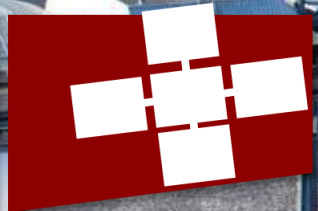


Serverless as a Bridge Between HPC and Clouds

Marcin Copik, Alexandru Calotoiu, Torsten Hoefler



IPOPS
2023 • St. Petersburg,
Florida USA

Serverless on servers.



Serverless on servers.



“But serverless is slow and expensive”

“But serverless is slow and expensive”

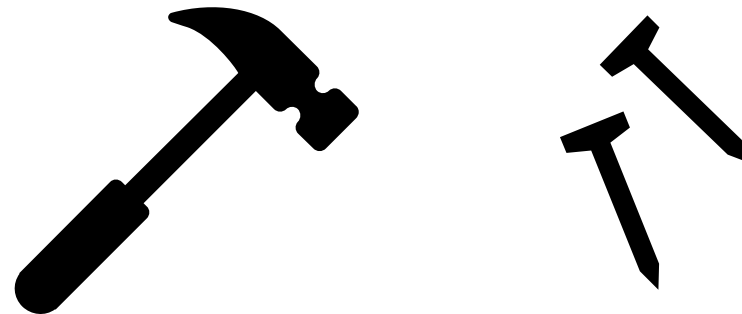
Scaling up the Prime Video audio/video monitoring service and reducing costs by 90%

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

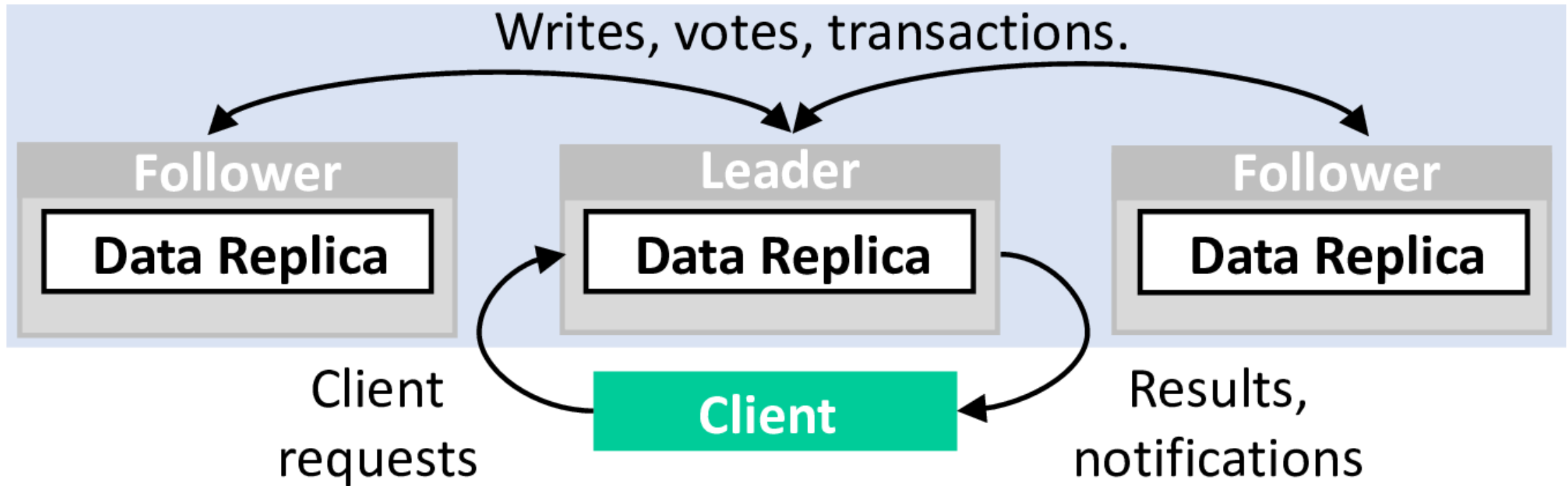
“But serverless is slow and expensive”

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

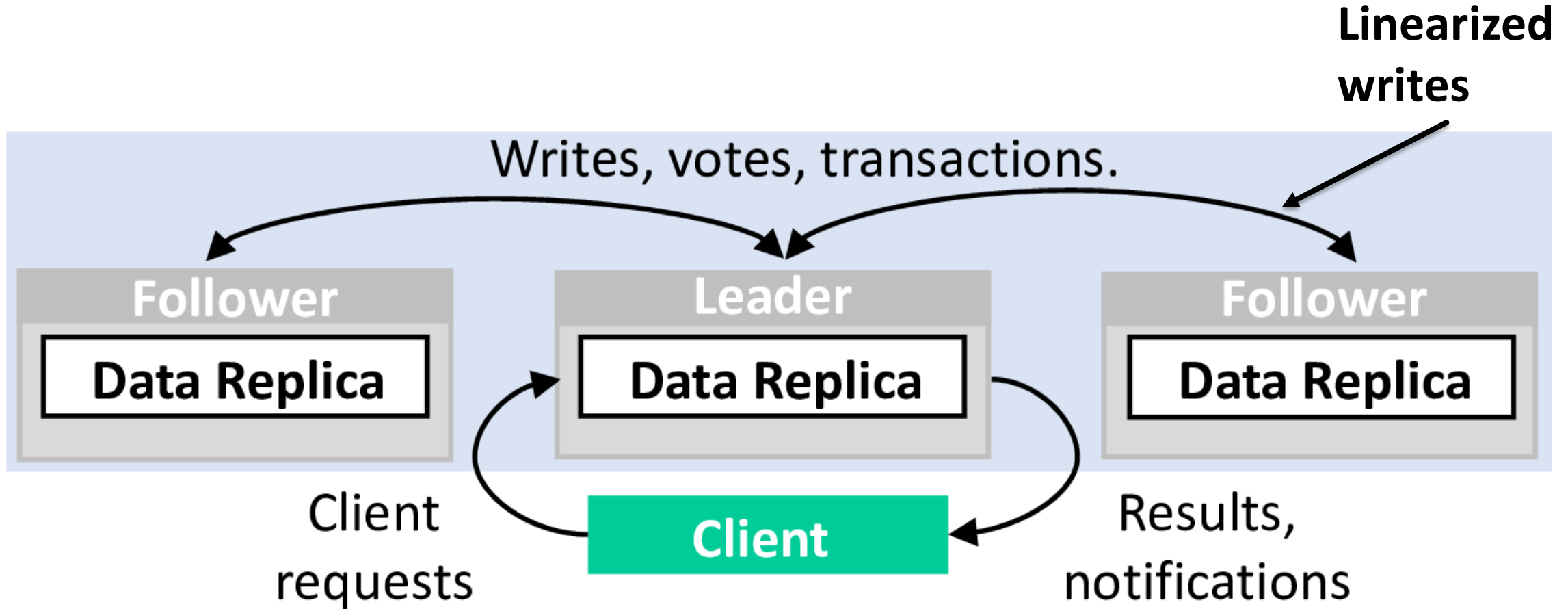
The move from a distributed microservices architecture to a monolith application helped achieve higher scale, resilience, and reduce costs.



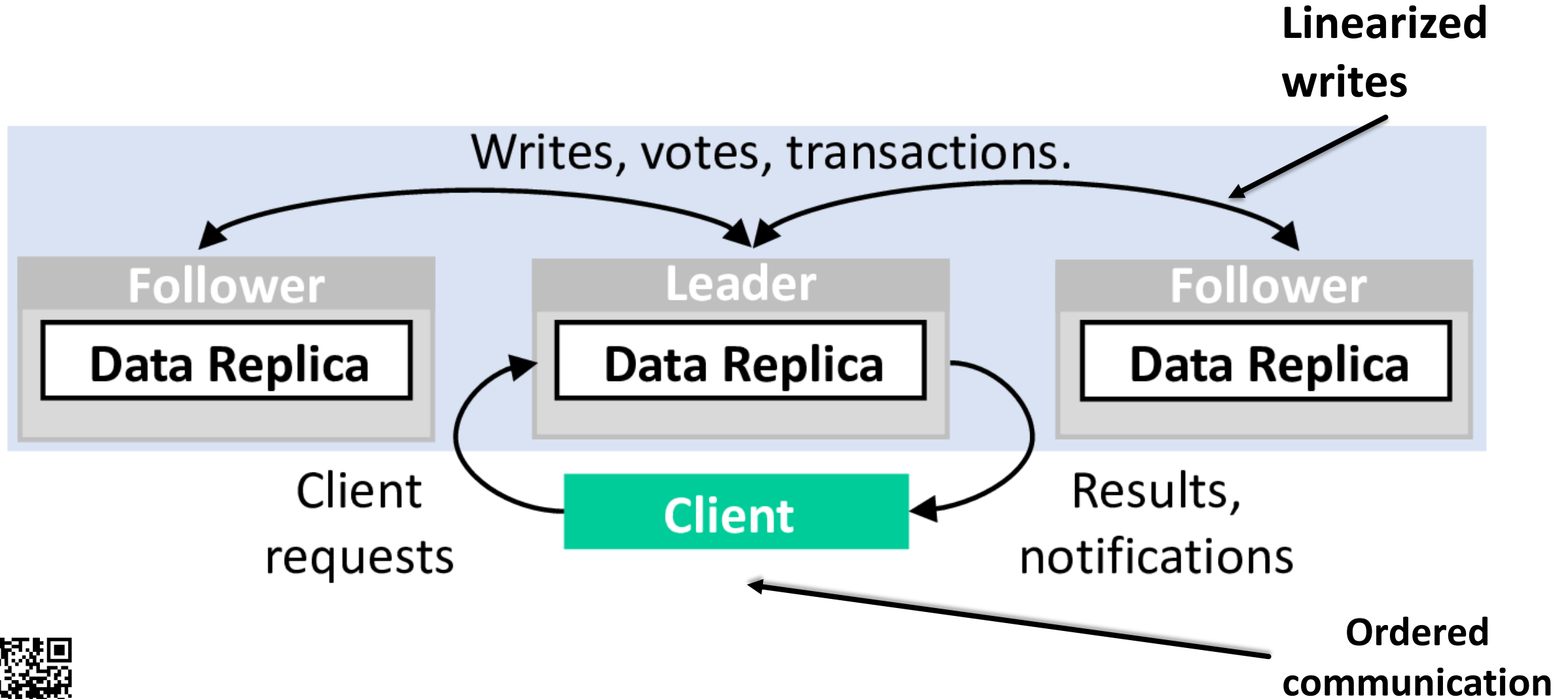
Building Serverless Services: FaaSKeeper



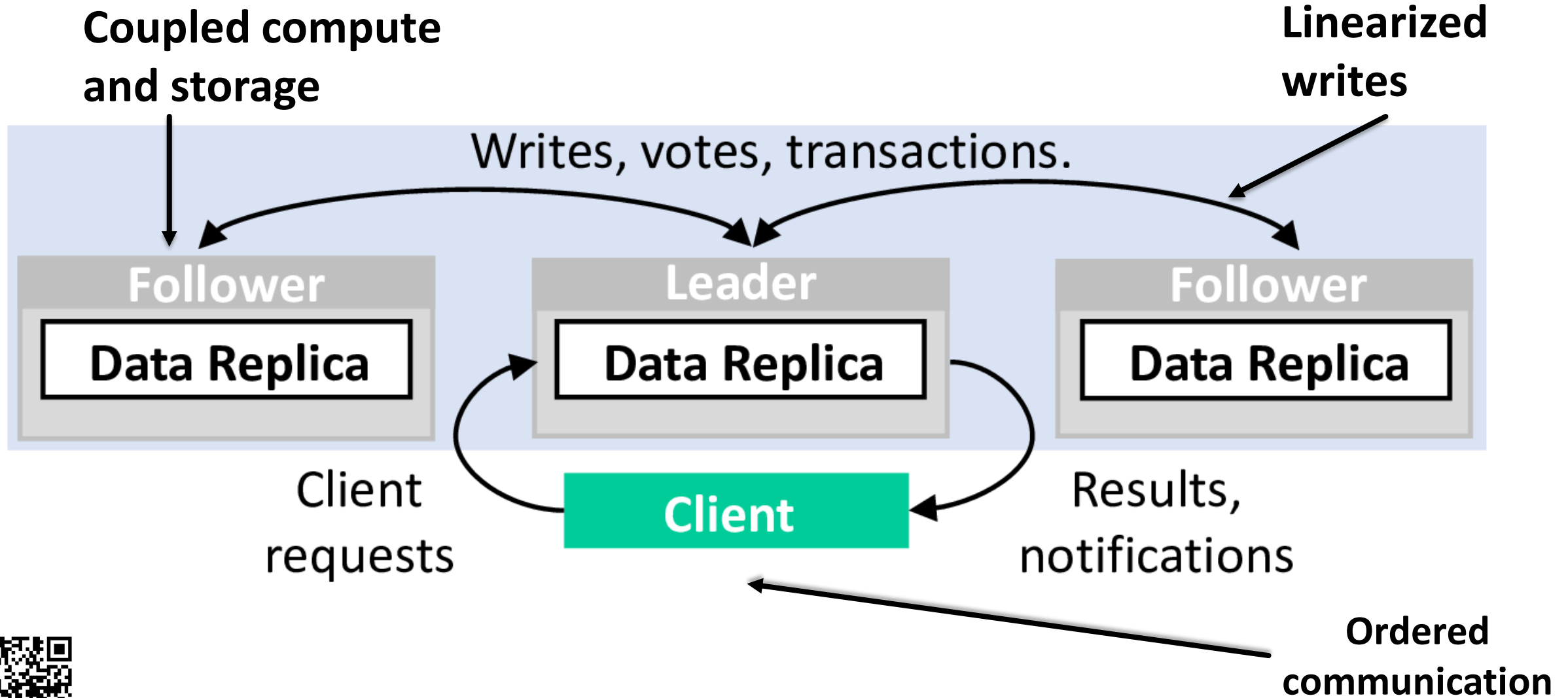
Building Serverless Services: FaaSKeeper



