functions", ICS'23**

# Communication in serverless



"FMI: Fast and Cheap Message Passing for Serverless Functions", ICS'23

# Communication in serverless



"FMI: Fast and Cheap Message Passing for Serverless Functions", ICS'23

# Communication in serverless



**Cloud Storage**

"FMI: Fast and Cheap Message Passing for Serverless Functions", ICS'23

# Communication in serverless

**High Latency**
**For Small Messages**



**S3**

**Cloud Storage**



"FMI: Fast and Cheap Message Passing for Serverless Functions", ICS'23

# Communication in serverless

High Latency
For Small Messages

Expensive for
Large Messages



S3



DynamoDB

**Cloud Storage**



"FMI: Fast and Cheap Message Passing for Serverless Functions", ICS'23

# Communication in serverless



High Latency
For Small Messages

**S3**

Expensive for
Large Messages

**DynamoDB**

Not Serverless

**Redis**

**Cloud Storage**

**"FMI: Fast and Cheap Message Passing for Serverless Functions", ICS'23**

# Communication in serverless

"FMI: Fast and Cheap Message Passing for Serverless Functions", ICS'23

# Communication in serverless



**Hole Puncher**

"FMI: Fast and Cheap Message Passing for Serverless Functions", ICS'23

# Communication in serverless



**Hole Puncher**

"FMI: Fast and Cheap Message Passing for Serverless Functions", ICS'23

# FMI on AWS Lambda

**"FMI: Fast and Cheap Message Passing for Serverless Functions", ICS'23**

# FMI on AWS Lambda



**S3**    **Redis**    **TCP**

allreduce         bcast         gather

reduce         scan         scatter

**"FMI: Fast and Cheap Message Passing for Serverless Functions", ICS'23**

12

# FMI on AWS Lambda

**"FMI: Fast and Cheap Message Passing for Serverless Functions", ICS'23**

# Serverless for High-Performance Applications

**Functions are expensive to invoke.**

**Communication is slow and restricted.**

Applications

SeBS

Programming

FMI

Runtime

rFaaS

Hardware

# Serverless for High-Performance Applications

Functions are expensive to invoke.

Communication is slow and restricted.

Serverless is hard to implement in practice.

Applications

SeBS

Programming

FMI

Runtime

rFaaS

Hardware

# Serverless for High-Performance Applications



Functions are expensive to invoke.

Communication is slow and restricted.

Serverless is hard to implement in practice.

Applications
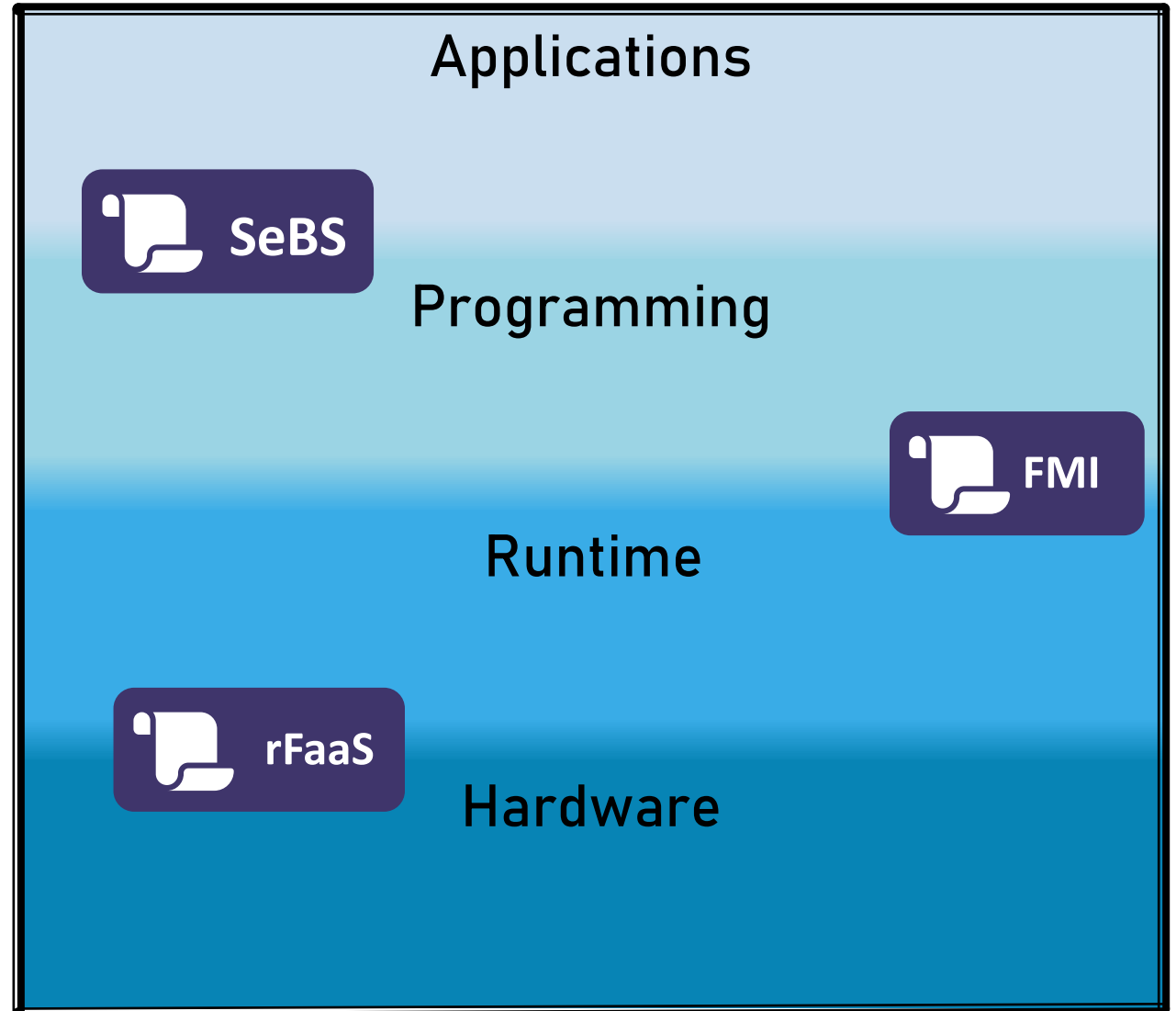
SeBS

PraaS

Programming

FMI

Runtime

rFaaS

Hardware

# Serverless Process



OS Process
Nano- and micro-second
latency of OS primitives.

State   IPC   Fork

# Serverless Process



**OS Process**
Nano- and micro-second latency of OS primitives.

**Serverless Function**
Millisecond latency of cloud proxies.

# Serverless Process



**OS Process**
Nano- and micro-second latency of OS primitives.

**Serverless Function**
Millisecond latency of cloud proxies.

"Process-as-a-Service: Elastic and Stateful Serverless with Cloud Processes", paper preprint.

# Serverless Process



**OS Process**
Nano- and micro-second latency of OS primitives.

**Serverless Function**
Millisecond latency of cloud proxies.

**Serverless Process**
Microsecond latency of PraaS backend.
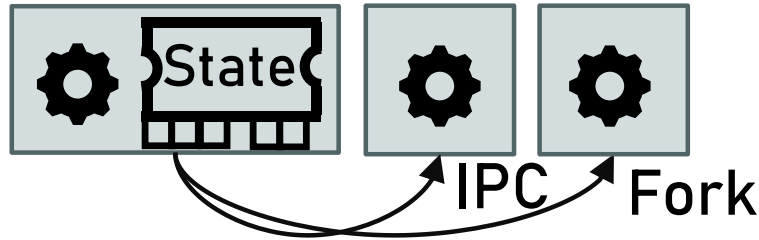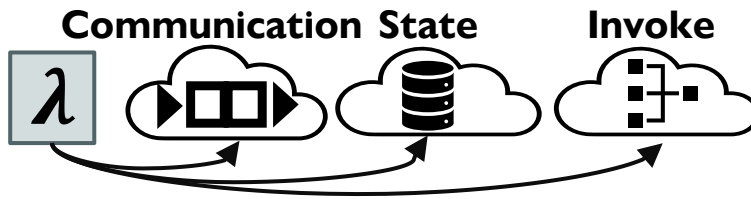
**"Process-as-a-Service: Elastic and Stateful Serverless with Cloud Processes", paper preprint.**
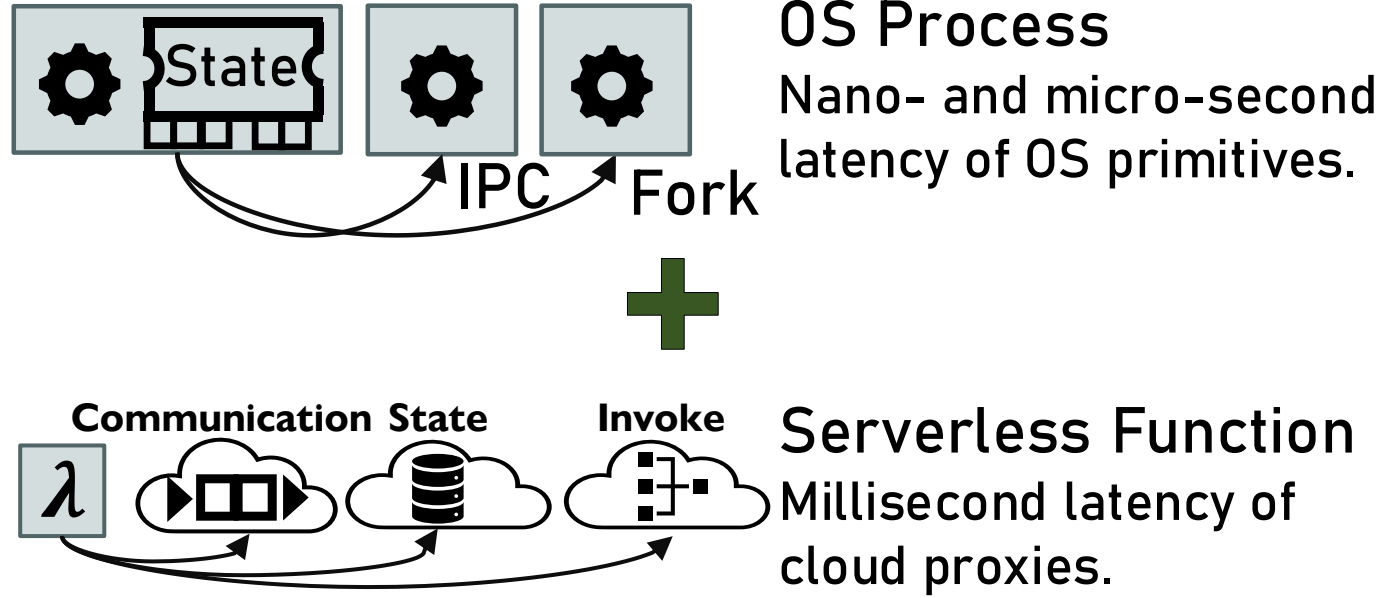
14