Building Serverless Services: FaaSKeeper

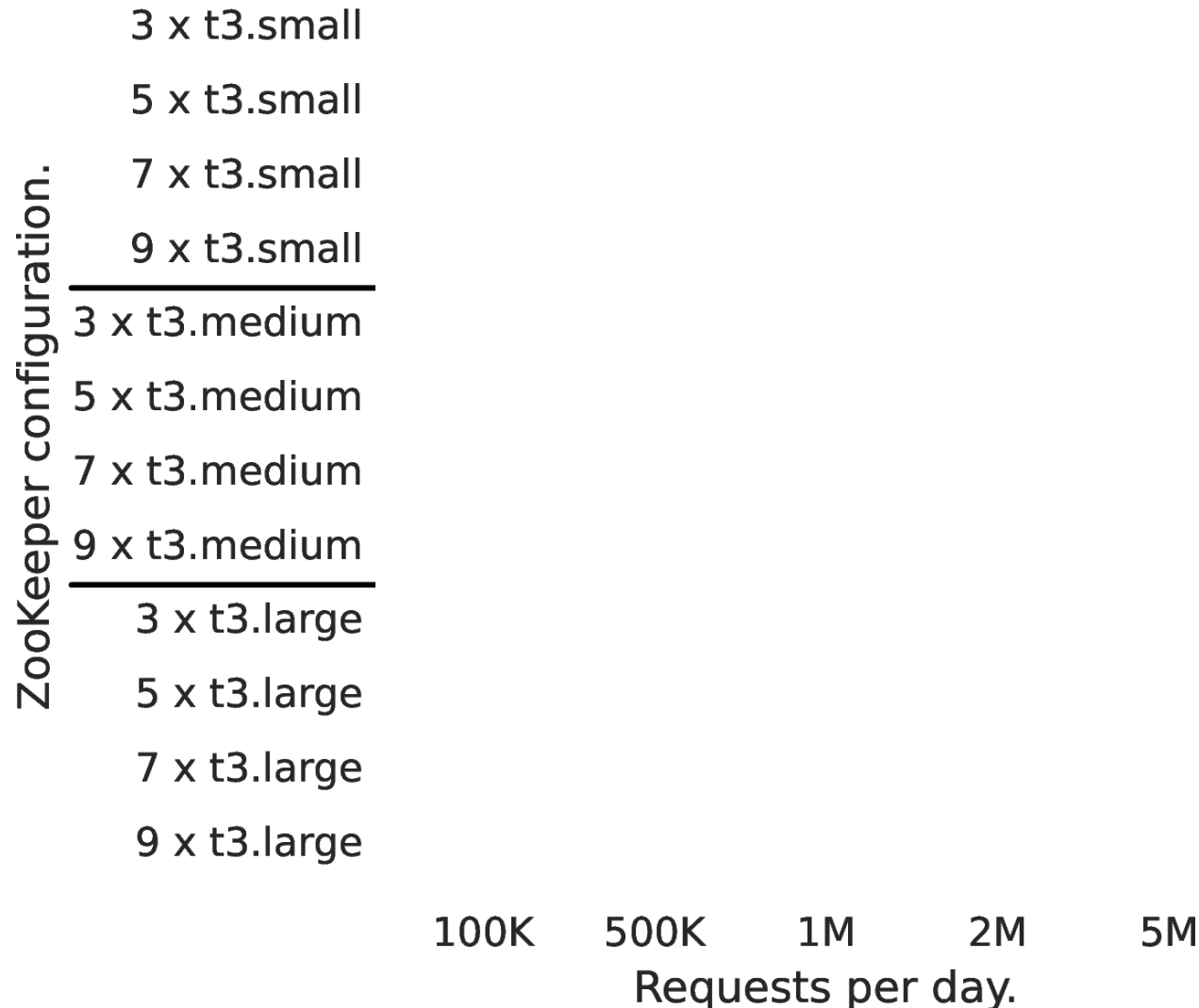


Building Serverless Services: FaaSKeeper



Building Serverless Services: FaaSKeeper

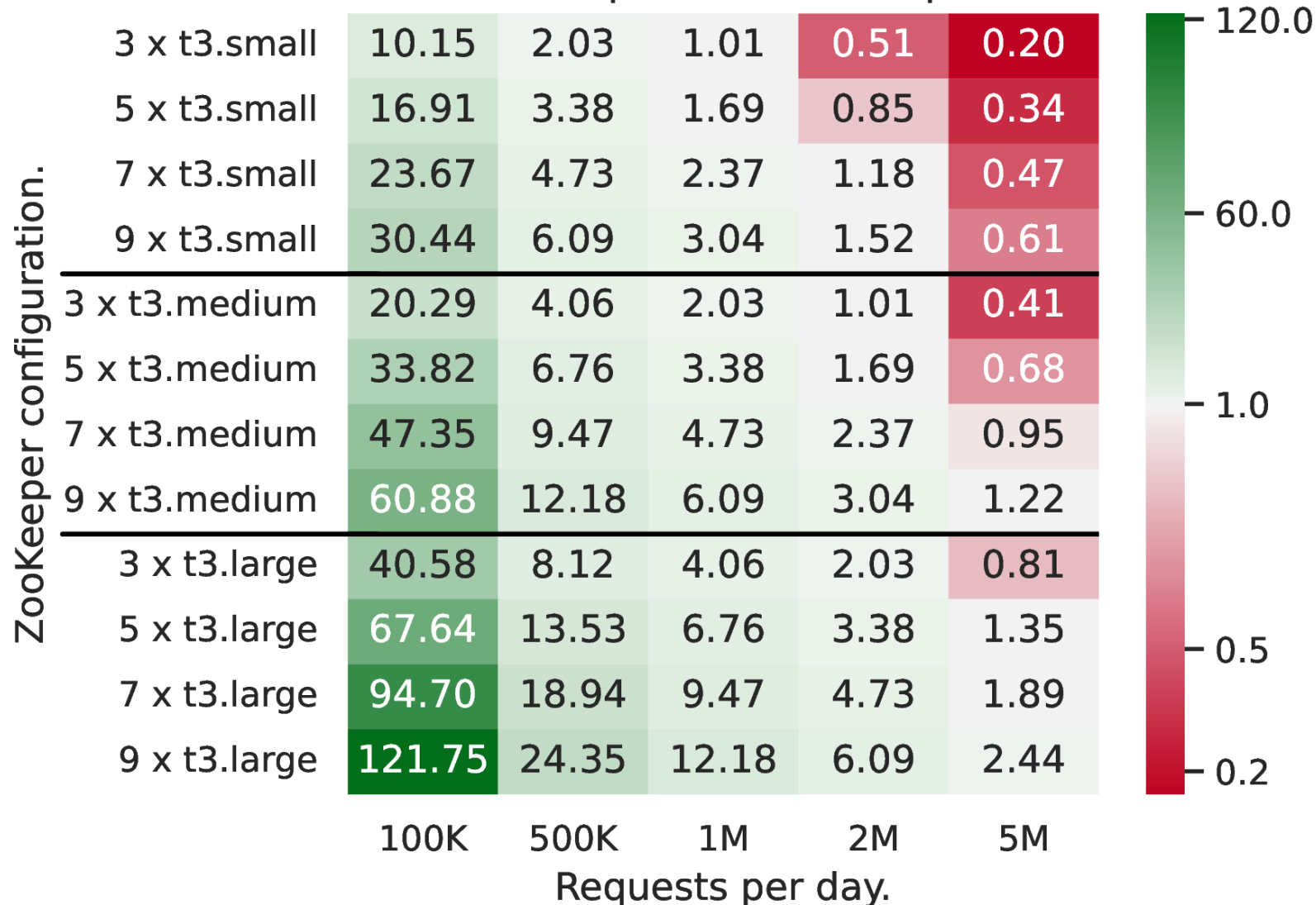
Cost ratio of ZooKeeper and FaaSKeeper, 90% reads.



“FaaSKeeper: Learning from Building Serverless Services with ZooKeeper as an Example”

Building Serverless Services: FaaSKeeper

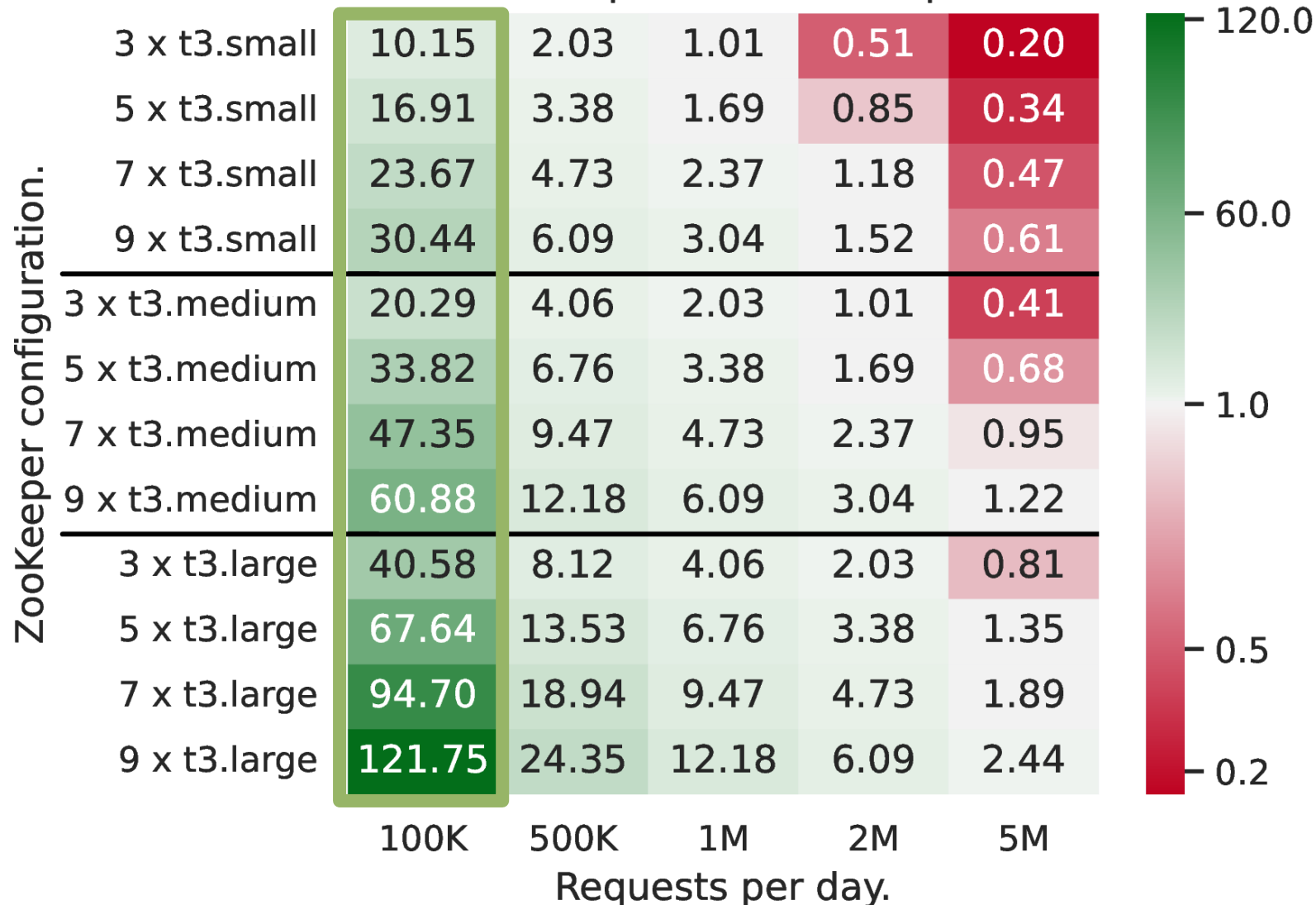
Cost ratio of ZooKeeper and FaaSKeeper, 90% reads.



“FaaSKeeper: Learning from Building Serverless Services with ZooKeeper as an Example”

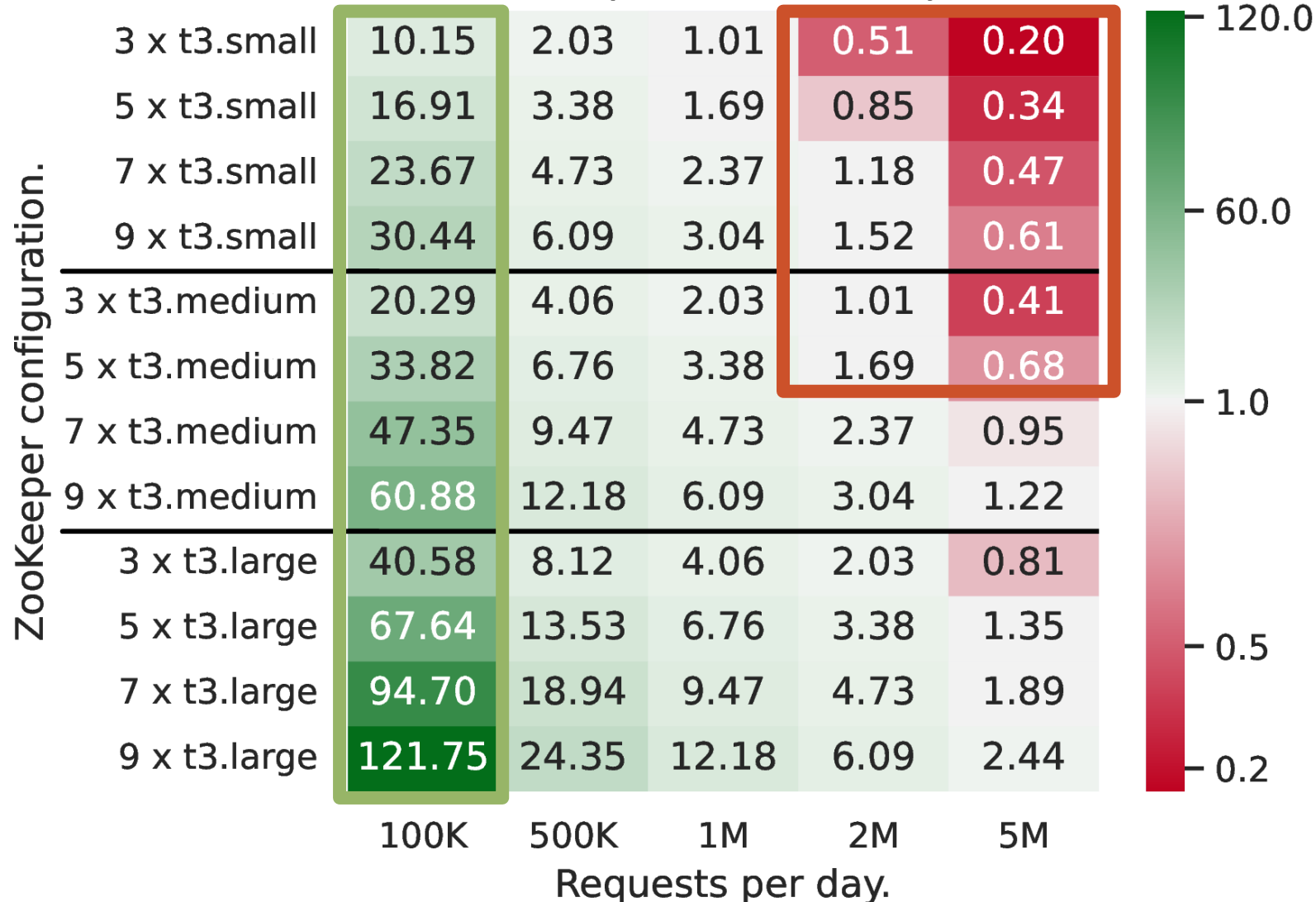
Building Serverless Services: FaaSKeeper

Cost ratio of ZooKeeper and FaaSKeeper, 90% reads.



Building Serverless Services: FaaSKeeper

Cost ratio of ZooKeeper and FaaSKeeper, 90% reads.



Tracking Wasted Money in HPC

Tracking Wasted Money in HPC

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
Argonne National Laboratory

Devesh Tiwari
Northeastern University

SC, 2020

FINAL REPORT WORKLOAD ANALYSIS OF BLUE WATERS (ACI 1650758)

Matthew D. Jones, Joseph P. White, Martins Innus, Robert L. DeLeon, Nikolay Simakov, Jeffrey T. Palmer, Steven M. Gallo, and Thomas R. Furlani (furlani@buffalo.edu), Center for Computational Research, University at Buffalo, SUNY

Michael Showerman, Robert Brunner, Andry Kot, Gregory Bauer, Brett Bode, Jeremy Enos, and William Kramer (wtkramer@illinois.edu), National Center for Supercomputing Applications (NCSA), University of Illinois at Urbana Champaign

arXiv, 2017

Comprehensive Workload Analysis and Modeling of a Petascale Supercomputer

Haihang You¹ and Hao Zhang²

¹ National Institute for Computational Sciences,
Oak Ridge National Laboratory, Oak Ridge, TN 37831, USA
² Department of Electrical Engineering and Computer Science,
University of Tennessee, Knoxville, TN 37996, USA
{hyou, haozhang}@utk.edu

JSSPP, 2012

HPC System Utilization - CPU

HPC System Utilization - CPU



Piz Daint, April 2022.

- XC50 nodes – CPU + GPU, 64 GB memory.
- XC40 nodes – CPU, 64/128 GB memory.

Query SLURM info every two minutes.

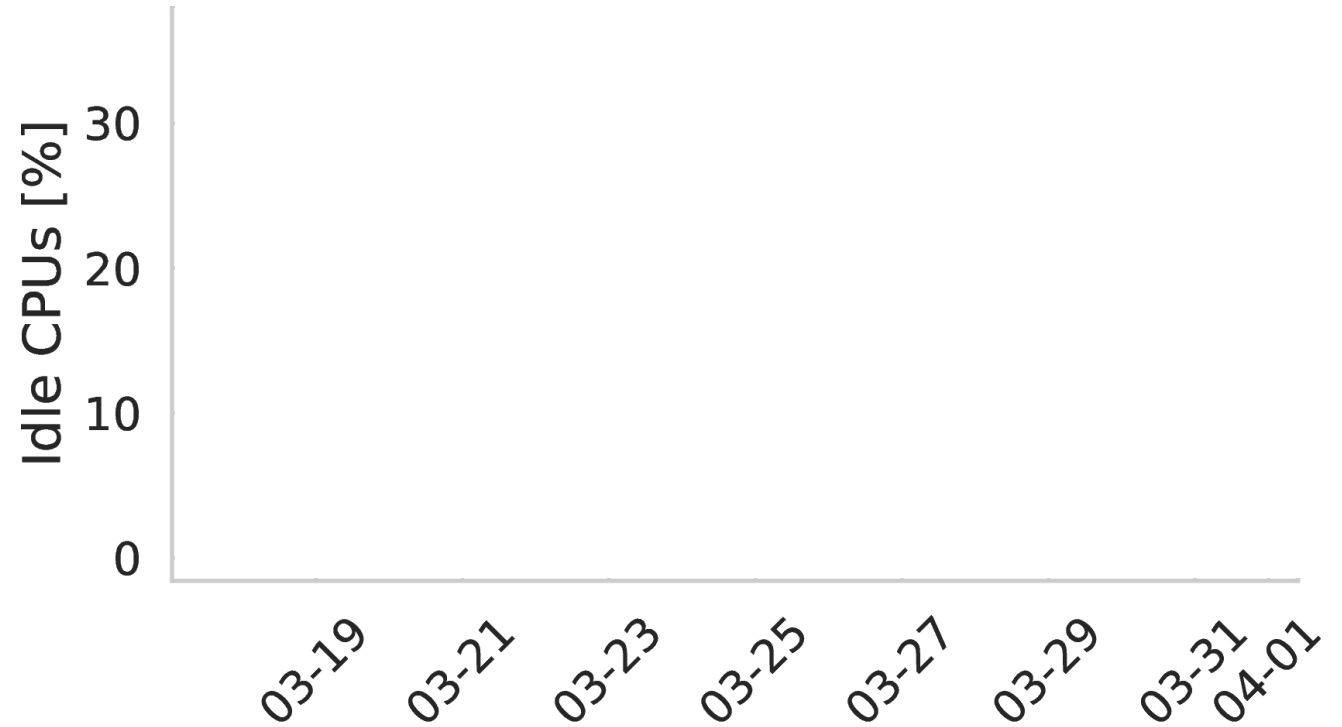
HPC System Utilization - CPU



Piz Daint, April 2022.

- XC50 nodes – CPU + GPU, 64 GB memory.
- XC40 nodes – CPU, 64/128 GB memory.

Query SLURM info every two minutes.



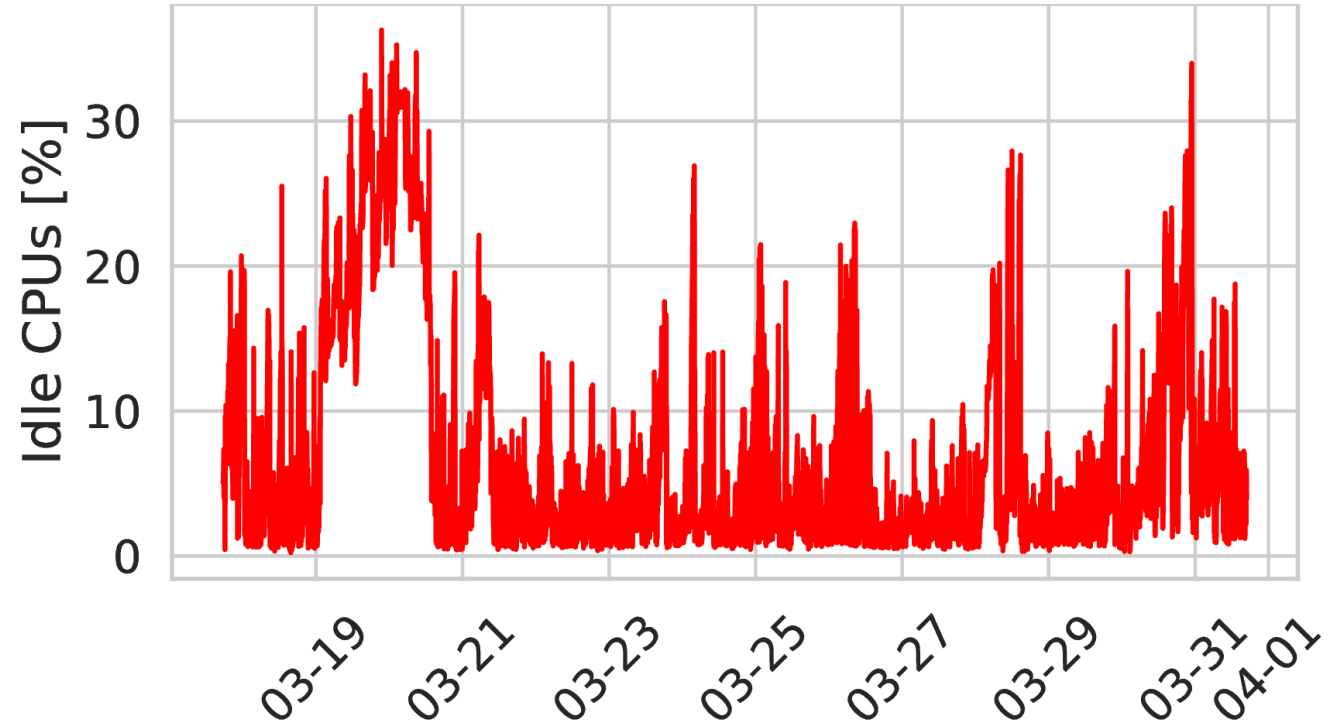
HPC System Utilization - CPU



Piz Daint, April 2022.

- XC50 nodes – CPU + GPU, 64 GB memory.
- XC40 nodes – CPU, 64/128 GB memory.

Query SLURM info every two minutes.



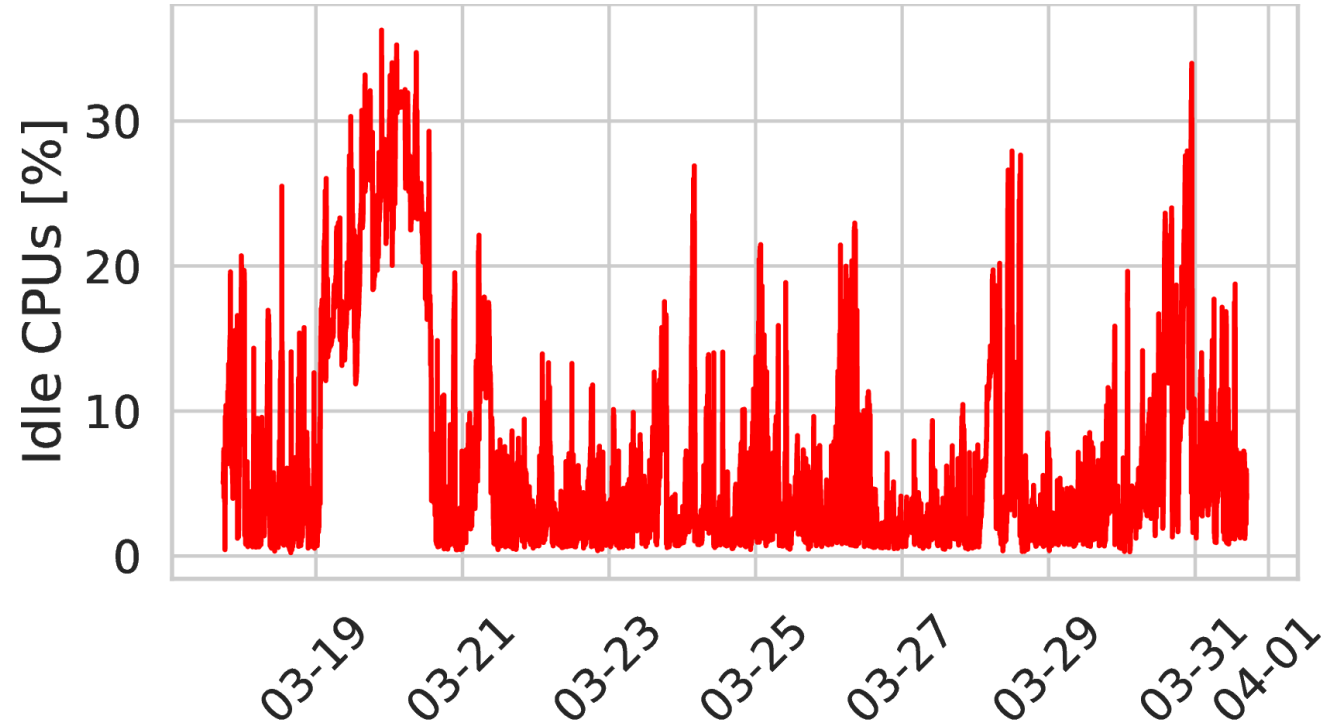
HPC System Utilization - CPU



Piz Daint, April 2022.

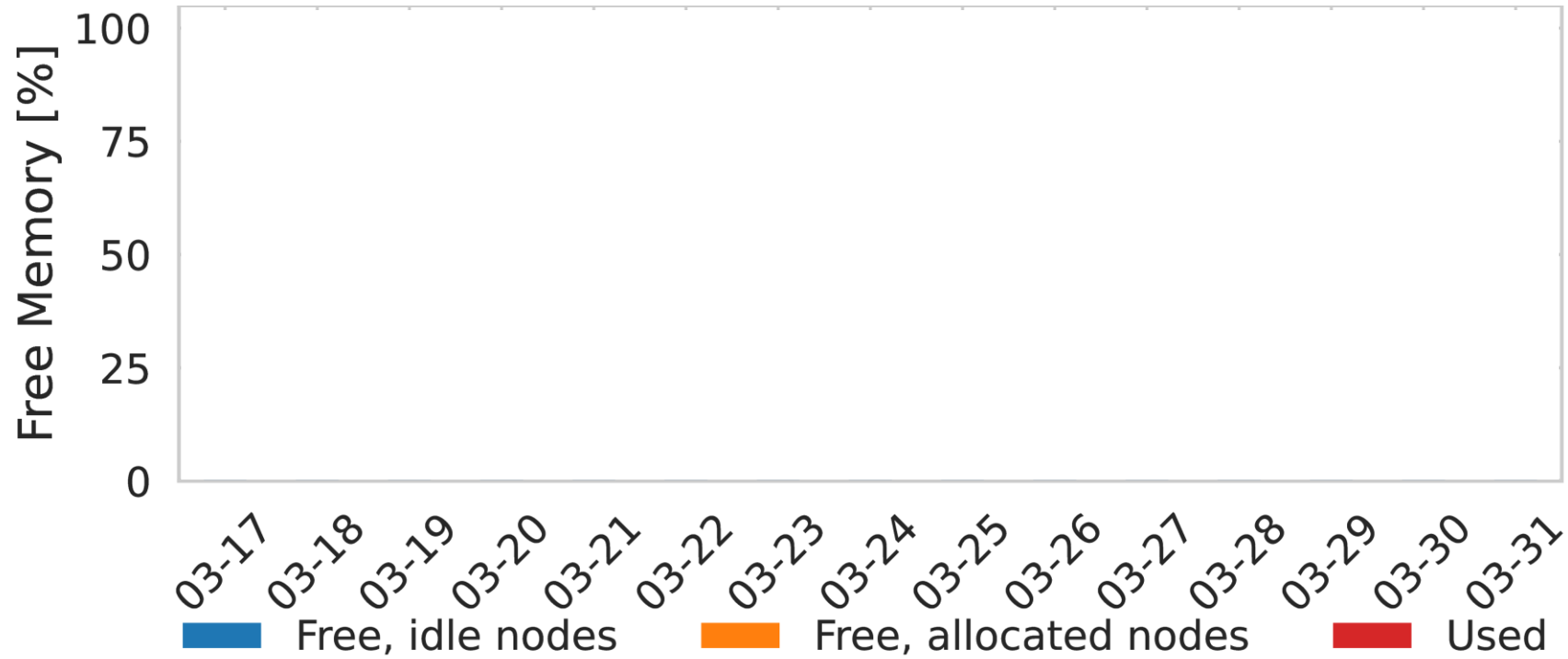
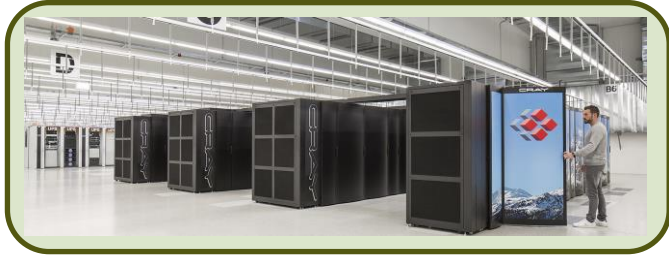
- XC50 nodes – CPU + GPU, 64 GB memory.
- XC40 nodes – CPU, 64/128 GB memory.

Query SLURM info every two minutes.

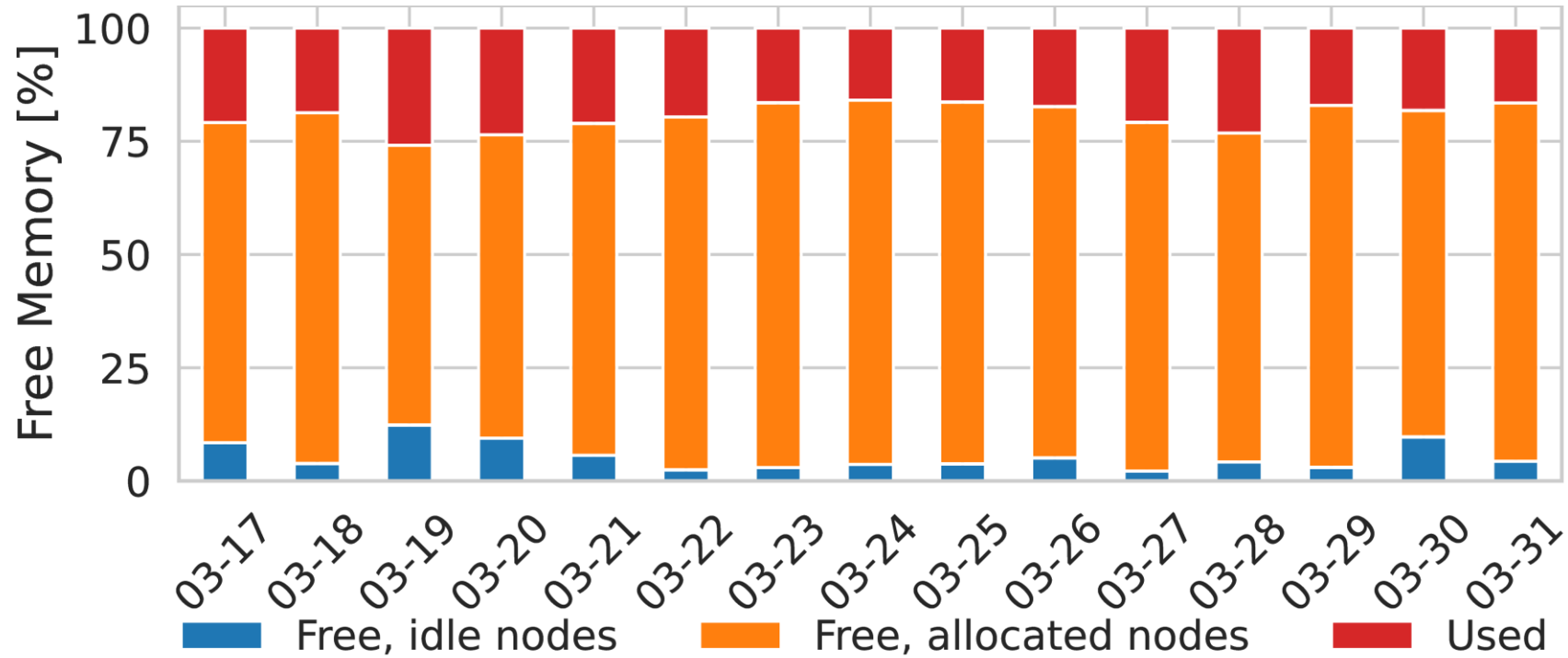
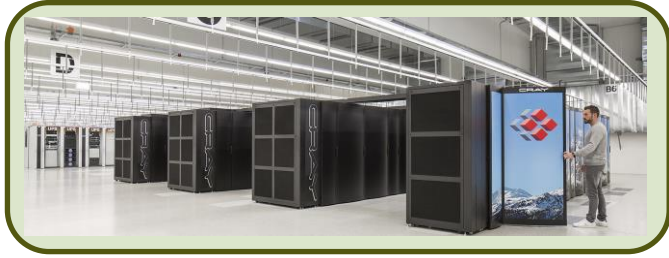


Nodes do not stay idle for long.