# Serverless Process



**OS Process**
Nano- and micro-second latency of OS primitives.

**Serverless Function**
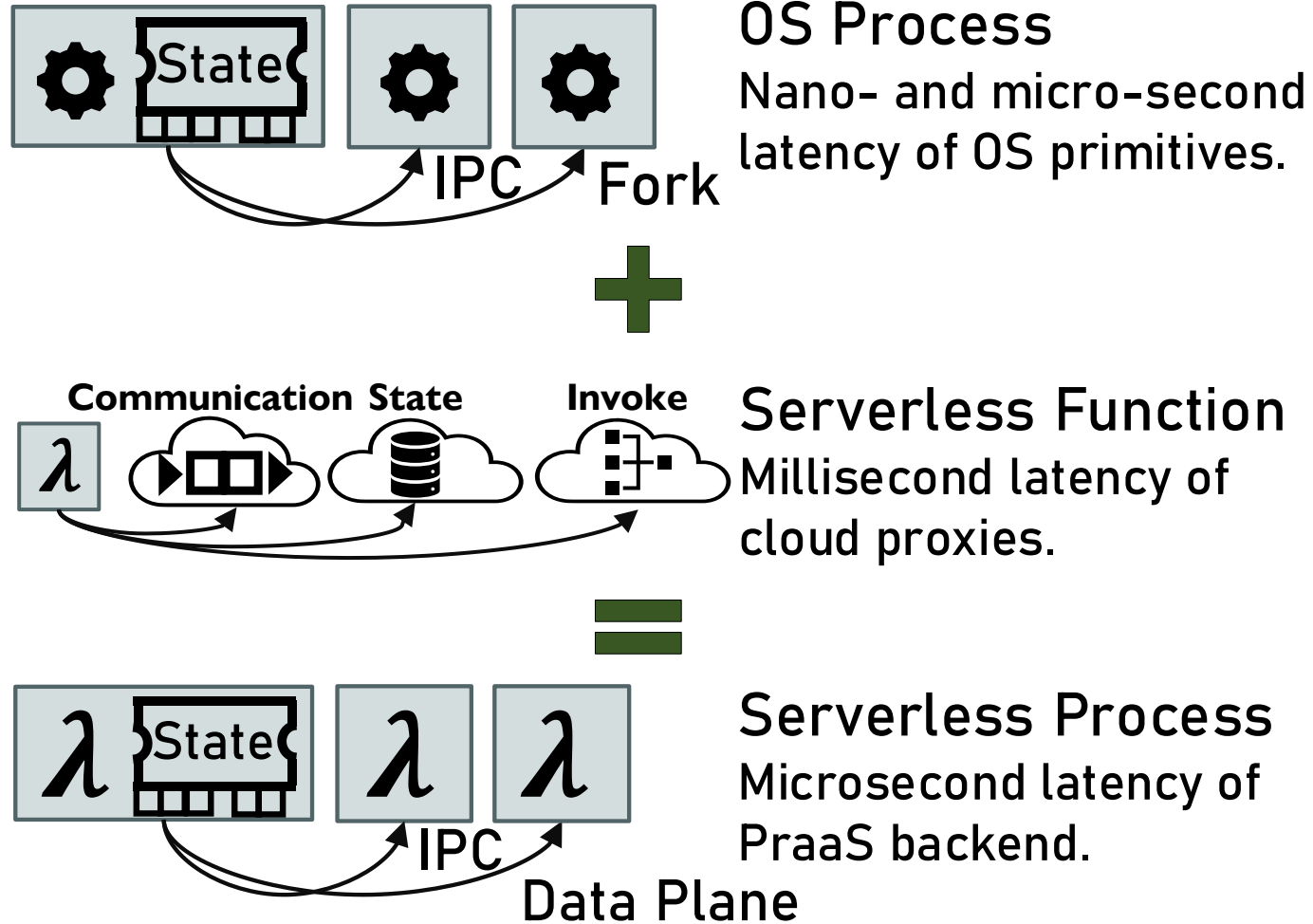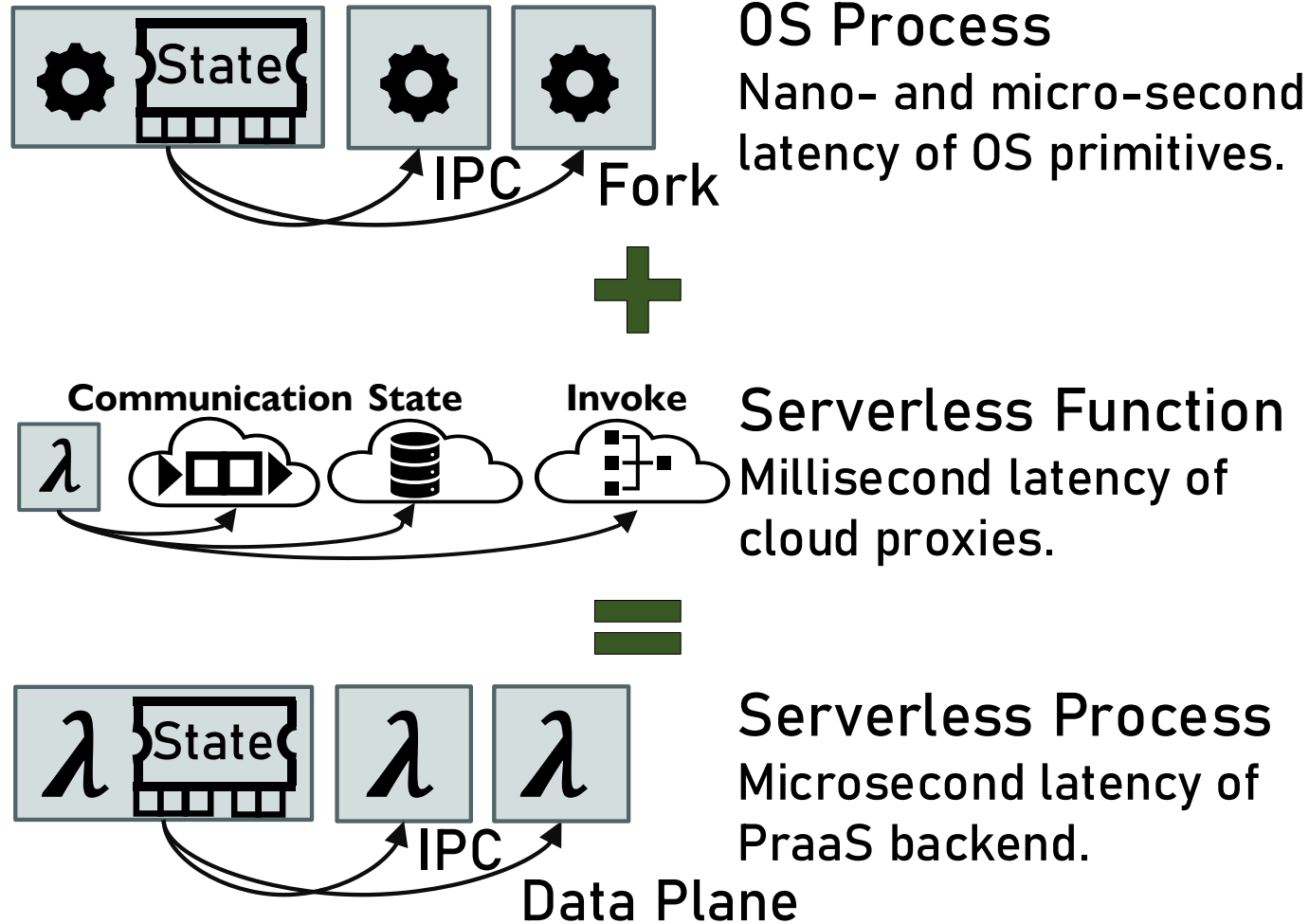Millisecond latency of cloud proxies.