HPC System Utilization - Memory



HPC System Utilization - Memory



HPC System Utilization - Memory



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

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	

MICRO, 2019

FINAL REPORT WORKLOAD ANALYSIS OF BLUE WATERS (ACI 1650758)

Matthew D. Jones, Joseph P. White, Martins Innus, Robert L. DeLeon, Nikolay Simakov, Jeffrey T. Palmer, Steven M. Gallo, and Thomas R. Furlani (furlani@buffalo.edu), Center for Computational Research, University at Buffalo, SUNY

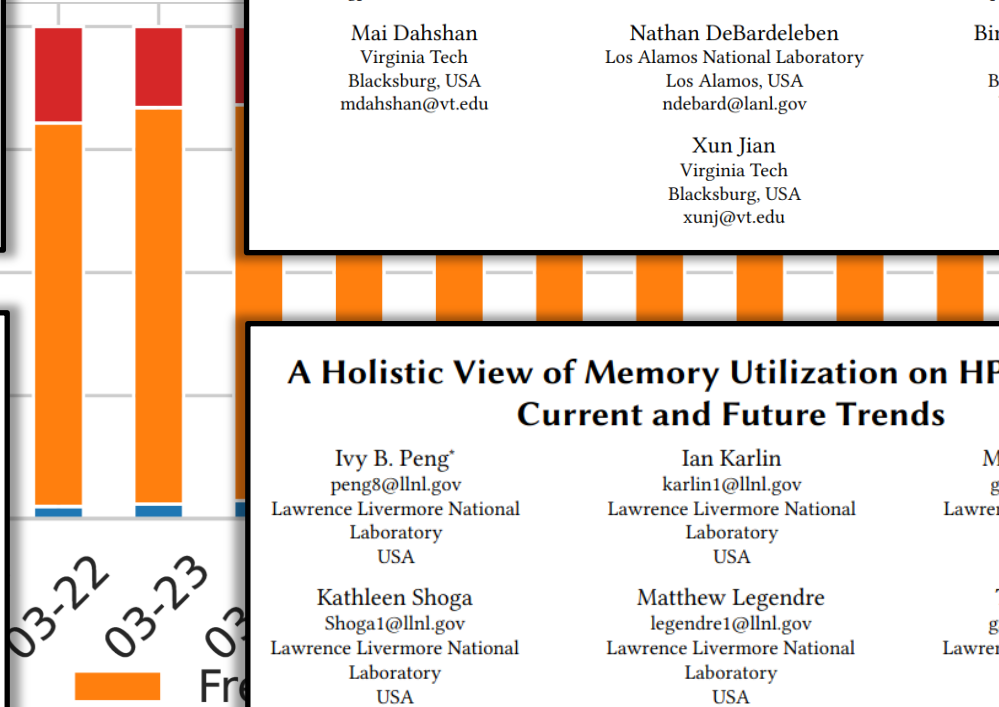
Michael Showerman, Robert Brunner, Andry Kot, Gregory Bauer, Brett Bode, Jeremy Enos, and William Kramer (wtkramer@illinois.edu), National Center for Supercomputing Applications (NCSA), University of Illinois at Urbana Champaign

arXiv, 2017

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

MEMSYS, 2021



FaaS in High-Performance Applications

FaaS in High-Performance Applications

Serverless is slow

FaaS in High-Performance Applications

Serverless is slow

Communication is slow
and restricted

FaaS in High-Performance Applications

Serverless is slow

Communication is slow
and restricted

Serverless is hard to
program.

FaaS in High-Performance Applications

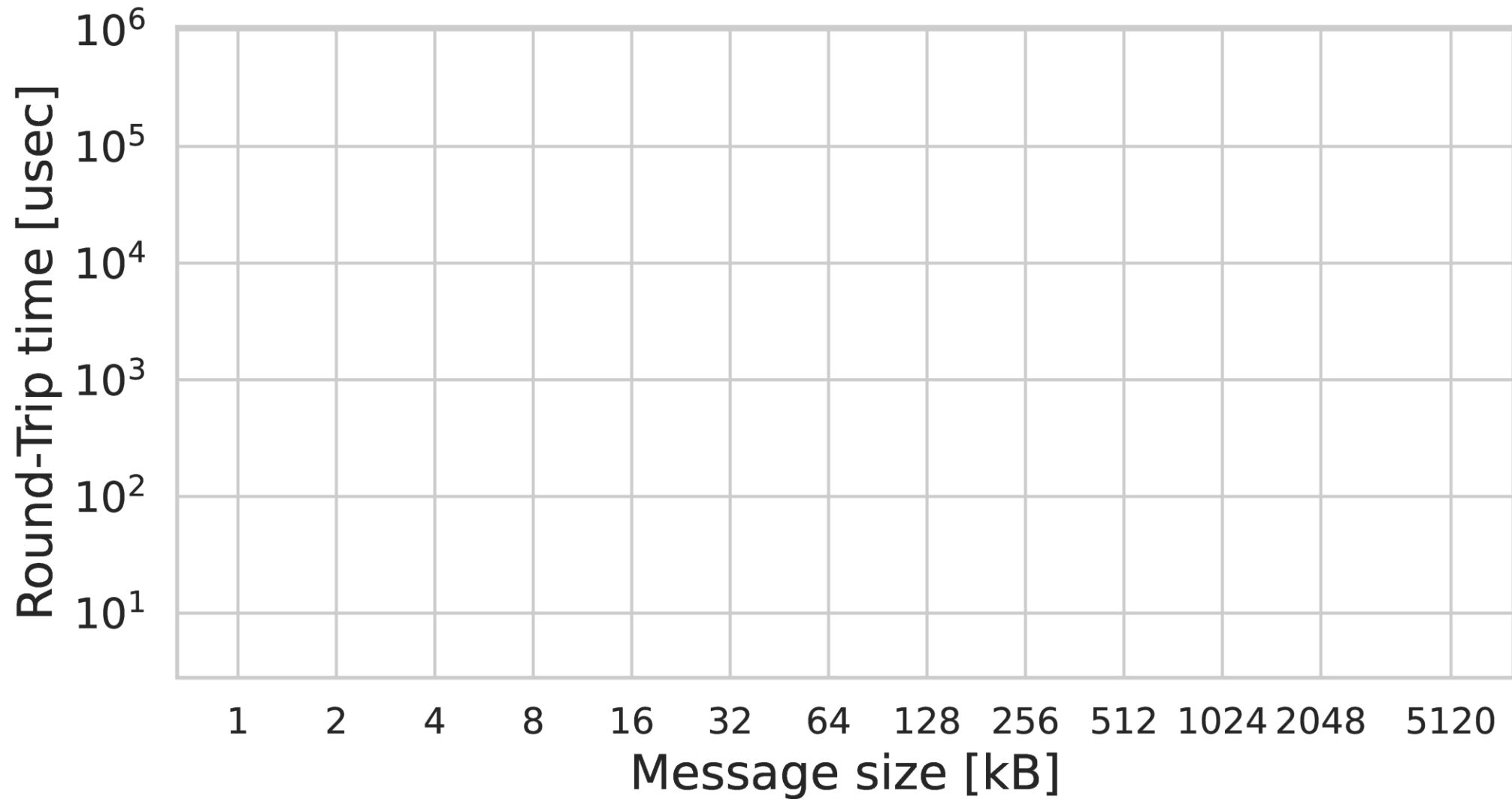
Serverless is slow

Communication is slow
and restricted

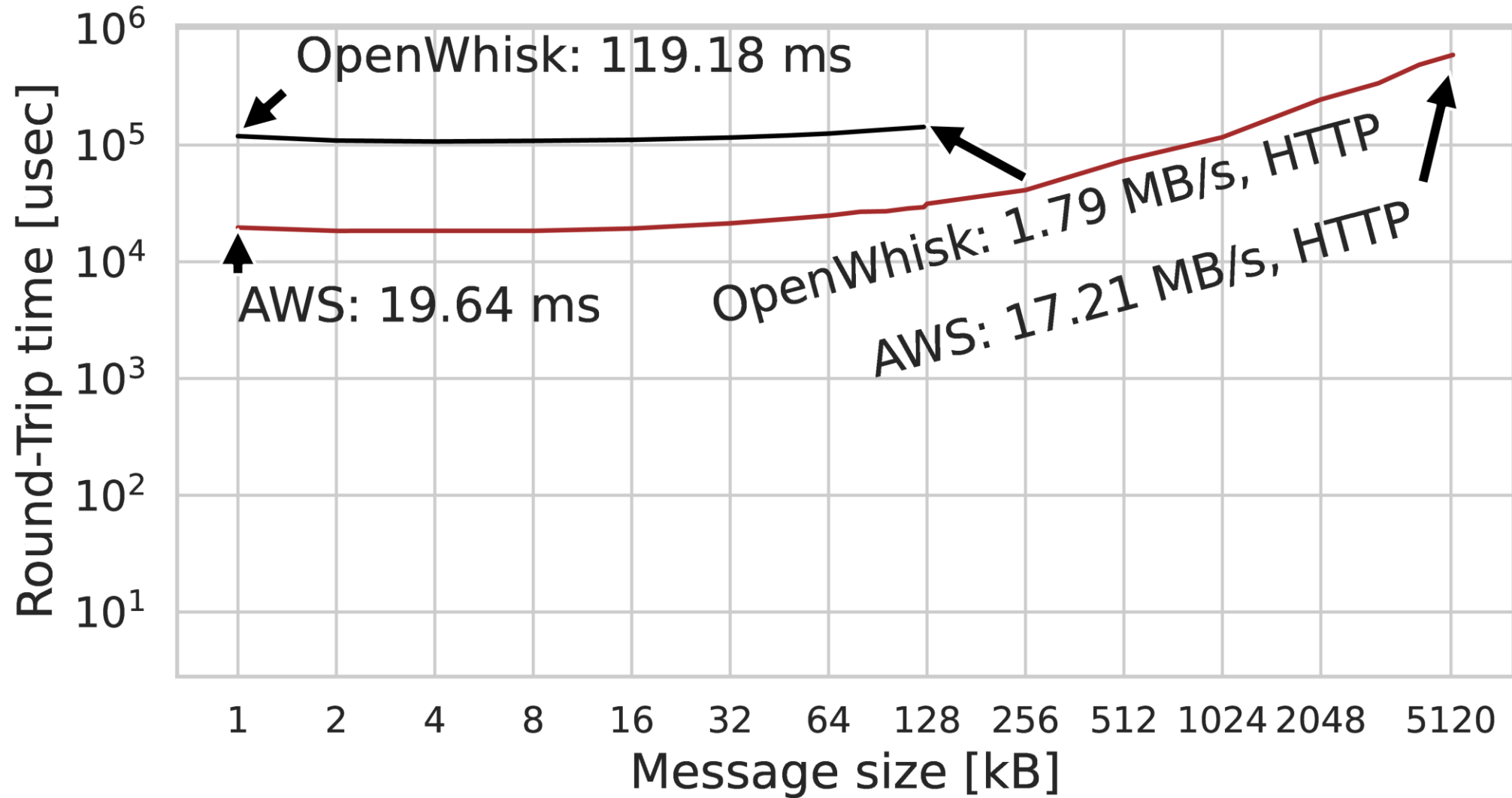
Answer:
rFaaS

Serverless is hard to
program.

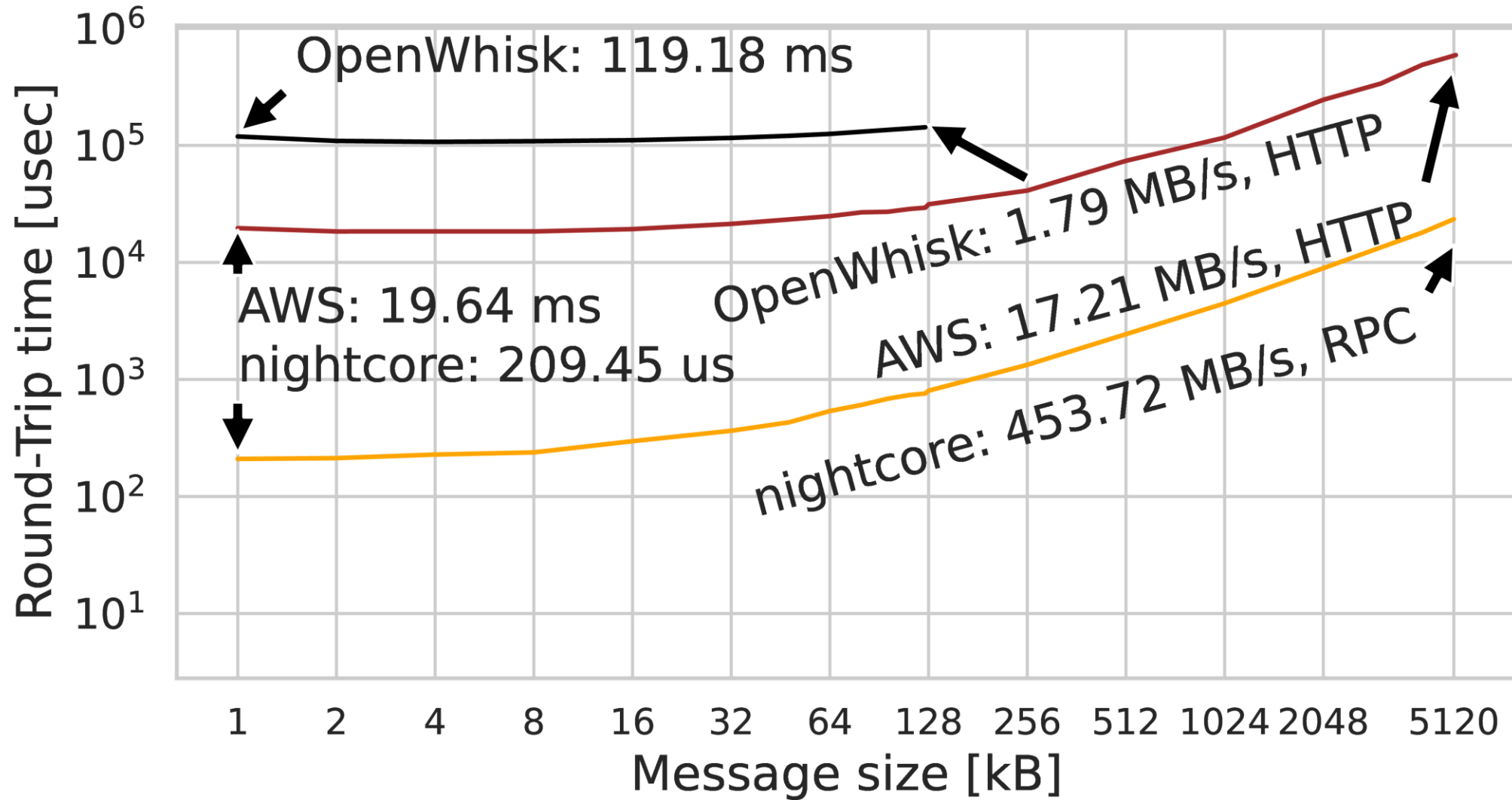
How fast are invocations in FaaS?



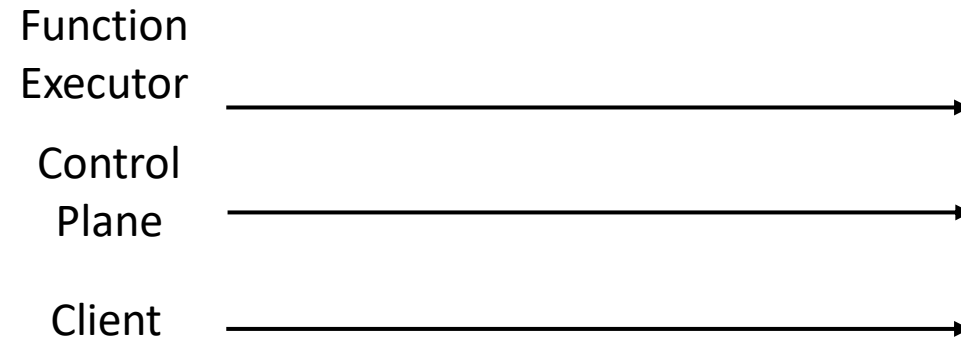
How fast are invocations in FaaS?