**Serverless Process**
Microsecond latency of PraaS backend.

**Works on AWS Fargate, Knative, Kubernetes.**

"Process-as-a-Service: Elastic and Stateful Serverless with Cloud Processes", paper preprint.

# Benchmark: LaTeX Microservice



Legend: Lambda, 443 MiB | Lambda, 1769 MiB | PraaS, 0.25 vCPU | PraaS, 1 vCPU | Lambda, 885 MiB | Lambda, 3538 MiB | PraaS, 0.50 vCPU | PraaS, 2 vCPU

Y-axis: Time [ms]

X-axis categories: get-file, Input 1; get-file, Input 2; get-file, Input 3; get-file, Input 4; get-pdf, Input 1; update-file, Input 1; update-file, Input 2; update-file, Input 3; update-file, Input 4

X-axis label: Service, benchmarking scenario

**"Process-as-a-Service: Elastic and Stateful Serverless with Cloud Processes", paper preprint.**

# Benchmark: LaTeX Microservice

**"Process-as-a-Service: Elastic and Stateful Serverless with Cloud Processes", paper preprint.**

# Benchmark: LaTeX Microservice



**"Process-as-a-Service: Elastic and Stateful Serverless with Cloud Processes"**, paper preprint.

# Serverless for High-Performance Applications

Functions are expensive to invoke.

Communication is slow and restricted.

Serverless is hard to program.

Applications

SeBS

Programming

PraaS

FMI

Runtime

rFaaS

Hardware

# Serverless for High-Performance Applications

Functions are expensive to invoke.

Communication is slow and restricted.

Serverless is hard to program.

How to port existing and complex systems?

Applications

FaaSKeeper

SeBS

Programming
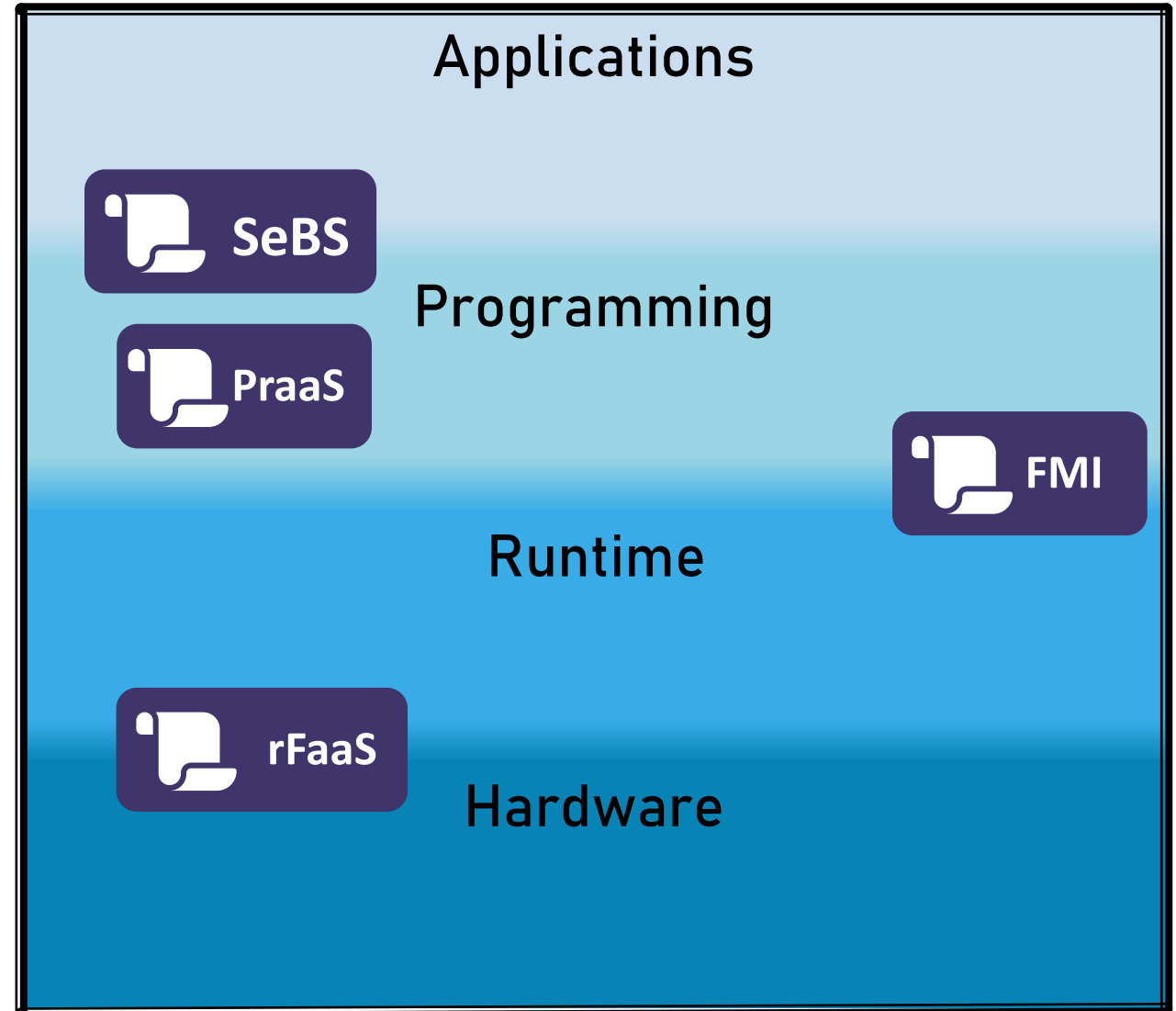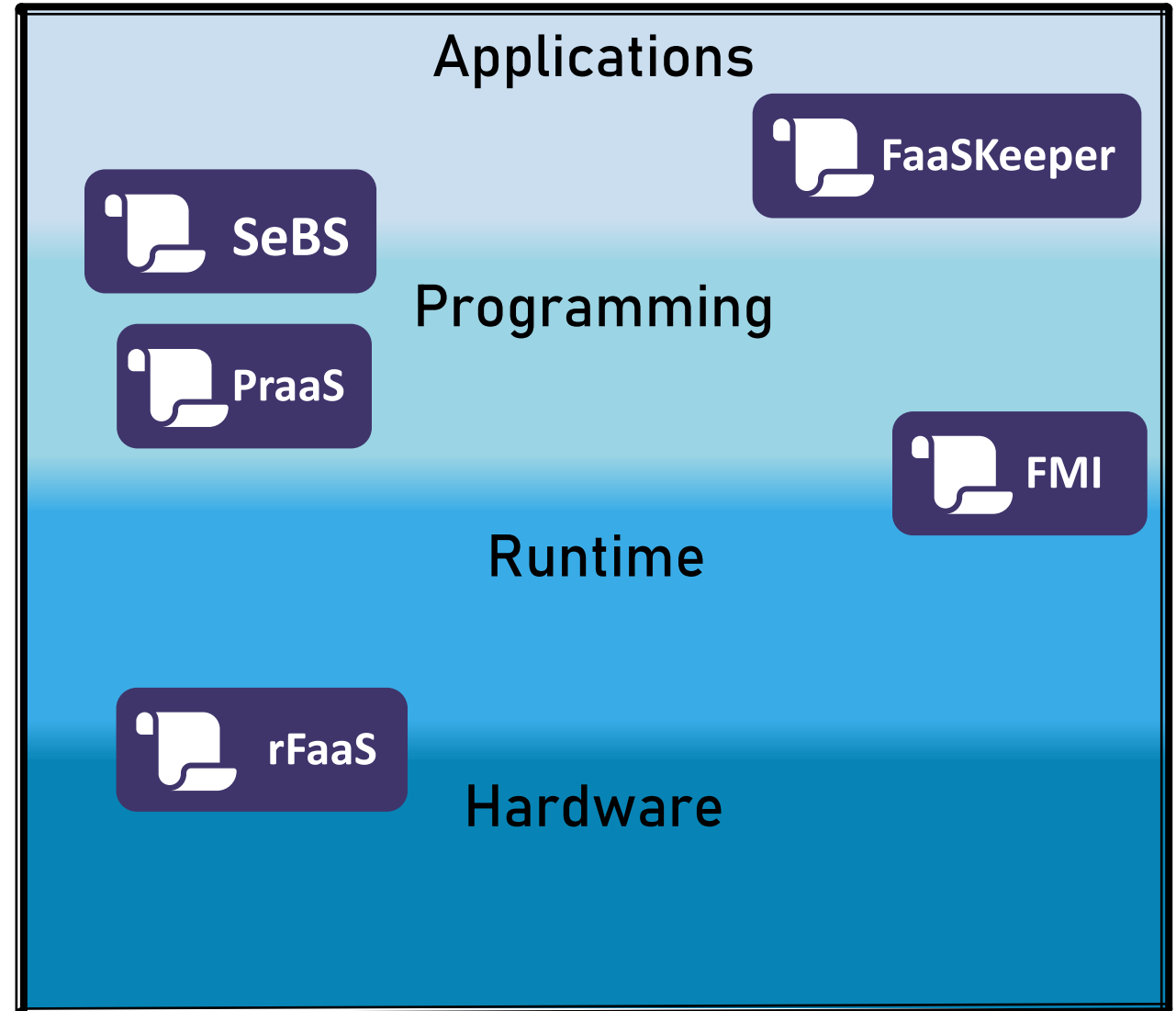
PraaS