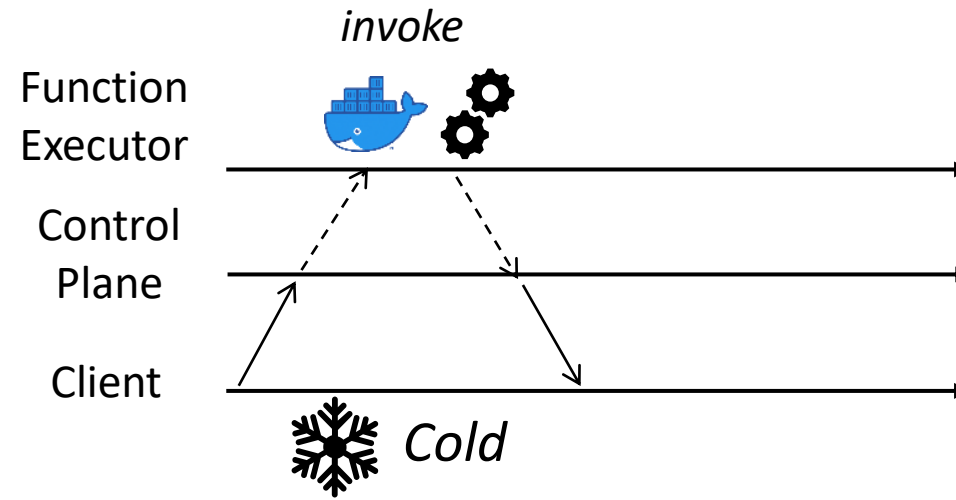
How fast are invocations in FaaS?



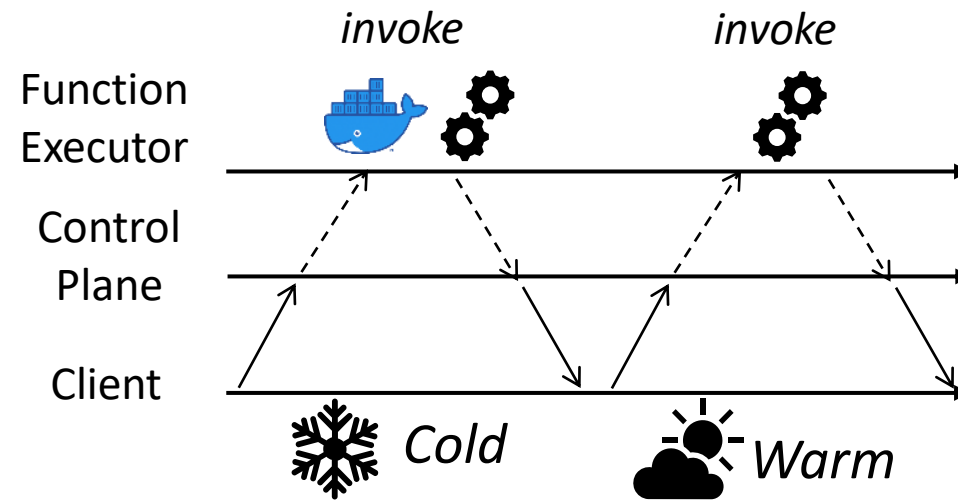
Invocations in FaaS and rFaaS



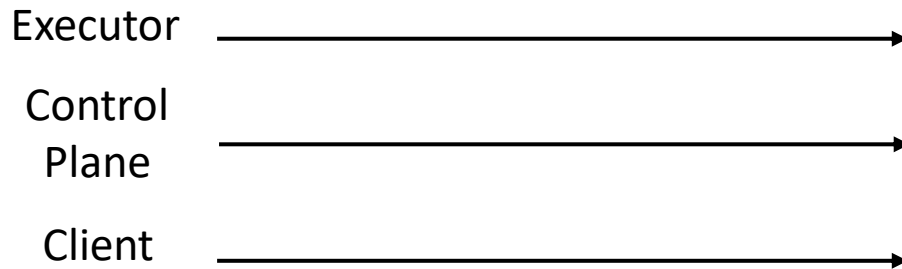
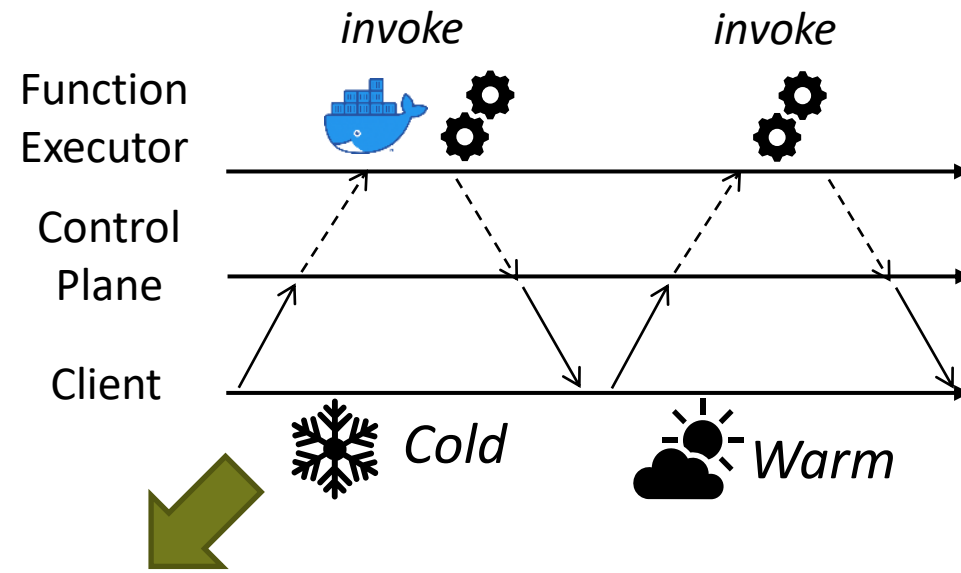
Invocations in FaaS and rFaaS



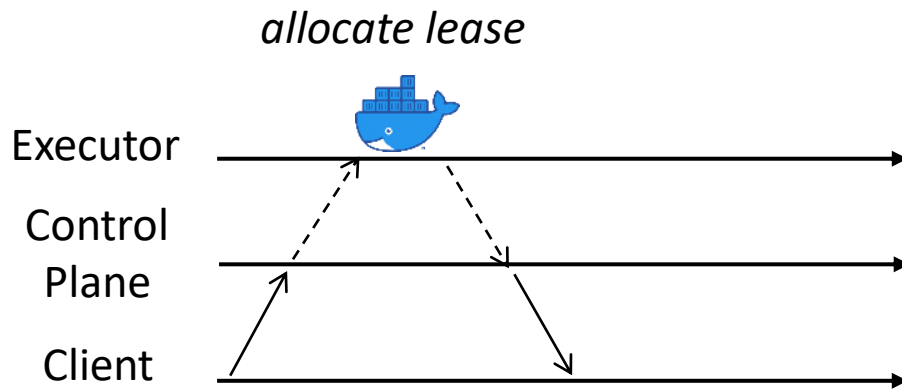
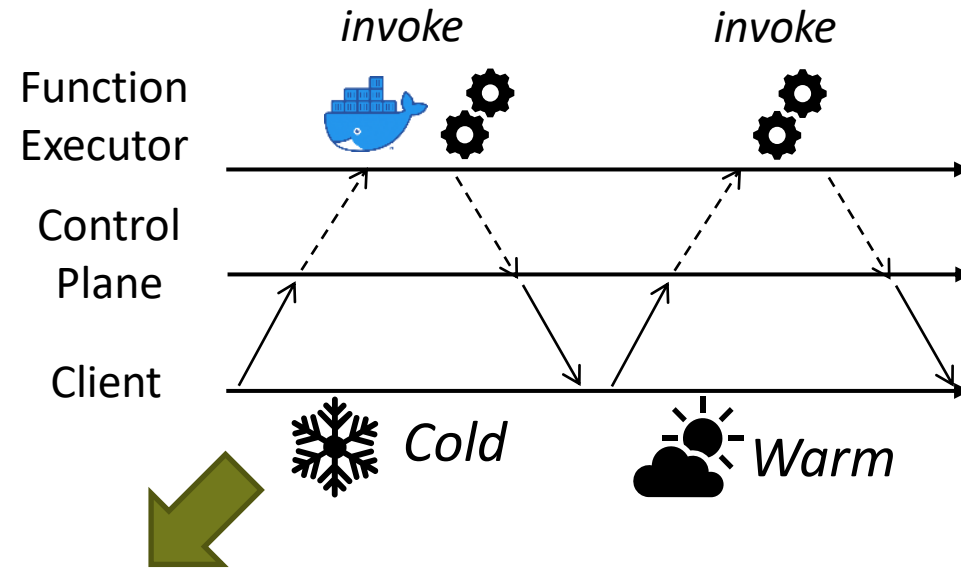
Invocations in FaaS and rFaaS



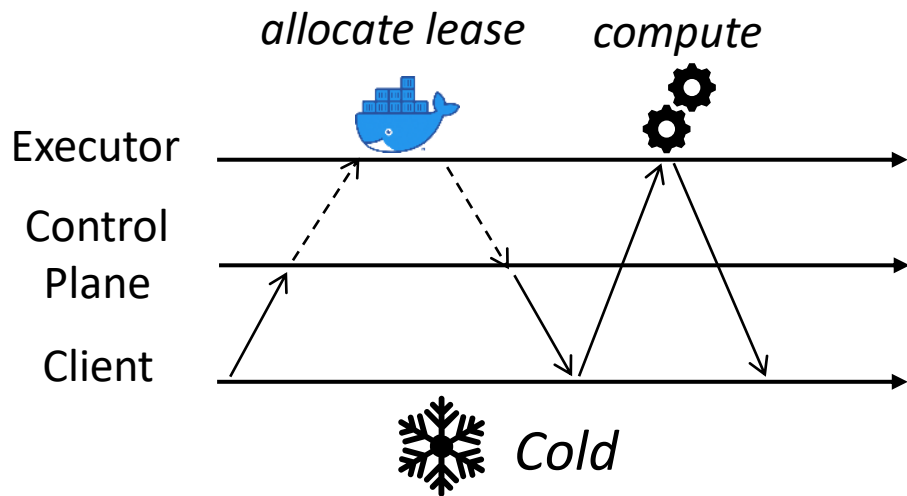
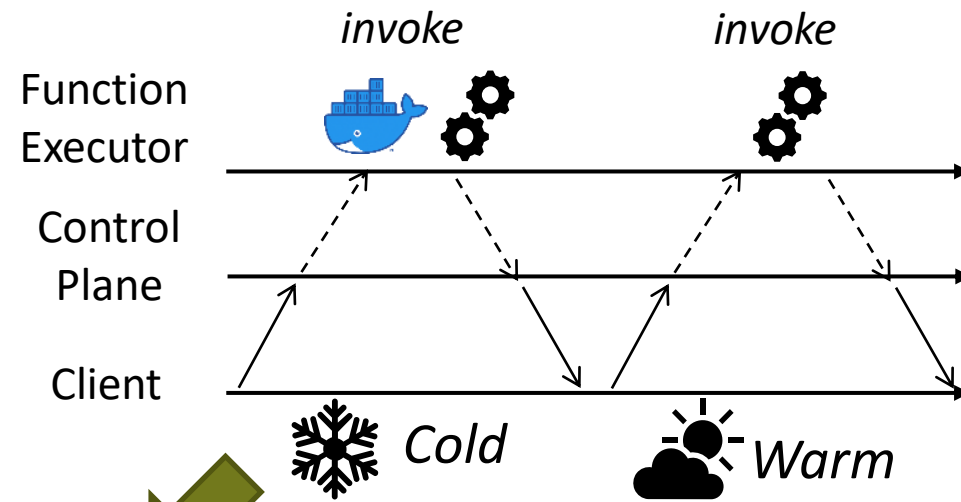
Invocations in FaaS and rFaaS



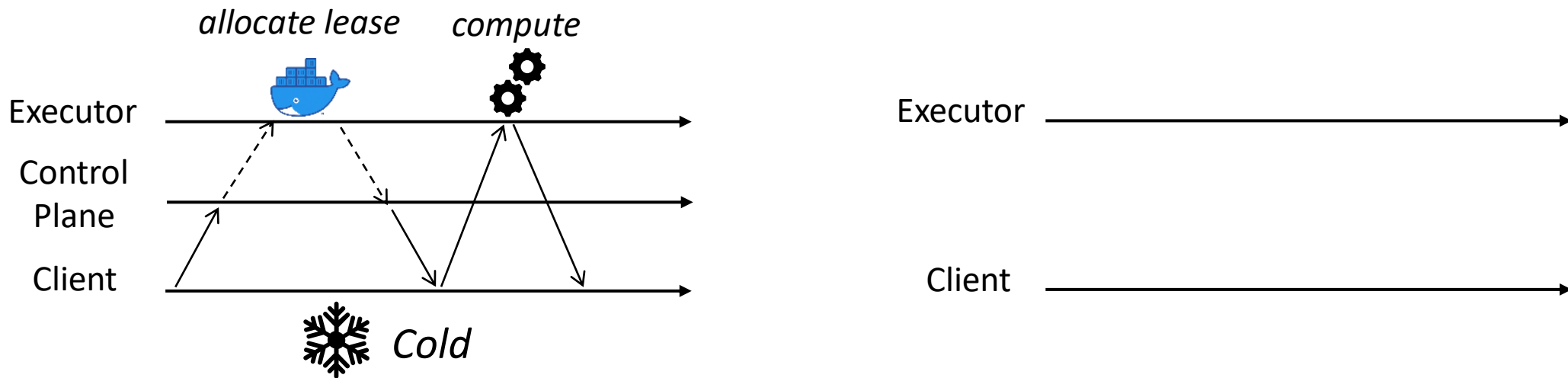
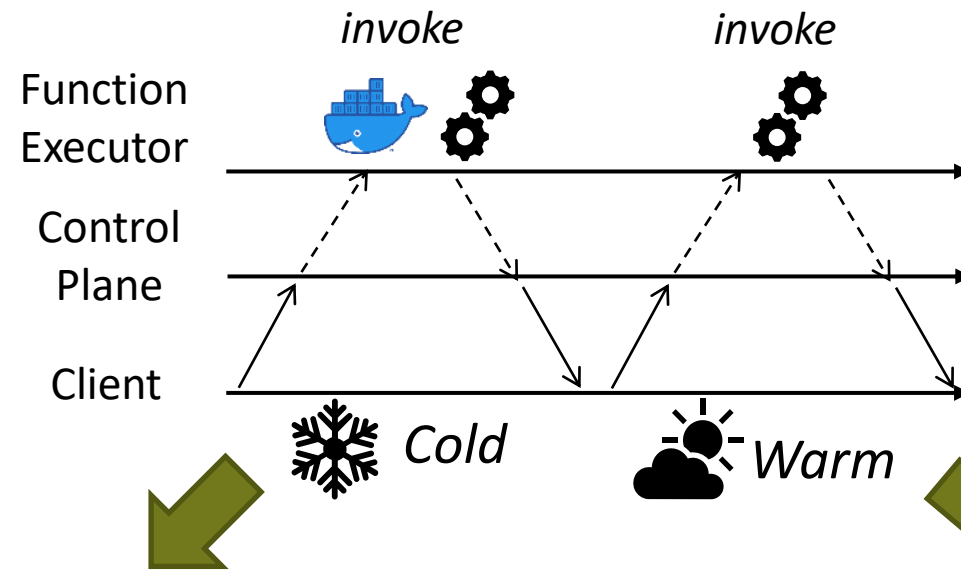
Invocations in FaaS and rFaaS



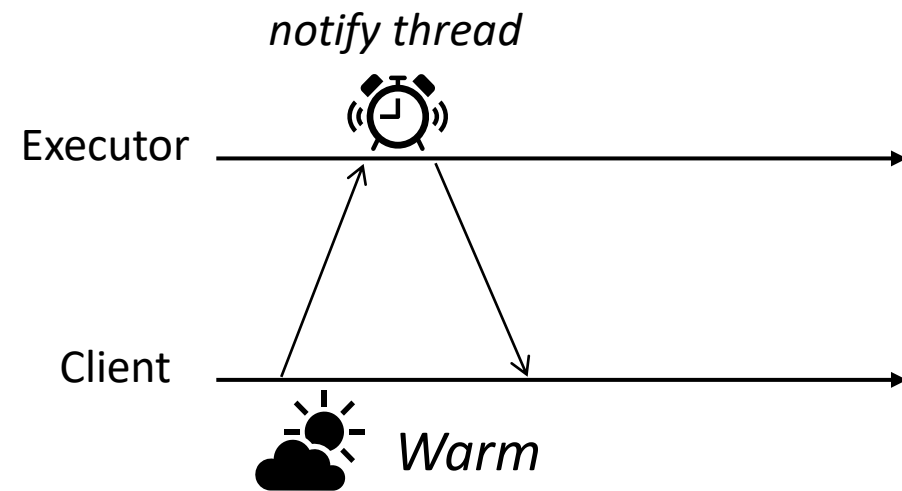
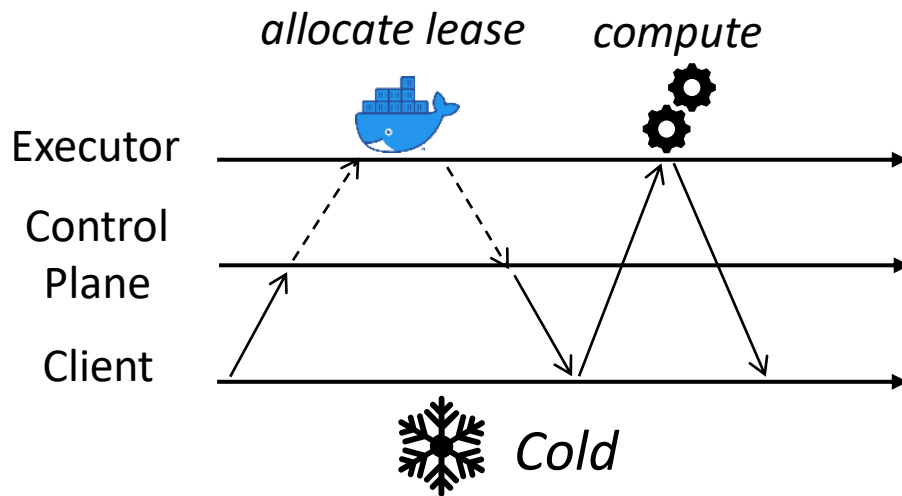
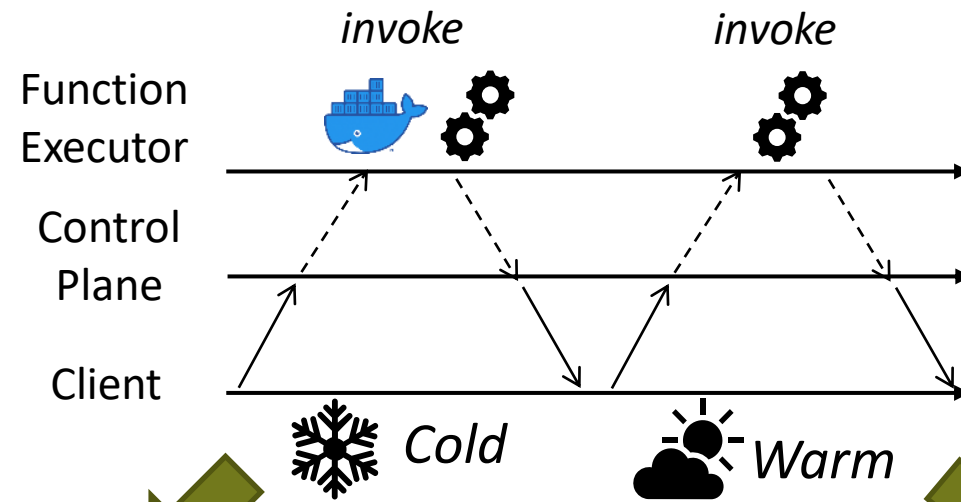
Invocations in FaaS and rFaaS



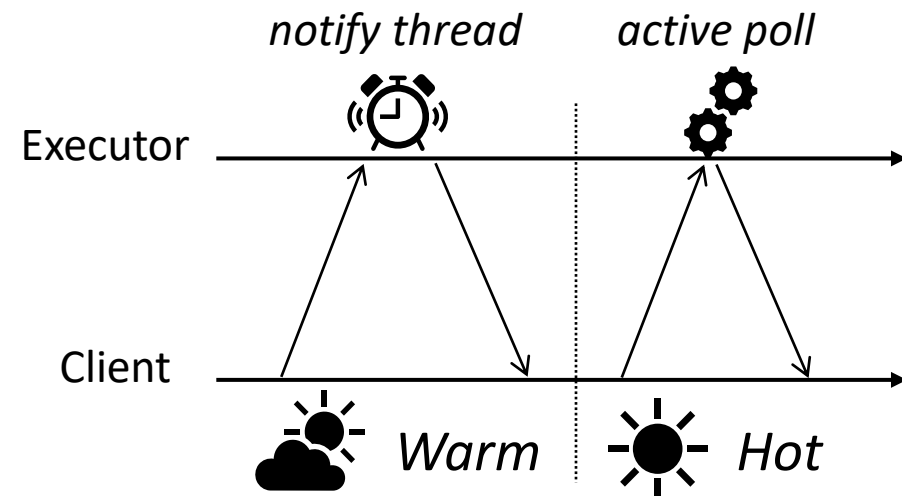
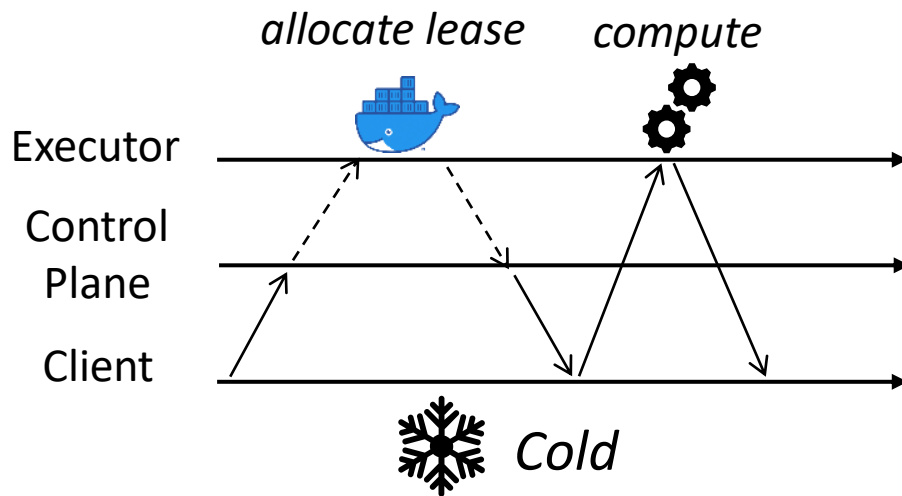
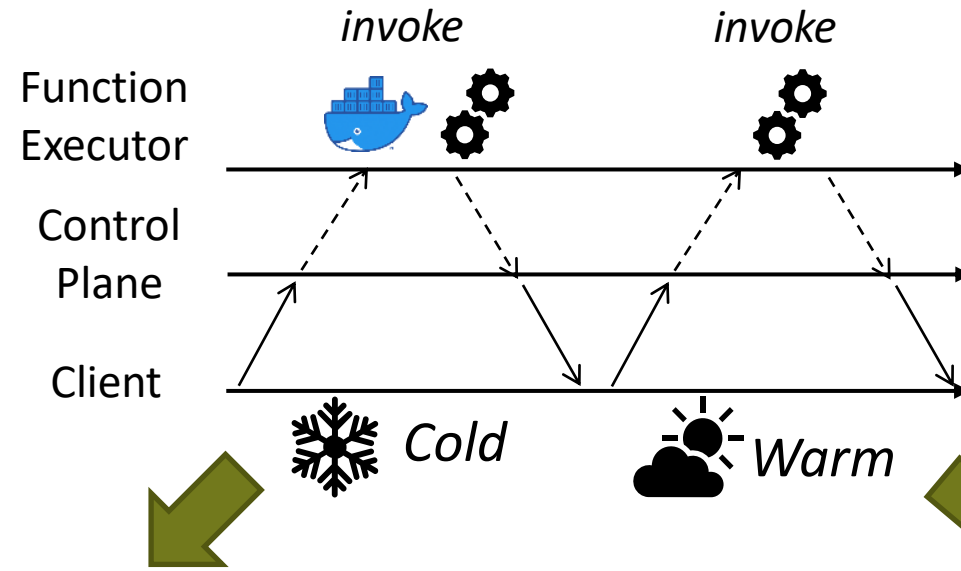
Invocations in FaaS and rFaaS



Invocations in FaaS and rFaaS

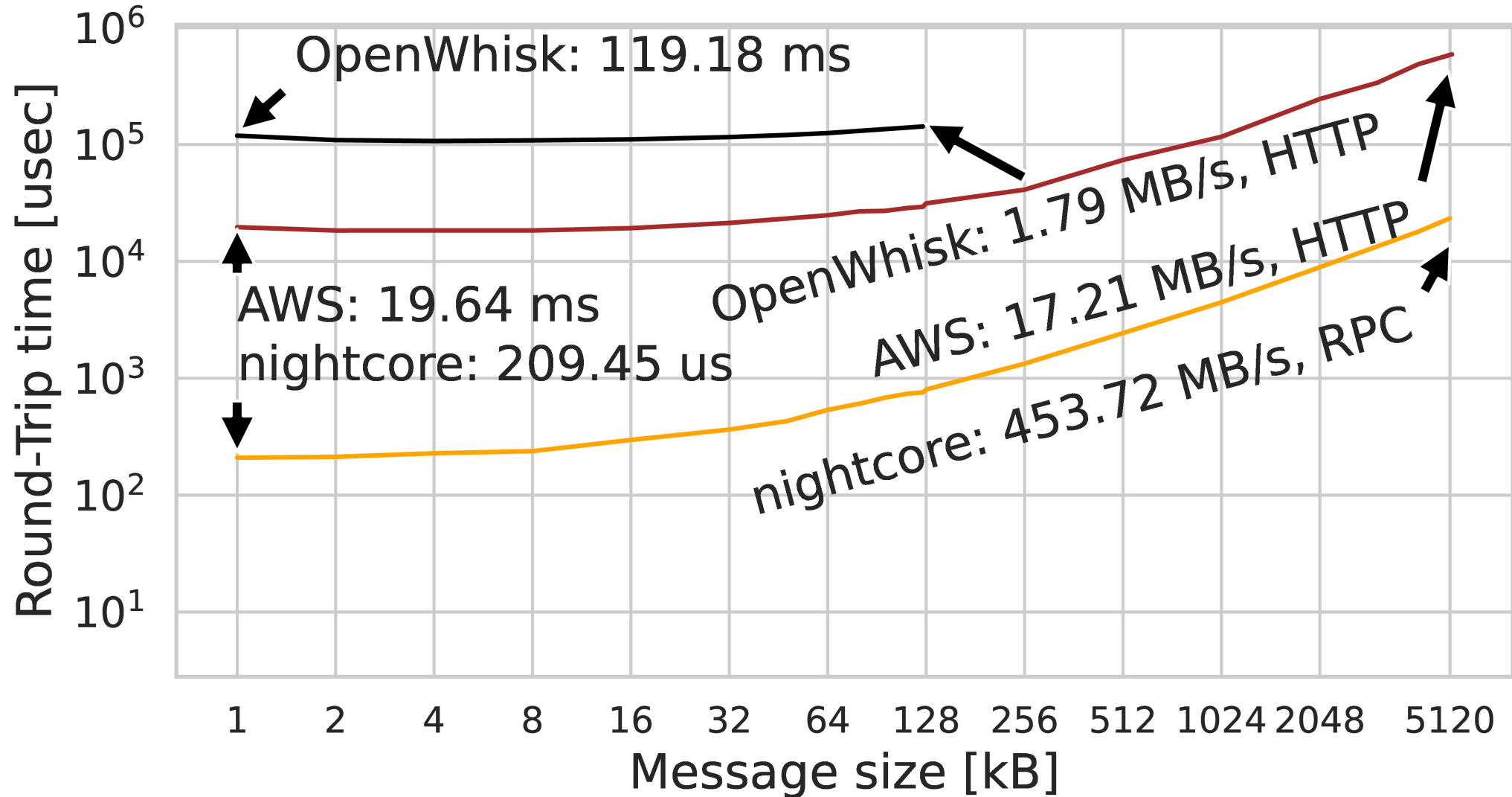


Invocations in FaaS and rFaaS



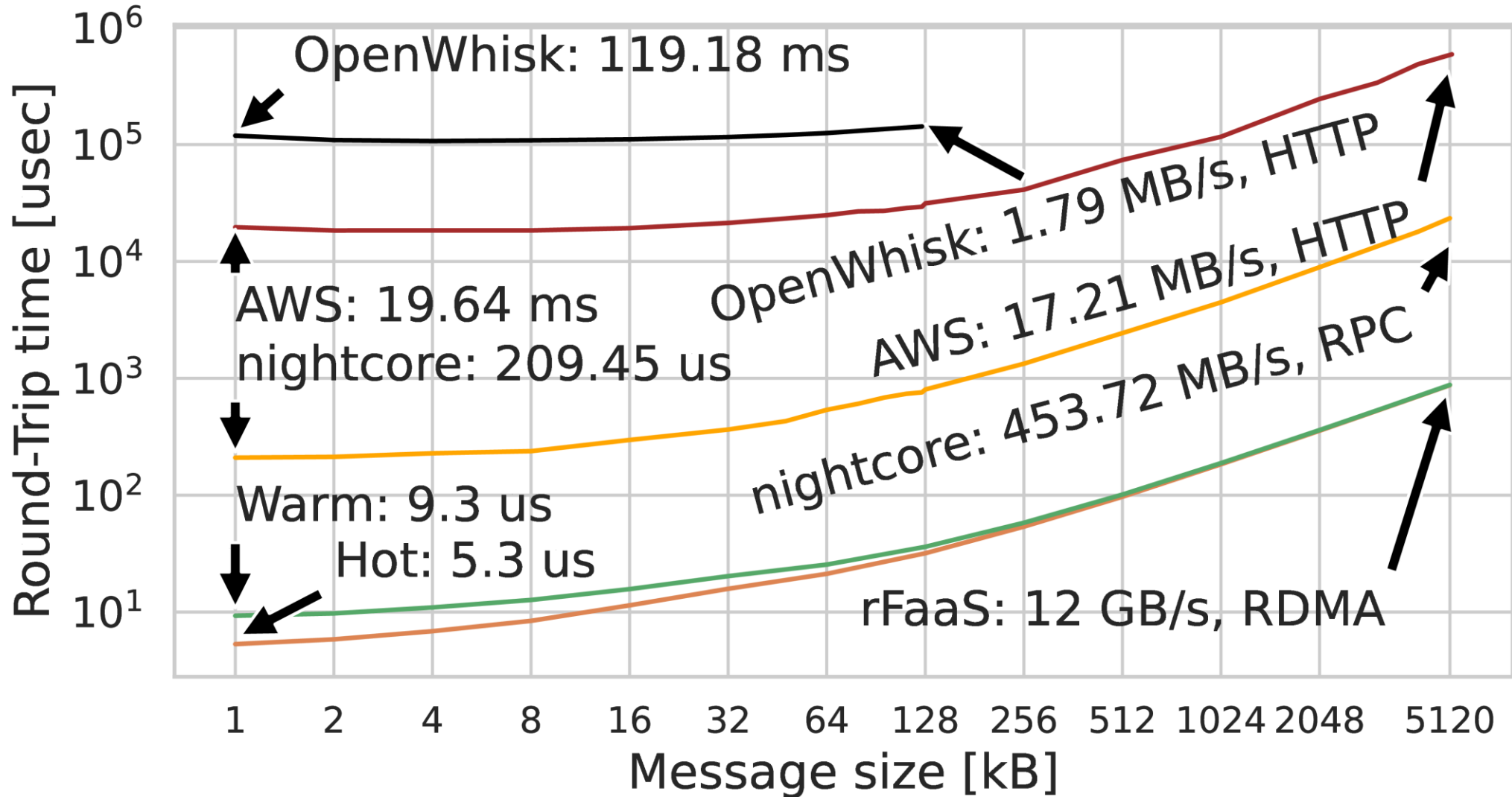
How fast are invocations in rFaaS?