FMI

Runtime

rFaaS

Hardware

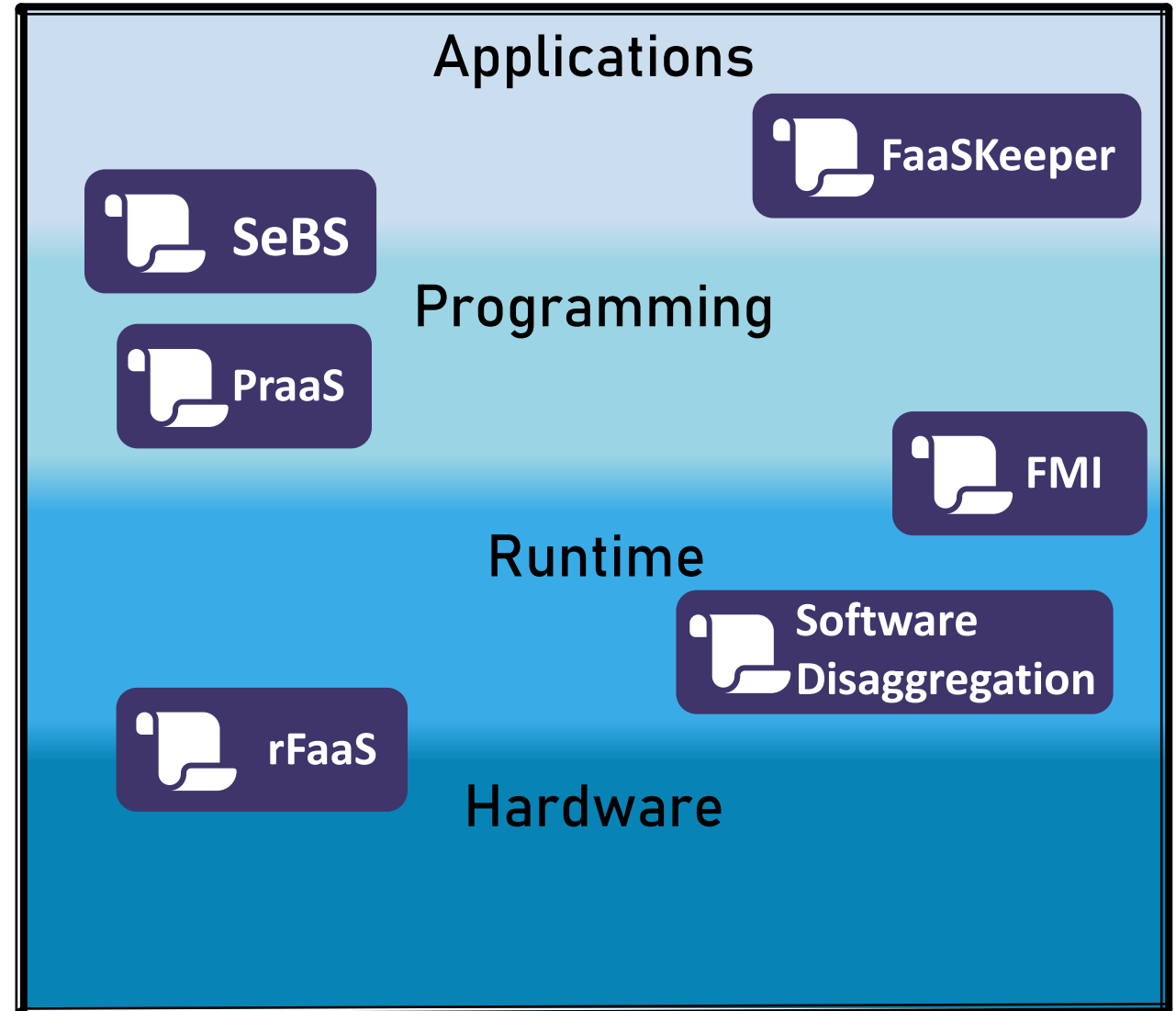# Serverless for High-Performance Applications

Functions are expensive to invoke.

Communication is slow and restricted.

Serverless is hard to program.

How to port existing and complex systems?

How can serverless improve HPC?

# Serverless Solutions for HPC

# Serverless Solutions for HPC

 spcl/serverless-benchmarks

 spcl/rFaaS

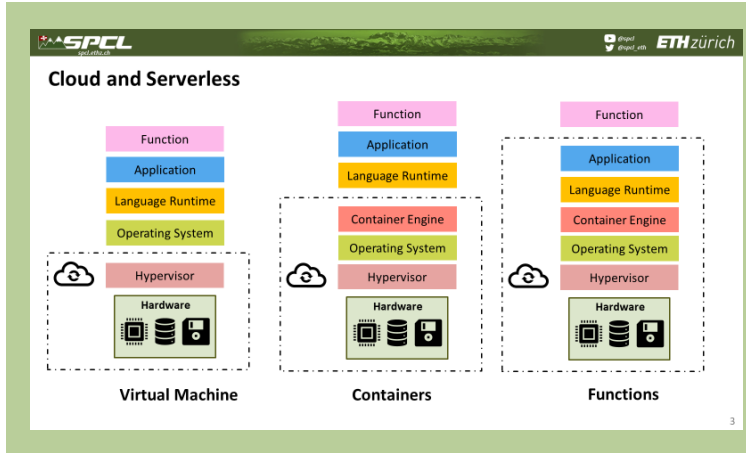 spcl/fmi

 spcl/PraaS

 spcl/FaaSKeeper

# Conclusions

## More of SPCL's research:

youtube.com/@spcl — **150+ Talks**

twitter.com/spcl_eth — **1.2K+ Followers**

github.com/spcl — **2K+ Stars**

**... or** spcl.ethz.ch



### Cloud and Serverless

# Conclusions

**Cloud and Serverless**

| Virtual Machine | Containers | Functions |
| --- | --- | --- |

Function
Application
Language Runtime
Operating System
Hypervisor
Hardware

Function
Application
Language Runtime
Container Engine
Operating System
Hypervisor
Hardware

Function
Application
Language Runtime
Container Engine
Operating System
Hypervisor
Hardware

**How fast rFaaS invocations are?**

36 CPU cores, 377 GB memory. Ethernet with RoCEv2 support.

OpenWhisk: 119.18 ms

AWS: 19.64 ms
nightcore: 209.45 us

Warm: 9.3 us
Hot: 5.3 us

OpenWhisk: 1.79 MB/s

AWS: 17.21 MB/s

nightcore: 453.72 MB/s

rFaaS: 12 GB/s

Round-Trip Time [usec] vs Function argument size [kB]

# Conclusions

# Conclusions

**More of SPCL's research:**

youtube.com/@spcl — 150+ Talks

twitter.com/spcl_eth — 1.2K+ Followers

github.com/spcl — 2K+ Stars

… or spcl.ethz.ch

# Conclusions

**More of SPCL's research:**



| youtube.com/@spcl | 150+ Talks |
| twitter.com/spcl_eth | 1.2K+ Followers |
| github.com/spcl | 2K+ Stars |

**... or spcl.ethz.ch**



**Cloud and Serverless**

Function
Application
Language Runtime
Operating System
Hypervisor
Hardware

Function
Application
Language Runtime
Container Engine
Operating System
Hypervisor
Hardware

Function
Application
Language Runtime
Container Engine
Operating System
Hypervisor
Hardware

Virtual Machine     Containers     Functions

**How fast rFaaS invocations are?**

36 CPU cores, 377 GB memory. Ethernet with RoCEv2 support.

OpenWhisk: 119.18 ms
OpenWhisk: 1.79 MB/s
AWS: 19.64 ms
nightcore: 209.45 us
AWS: 17.21 MB/s
nightcore: 453.72 MB/s
Warm: 9.3 us
Hot: 5.3 us
rFaaS: 12 GB/s

Round-Trip Time [usec] vs Function argument size [kB]

**FMI on AWS Lambda**

allreduce    bcast    gather
reduce    scan    scatter

**Serverless Process**

**OS Process**
Nano- and micro-second latency of OS primitives.
State — IPC — Fork

**+**

**Serverless Function**
Millisecond latency of cloud proxies.
Communication State — Invoke

**=**

**Serverless Process**
Microsecond latency of PraaS backend.
λ State λ λ — IPC
Data Plane

**Poster**     **All projects.**

     

This work has received funding from the European Research Council (ERC), Swiss National Science Foundation (SNF), and from Amazon Web Services through the AWS Cloud Credits for Research program.
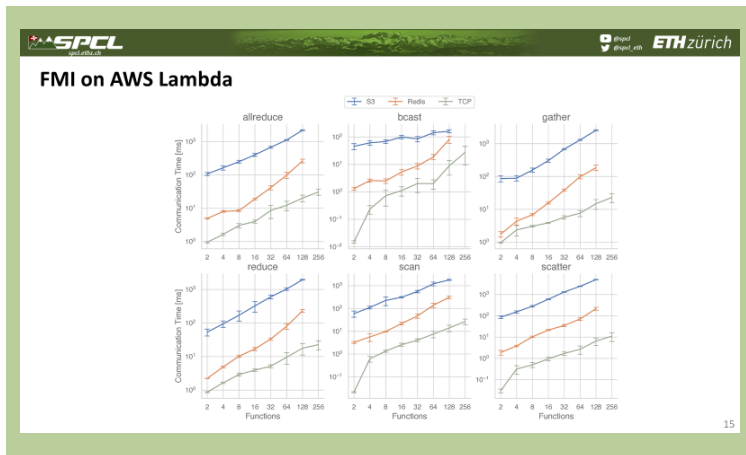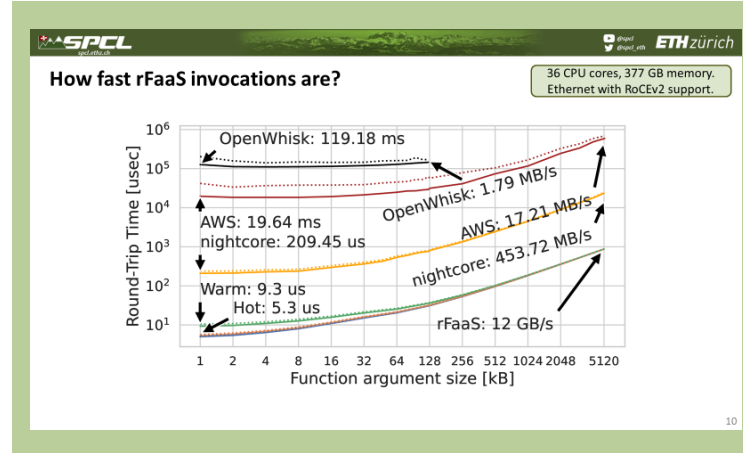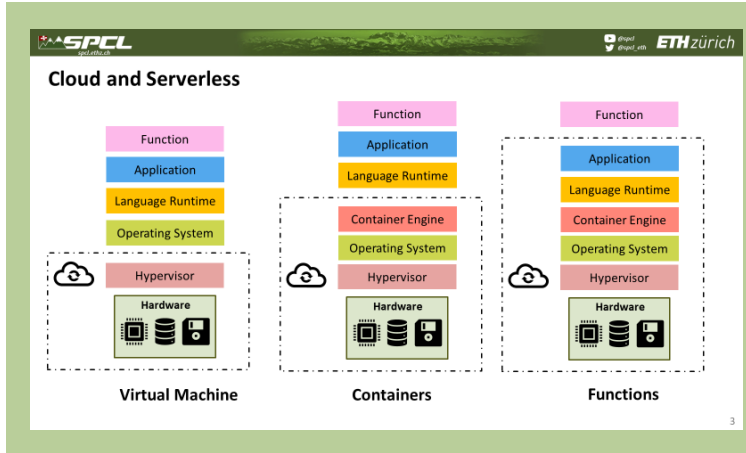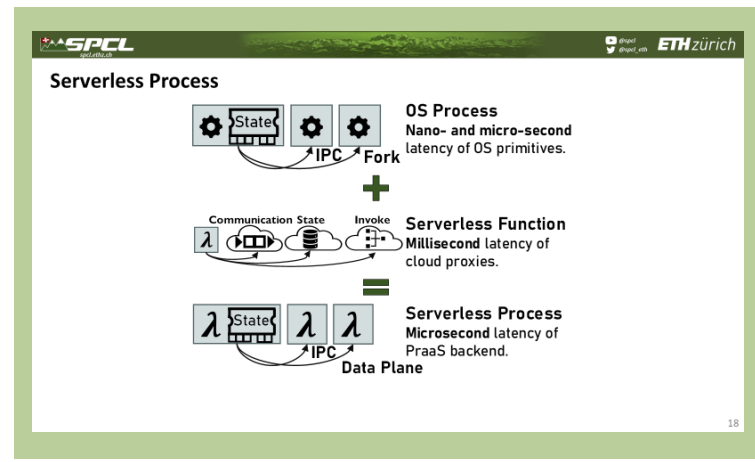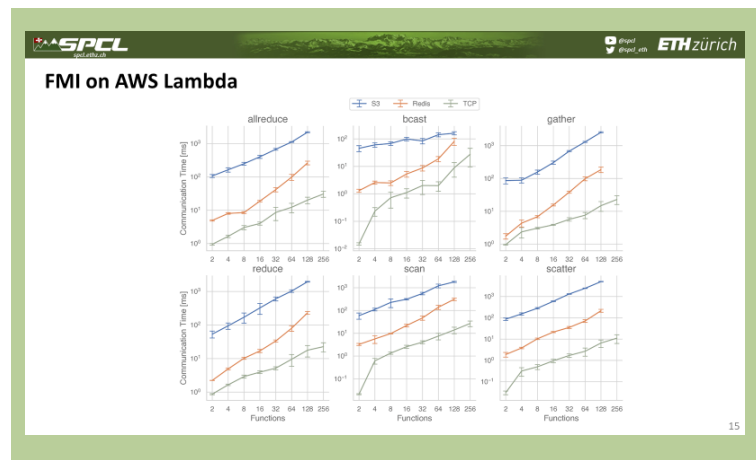
# "But serverless is slow and expensive"
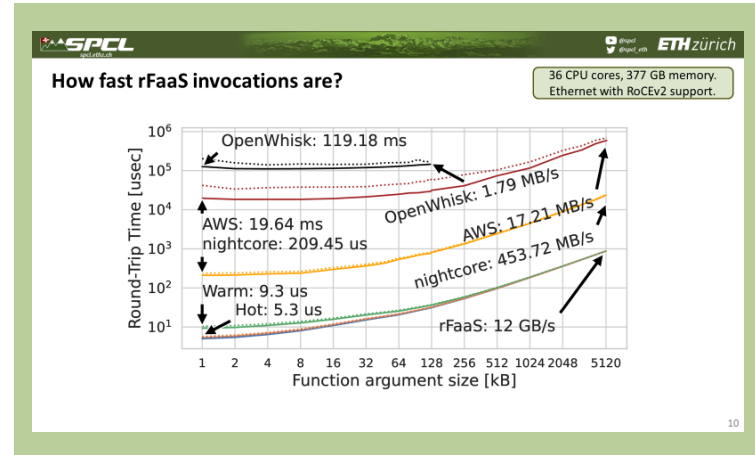
# "But serverless is slow and expensive"



Scaling up the Prime Video audio/video monitoring service and reducing costs by 90%

The move from a distributed microservices architecture to a monolith application helped achieve higher scale, resilience, and reduce costs.

# "But serverless is slow and expensive"



Scaling up the Prime Video audio/video monitoring service and reducing costs by 90%

The move from a distributed microservices architecture to a monolith application helped achieve higher scale, resilience, and reduce costs.

# Benchmarking Serverless

**"SeBS: a Serverless Benchmark Suite for Function-as-a-Service Computing", ACM/IFIP Middleware 2021**

# Benchmarking Serverless



**"SeBS: a Serverless Benchmark Suite for Function-as-a-Service Computing", ACM/IFIP Middleware 2021**

# Benchmarking Serverless



| Type | Name | Language |
|------|------|----------|
| **Webapps** | uploader | **Python**, Node.js |
| **Multimedia** | thumbnailer | **Python**, **Node.js**, C++ |
| **Utilities** | compression | **Python** |
| **Inference** | image-recognition | **Python**, C++ |
| **Scientific** | graph-bfs | **Python** |

**"SeBS: a Serverless Benchmark Suite for Function-as-a-Service Computing", ACM/IFIP Middleware 2021**

# Benchmarking Serverless



| Type | Name | Language |
|------|------|----------|
| Webapps | uploader | **Python**, Node.js |
| Multimedia | thumbnailer | **Python**, **Node.js**, C++ |
| Utilities | compression | **Python** |
| Inference | image-recognition | **Python**, C++ |
| Scientific | graph-bfs | **Python** |

| | Results, methods, and insights |
|---|---|
| | **High-memory allocations increase cold startup overheads on GCP.** |
| | **GCP functions experience reliability and availability issues.** |
| | I/O-bound functions experience very high latency variations. |
| | AWS Lambda achieves the best performance on all workloads. |
| | Irregular performance of concurrent Azure Function executions. |
| | AWS Lambda performance is not competitive against VMs assuming comparable resources. |
| | **Break-even analysis for IaaS and FaaS deployment.** |
| | Resource underutilization due to high granularity of pricing models. |
| | High costs of Azure Functions due to unconfigurable deployment. |
| | The function output size can be a dominating factor in pricing. |
| | Accurate methodology for estimation of invocation latency. |
| | Warm latencies are consistent and depend linearly on payload size. |
| | Highly variable and unpredictable cold latencies on Azure and GCP. |
| | **AWS Lambda container eviction is agnostic to function properties.** |
| | **Analytical models of AWS Lambda container eviction policy.** |