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



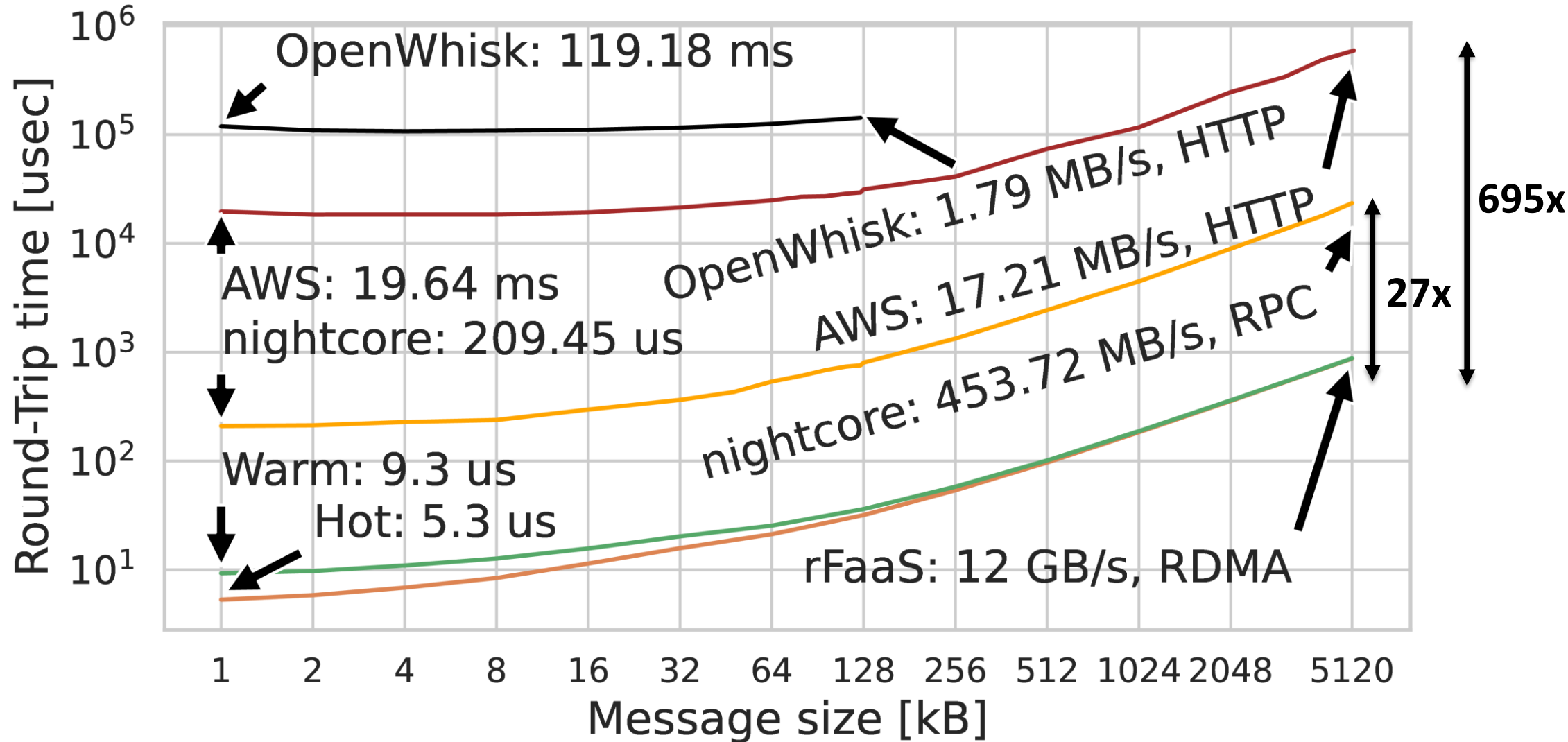
How fast are invocations in rFaaS?

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



How fast are invocations in rFaaS?

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



FaaS in High-Performance Applications

Serverless is slow

Communication is slow
and restricted

Answer:
rFaaS

Serverless is hard to
program.

FaaS in High-Performance Applications

Serverless is slow

Answer:
rFaaS

Communication is slow
and restricted

Answer:
FMI

Serverless is hard to
program.