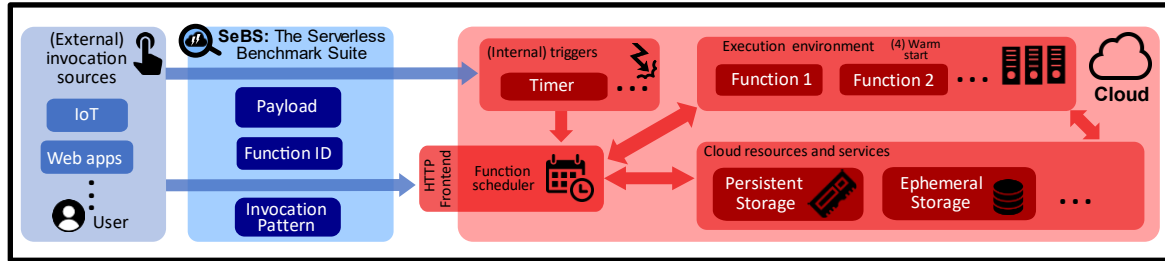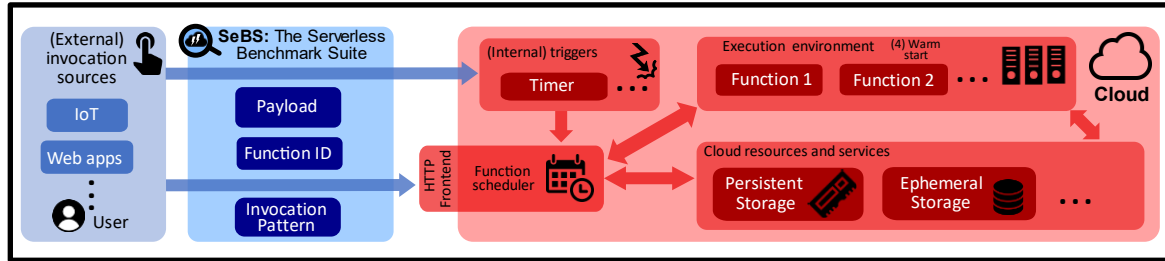**"SeBS: a Serverless Benchmark Suite for Function-as-a-Service Computing", ACM/IFIP Middleware 2021**
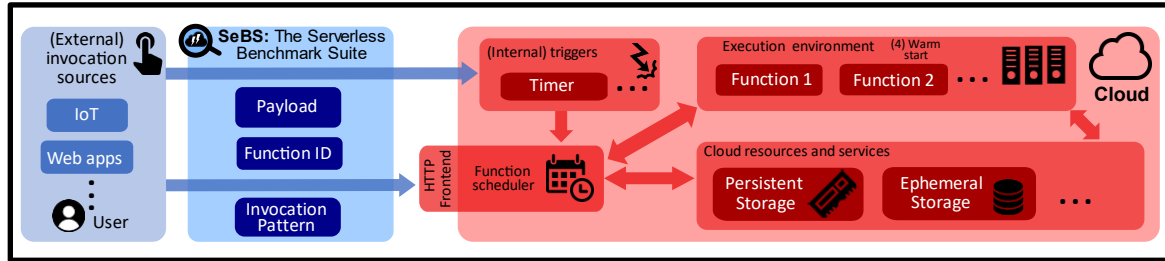
# Benchmarking Serverless



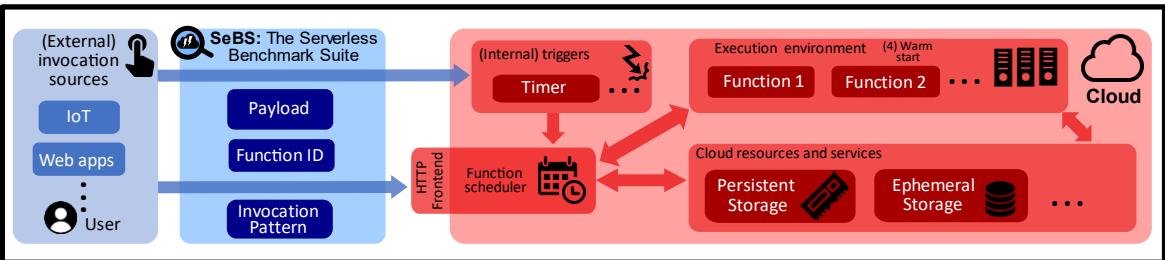| Type | Name | Language |
|------|------|----------|
| **Webapps** | uploader | **Python**, Node.js |
| **Multimedia** | thumbnailer | **Python**, **Node.js**, C++ |
| **Utilities** | compression | **Python** |
| **Inference** | image-recognition | **Python**, C++ |
| **Scientific** | graph-bfs | **Python** |

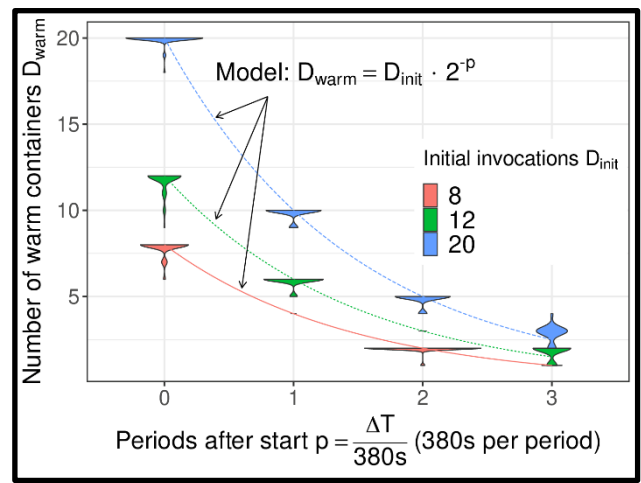| Results, methods, and insights |
|---|
| **High-memory allocations increase cold startup overheads on GCP.** |
| **GCP functions experience reliability and availability issues.** |
| I/O-bound functions experience very high latency variations. |
| AWS Lambda achieves the best performance on all workloads. |
| Irregular performance of concurrent Azure Function executions. |
| AWS Lambda performance is not competitive against VMs assuming comparable resources. |
| **Break-even analysis for IaaS and FaaS deployment.** |
| Resource underutilization due to high granularity of pricing models. |
| High costs of Azure Functions due to unconfigurable deployment. |
| The function output size can be a dominating factor in pricing. |
| Accurate methodology for estimation of invocation latency. |
| Warm latencies are consistent and depend linearly on payload size. |
| Highly variable and unpredictable cold latencies on Azure and GCP. |
| **AWS Lambda container eviction is agnostic to function properties.** |
| **Analytical models of AWS Lambda container eviction policy.** |



Model: $D_{warm} = D_{init} \cdot 2^{-p}$

Initial invocations $D_{init}$
- 8
- 12
- 20

Number of warm containers $D_{warm}$

Periods after start $p = \dfrac{\Delta T}{380s}$ (380s per period)

**"SeBS: a Serverless Benchmark Suite for Function-as-a-Service Computing", ACM/IFIP Middleware 2021**

# Serverless for High-Performance Applications

Functions are expensive to invoke.

Communication is slow and restricted.

Serverless is hard to program.

Applications

SeBS

Programming

PraaS

FMI

Runtime

rFaaS

Hardware

# Serverless for High-Performance Applications

Functions are expensive to invoke.

Communication is slow and restricted.

Serverless is hard to program.

How to port existing and complex systems?

Applications

FaaSKeeper

SeBS

Programming

PraaS

FMI

Runtime

rFaaS

Hardware

# Serverless for High-Performance Applications

Functions are expensive to invoke.

Communication is slow and restricted.

Serverless is hard to program.

How to port existing and complex systems?

How can serverless improve HPC?

Applications

FaaSKeeper

SeBS

Programming

PraaS

FMI

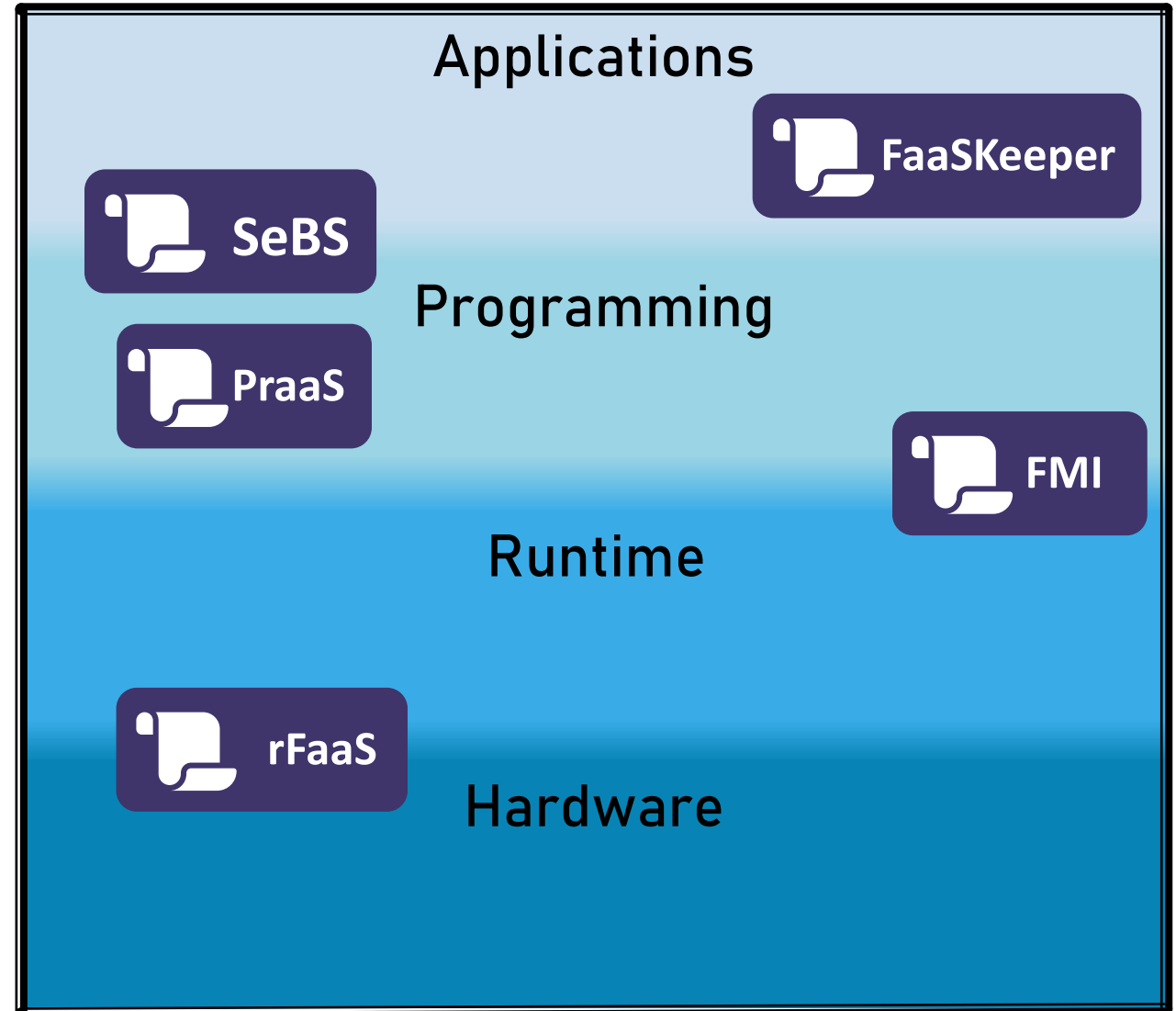Runtime

Software Disaggregation

rFaaS

Hardware

21

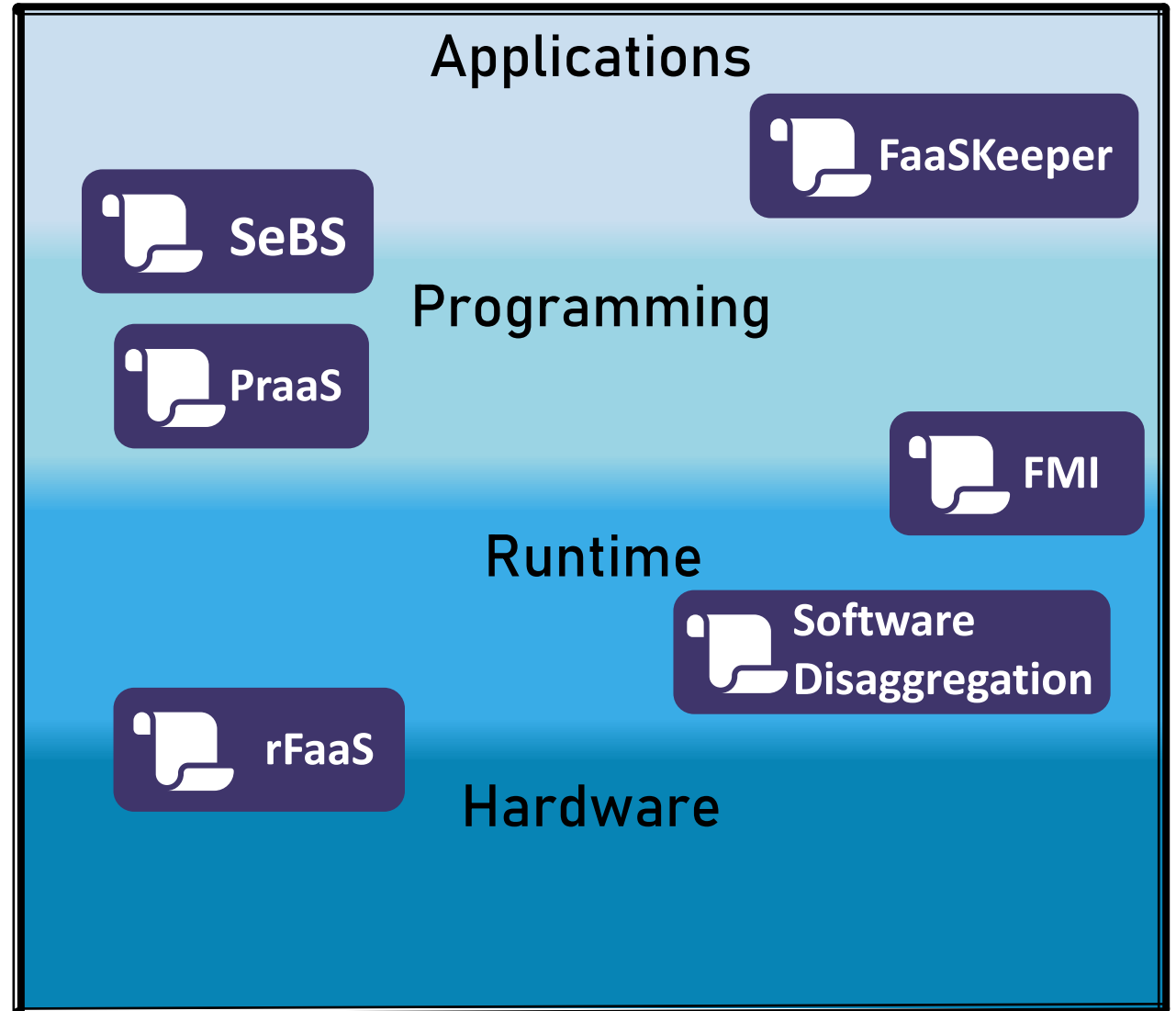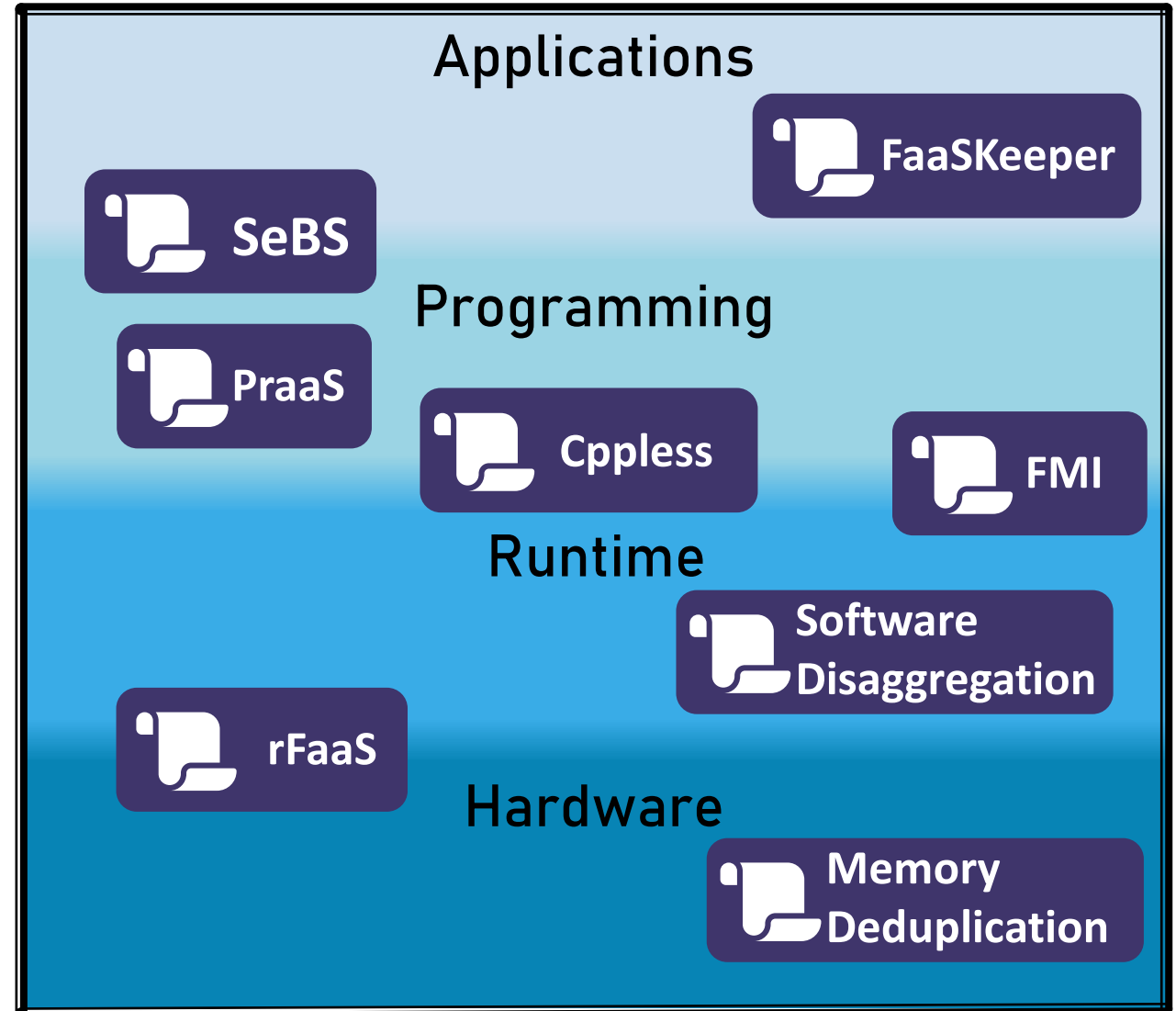# Serverless for High-Performance Applications

Functions are expensive to invoke.

Communication is slow and restricted.

Serverless is hard to program.

How to port existing and complex systems?

How can serverless improve HPC?



Applications

FaaSKeeper

SeBS

Programming

PraaS

Cppless

FMI

Runtime

Software Disaggregation

rFaaS

Hardware

Memory Deduplication

# Conclusions

Serverless Invocations

OpenWhisk: 119.18 ms

OpenWhisk: 1.79 MB/s, HTTP

AWS: 19.64 ms
nightcore: 209.45 us

AWS: 17.21 MB/s, HTTP

nightcore: 453.72 MB/s, RPC

Warm: 9.3 us
Hot: 5.3 us

rFaaS: 12 GB/s, RDMA

"rFaaS: Enabling High Performance Serverless with RDMA and Leases", IEEE IPDPS 2023

**mcopik.github.io/projects/praas**

# Conclusions

**More of SPCL's research:**



Serverless Invocations

OpenWhisk: 119.18 ms
AWS: 19.64 ms
nightcore: 209.45 us
Warm: 9.3 us
Hot: 5.3 us
OpenWhisk: 1.79 MB/s, HTTP
AWS: 17.21 MB/s, HTTP
nightcore: 453.72 MB/s, RPC
rFaaS: 12 GB/s, RDMA

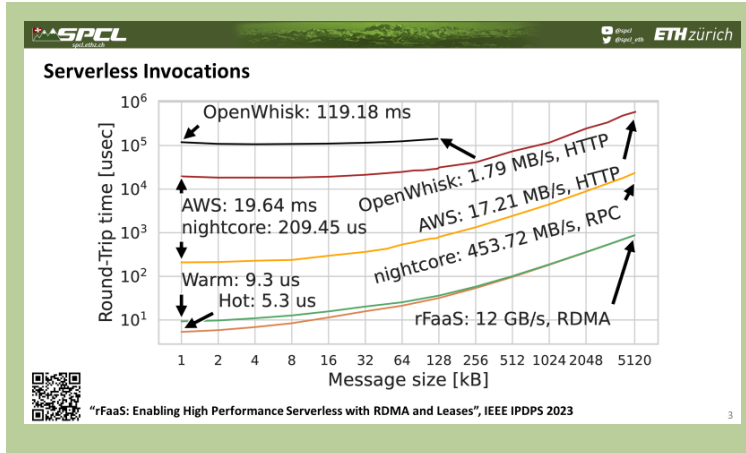"rFaaS: Enabling High Performance Serverless with RDMA and Leases", IEEE IPDPS 2023

| youtube.com/@spcl | 150+ Talks |
| twitter.com/spcl_eth | 1.2K+ Followers |
| github.com/spcl | 2K+ Stars |

**... or spcl.ethz.ch**

**mcopik.github.io/projects/praas**

# Conclusions

## More of SPCL's research:

youtube.com/@spcl | 150+ Talks

twitter.com/spcl_eth | 1.2K+ Followers

github.com/spcl | 2K+ Stars

... or spcl.ethz.ch



Serverless Invocations

OpenWhisk: 119.18 ms

OpenWhisk: 1.79 MB/s, HTTP

AWS: 19.64 ms
nightcore: 209.45 us

AWS: 17.21 MB/s, HTTP

nightcore: 453.72 MB/s, RPC

Warm: 9.3 us
Hot: 5.3 us

rFaaS: 12 GB/s, RDMA

"rFaaS: Enabling High Performance Serverless with RDMA and Leases", IEEE IPDPS 2023

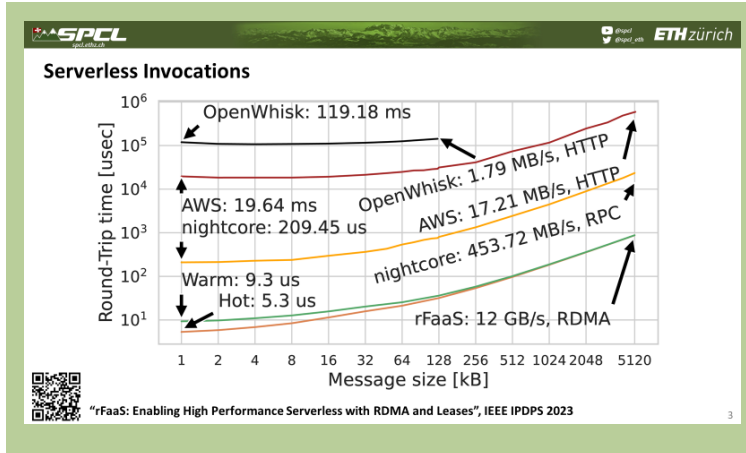mcopik.github.io/projects/praas

# Conclusions

**More of SPCL's research:**



### Serverless Invocations

OpenWhisk: 119.18 ms

OpenWhisk: 1.79 MB/s, HTTP

AWS: 19.64 ms
nightcore: 209.45 us

AWS: 17.21 MB/s, HTTP

nightcore: 453.72 MB/s, RPC

Warm: 9.3 us
Hot: 5.3 us

rFaaS: 12 GB/s, RDMA

"rFaaS: Enabling High Performance Serverless with RDMA and Leases", IEEE IPDPS 2023

youtube.com/@spcl — **150+ Talks**

twitter.com/spcl_eth — **1.2K+ Followers**

github.com/spcl — **2K+ Stars**

**... or spcl.ethz.ch**

**mcopik.github.io/projects/praas**

# Conclusions

## More of SPCL's research:


Serverless Invocations

"rFaaS: Enabling High Performance Serverless with RDMA and Leases", IEEE IPDPS 2023

youtube.com/@spcl — 150+ Talks

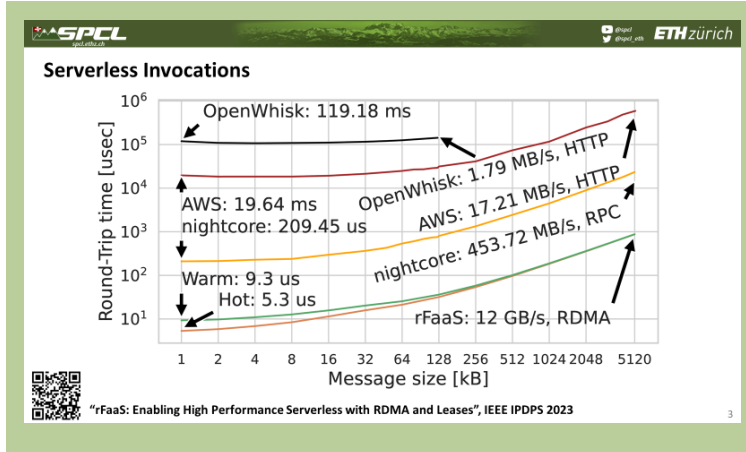twitter.com/spcl_eth — 1.2K+ Followers

github.com/spcl — 2K+ Stars

... or spcl.ethz.ch

mcopik.github.io/projects/praas

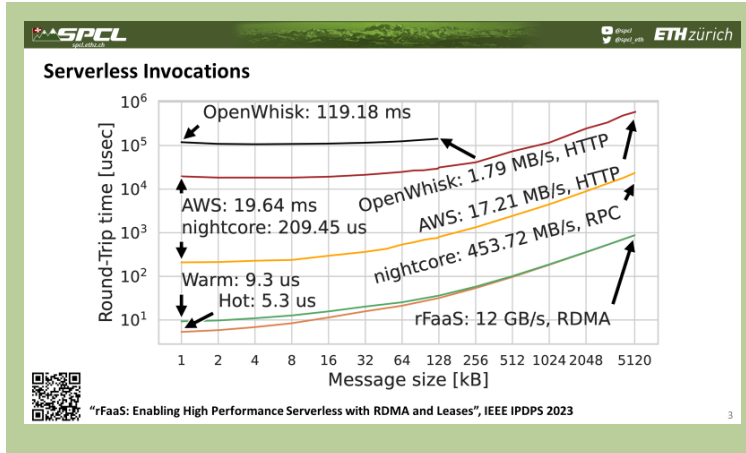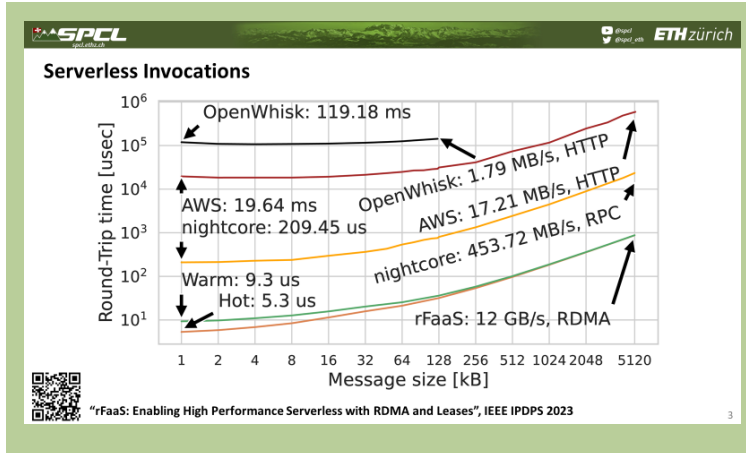# Conclusions

**More of SPCL's research:**

youtube.com/@spcl — **150+ Talks**

twitter.com/spcl_eth — **1.2K+ Followers**

github.com/spcl — **2K+ Stars**

**... or spcl.ethz.ch**



### Serverless Invocations

OpenWhisk: 119.18 ms

OpenWhisk: 1.79 MB/s, HTTP

AWS: 19.64 ms
nightcore: 209.45 us

AWS: 17.21 MB/s, HTTP

nightcore: 453.72 MB/s, RPC

Warm: 9.3 us
Hot: 5.3 us

rFaaS: 12 GB/s, RDMA

Round-Trip time [usec] vs Message size [kB]

"rFaaS: Enabling High Performance Serverless with RDMA and Leases", IEEE IPDPS 2023

**spcl/praas**

**Paper preprint**

mcopik.github.io/projects/praas