Communication in serverless



“FMI: Fast and Cheap Message Passing for Serverless Functions”, ICS’23

Communication in serverless

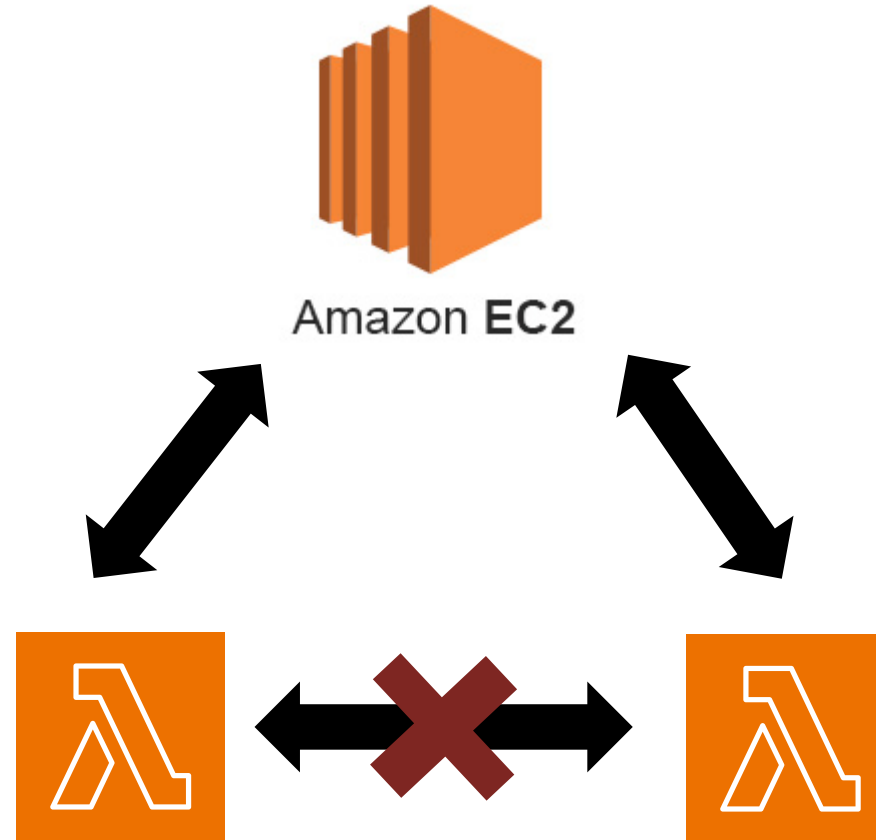


Communication in serverless



“FMI: Fast and Cheap Message Passing for Serverless Functions”, ICS’23

Communication in serverless

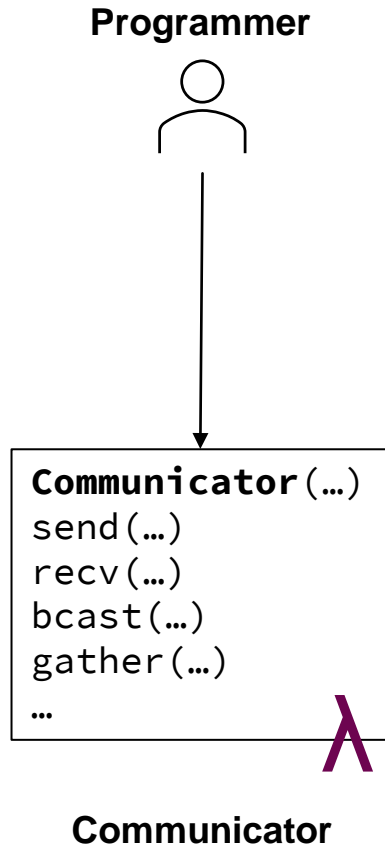


FMI: MPI for serverless

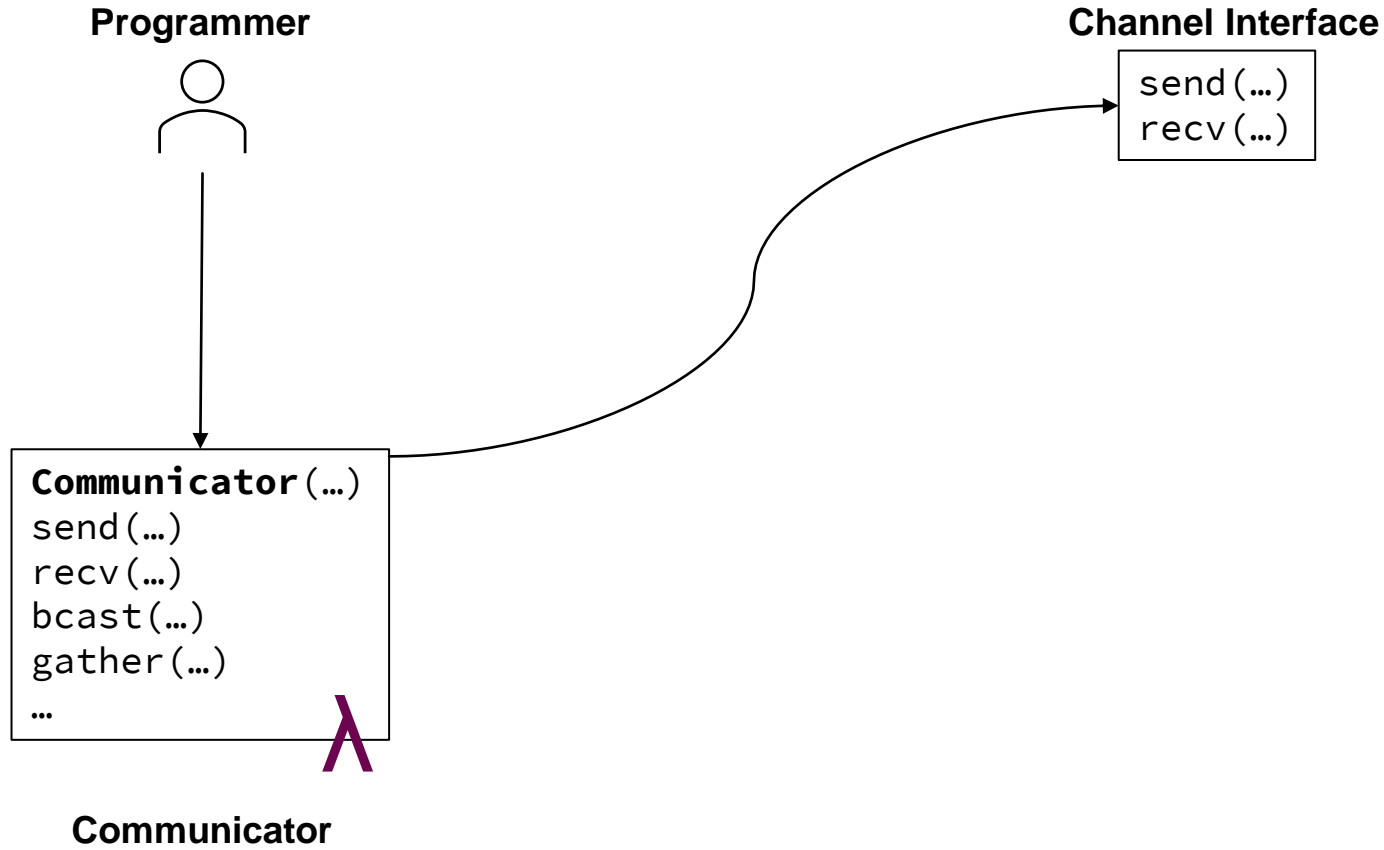


“FMI: Fast and Cheap Message Passing for Serverless Functions”, ICS’23

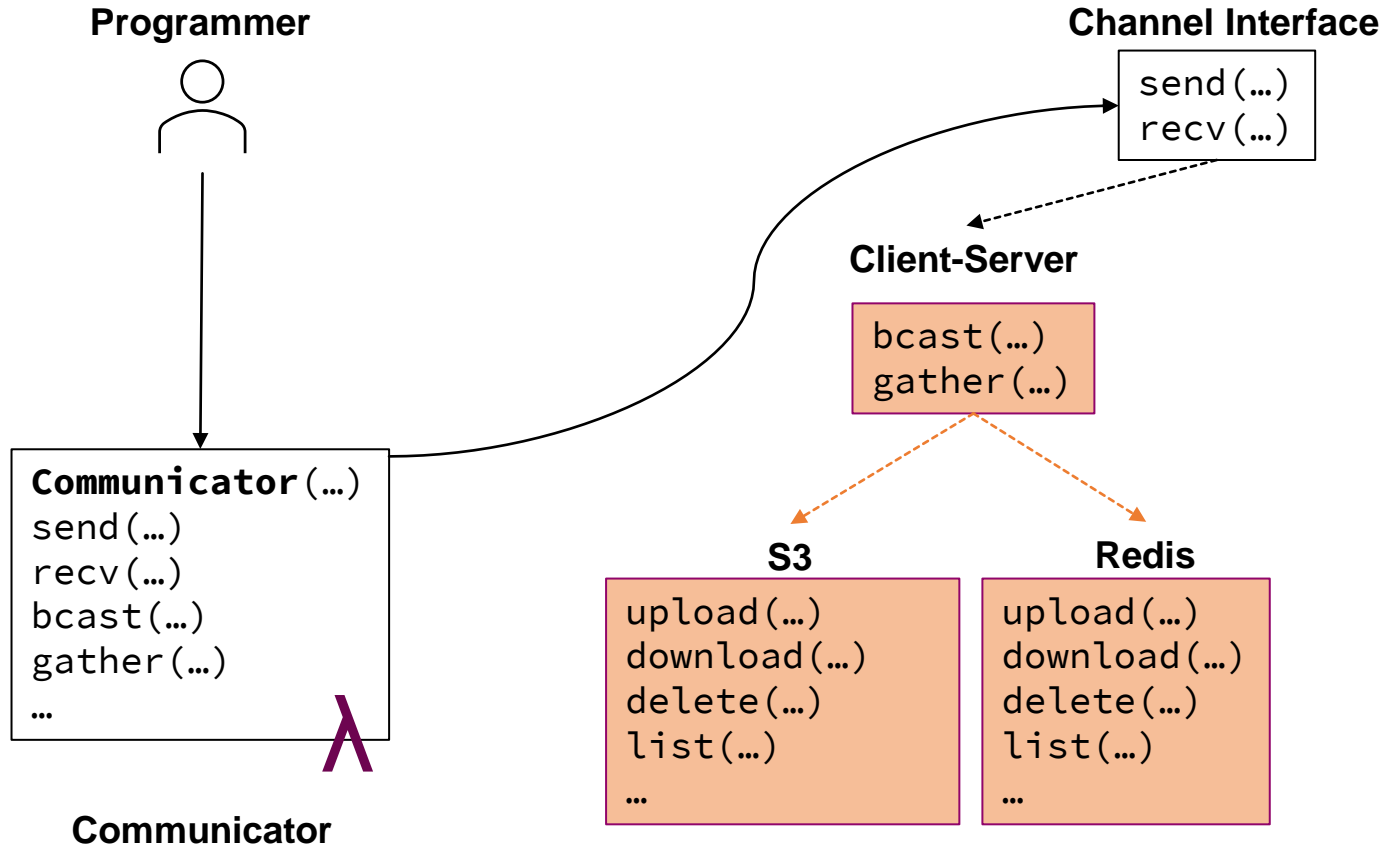
FMI: MPI for serverless



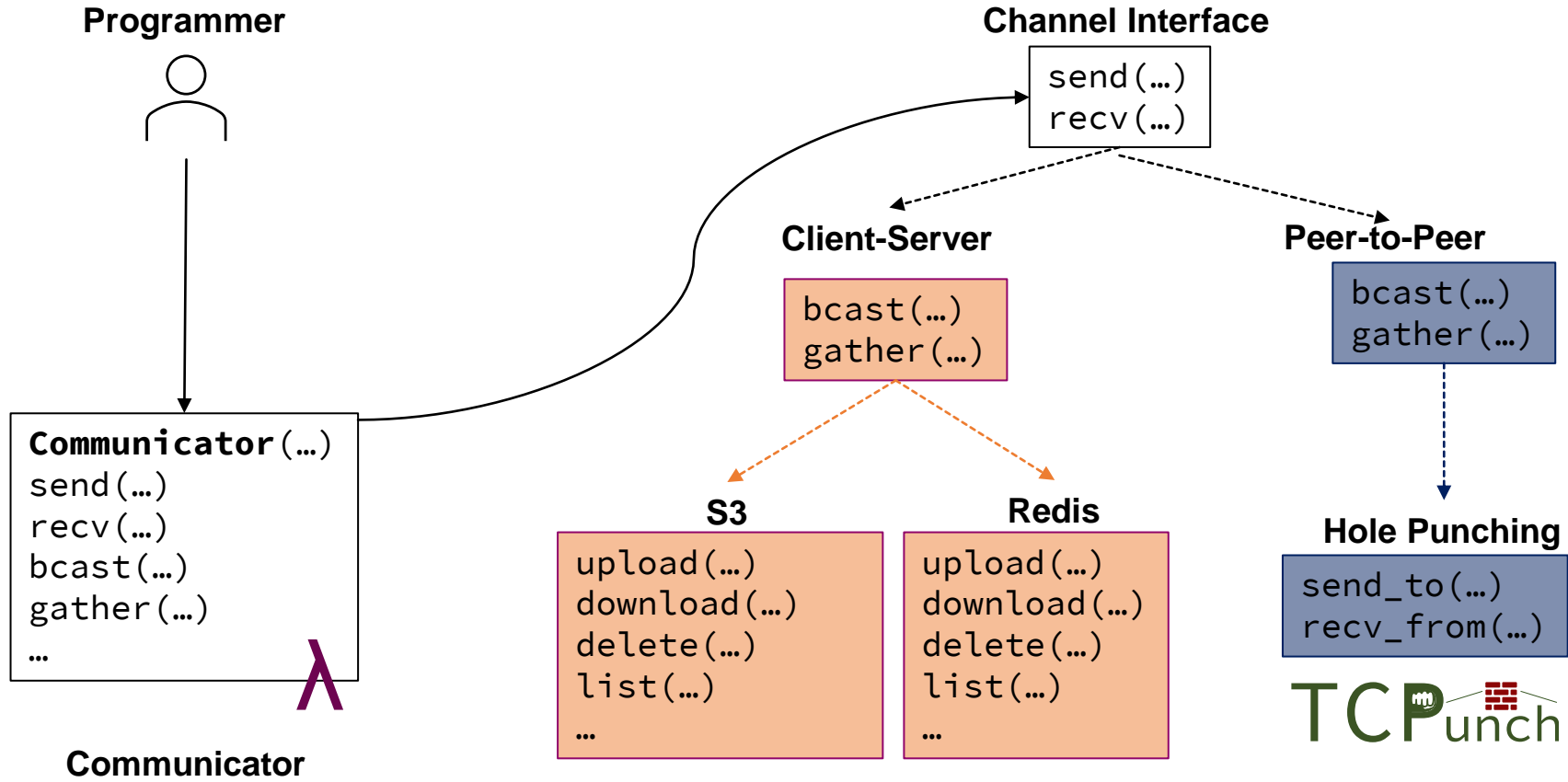
FMI: MPI for serverless



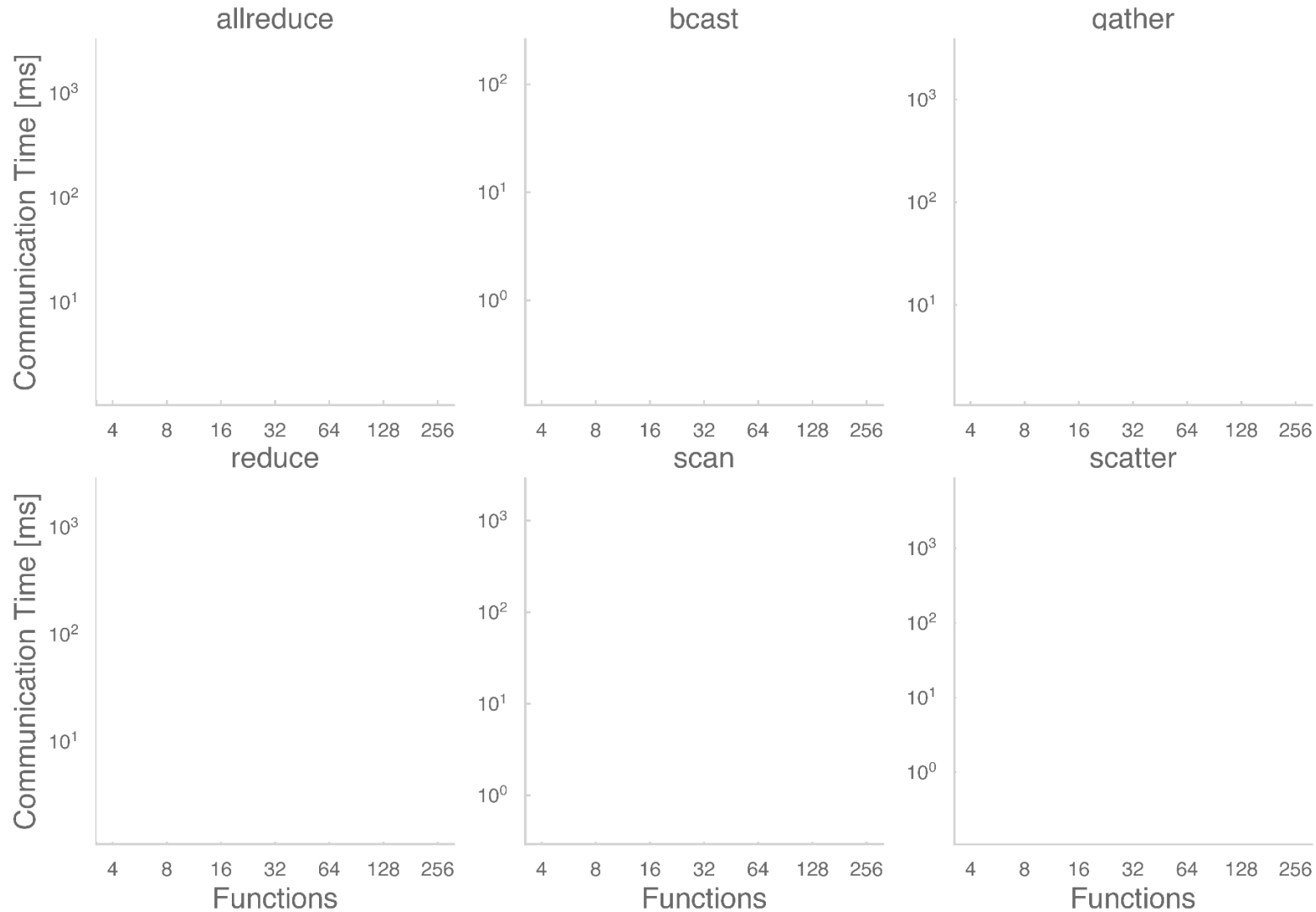
FMI: MPI for serverless



FMI: MPI for serverless

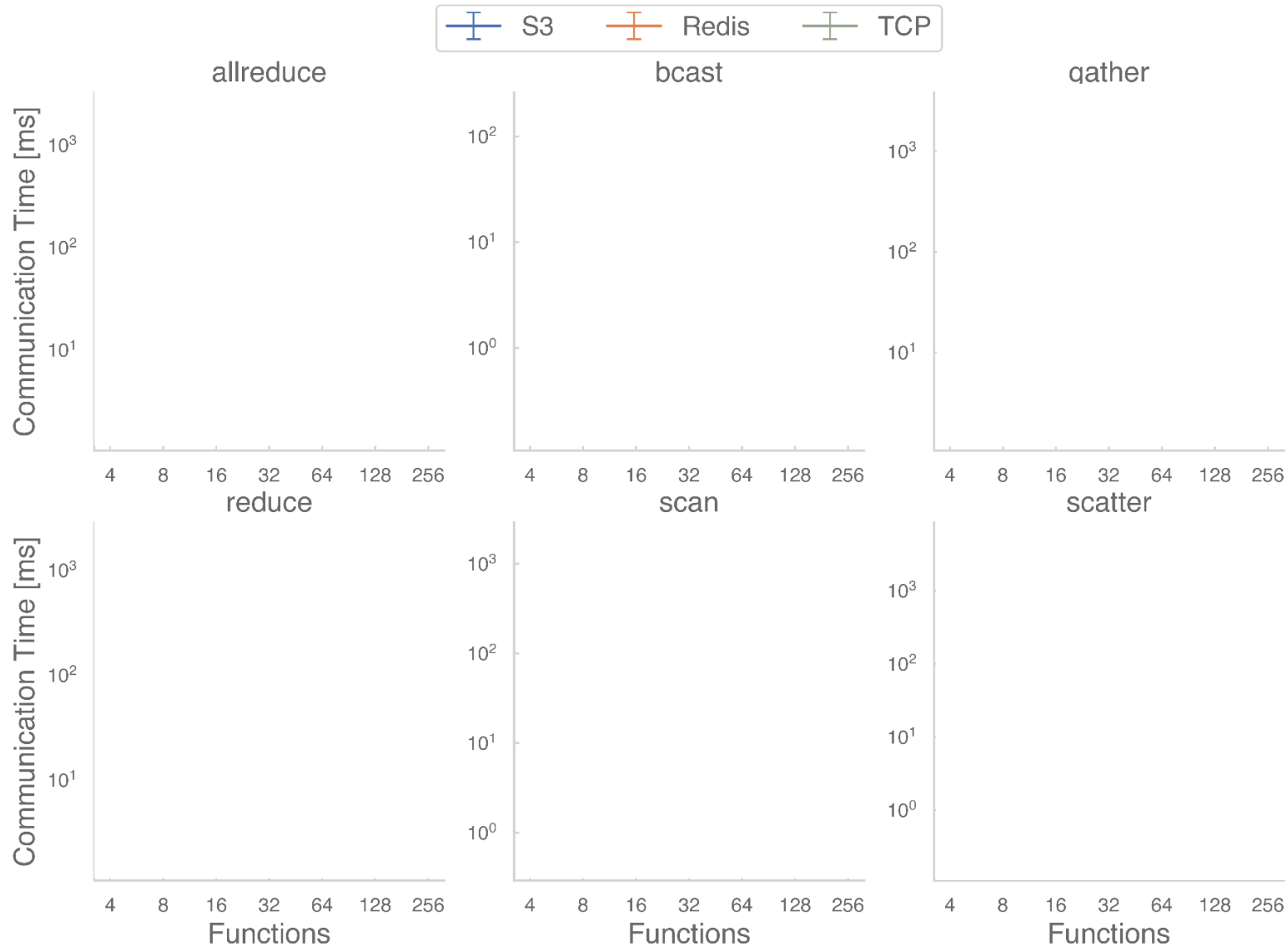


FMI on AWS Lambda

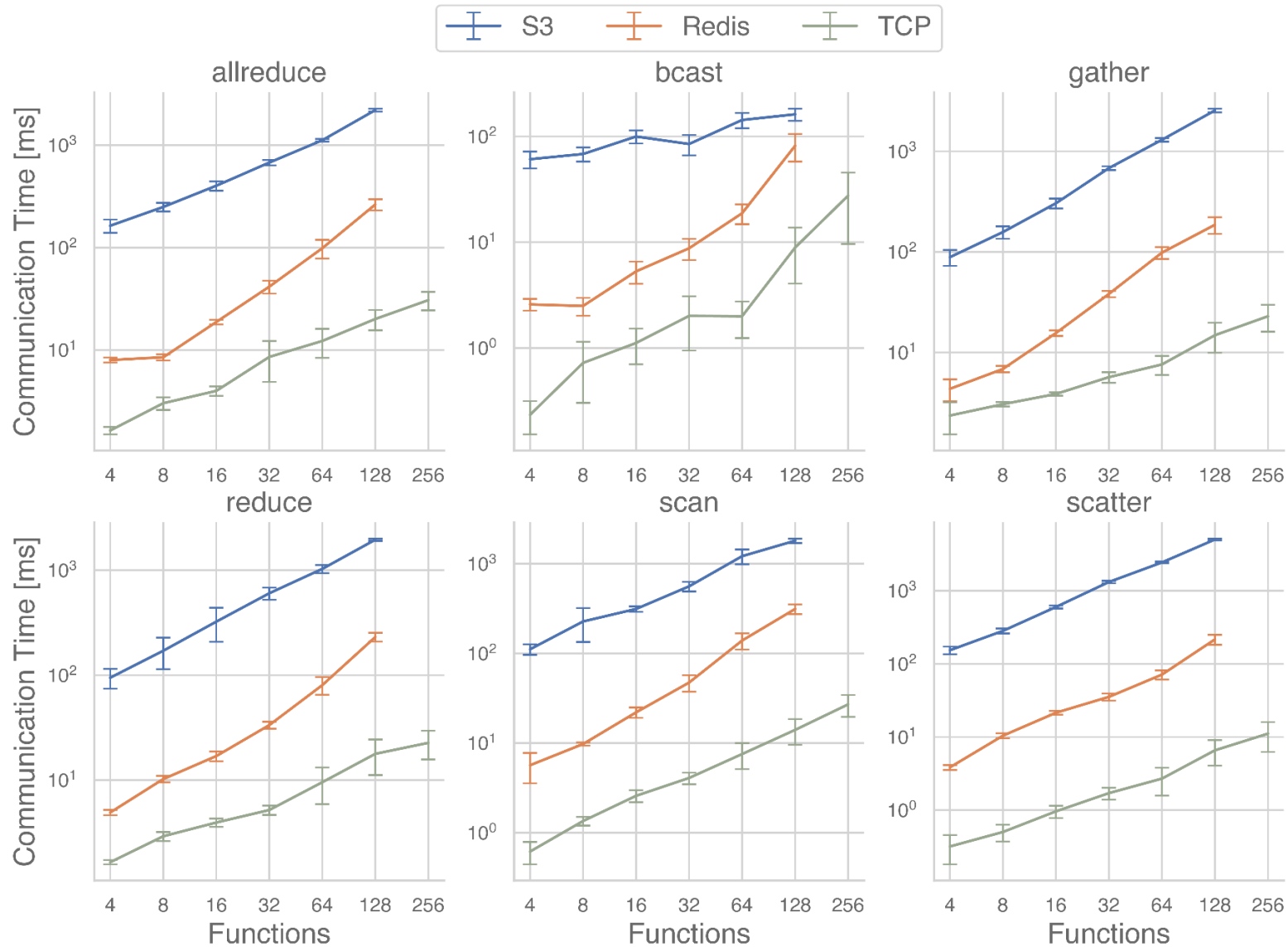


“FMI: Fast and Cheap Message Passing for Serverless Functions”, ICS’23

FMI on AWS Lambda



FMI on AWS Lambda



FaaS in High-Performance Applications

Serverless is slow

Answer:
rFaaS

Communication is slow
and restricted

Answer:
FMI

Serverless is hard to
program.

FaaS in High-Performance Applications

Serverless is slow

Answer:
rFaaS

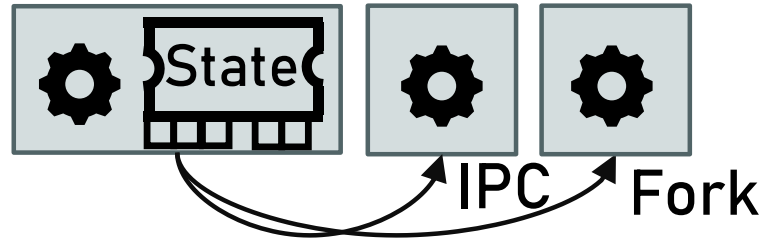
Communication is slow
and restricted

Answer:
FMI

Serverless is hard to
program.

Answer: Serverless
Processes

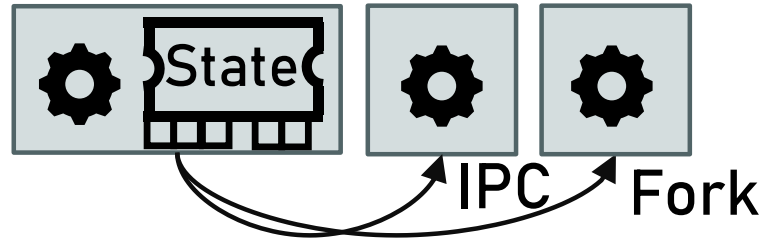
Serverless Process



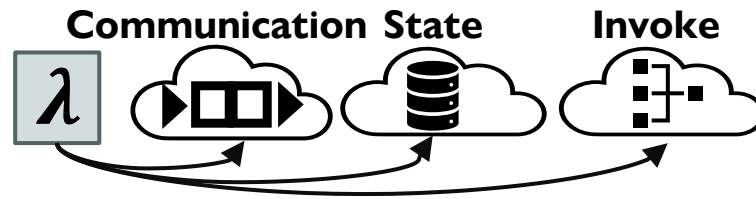
OS Process

Nano- and micro-second
latency of OS primitives.

Serverless Process

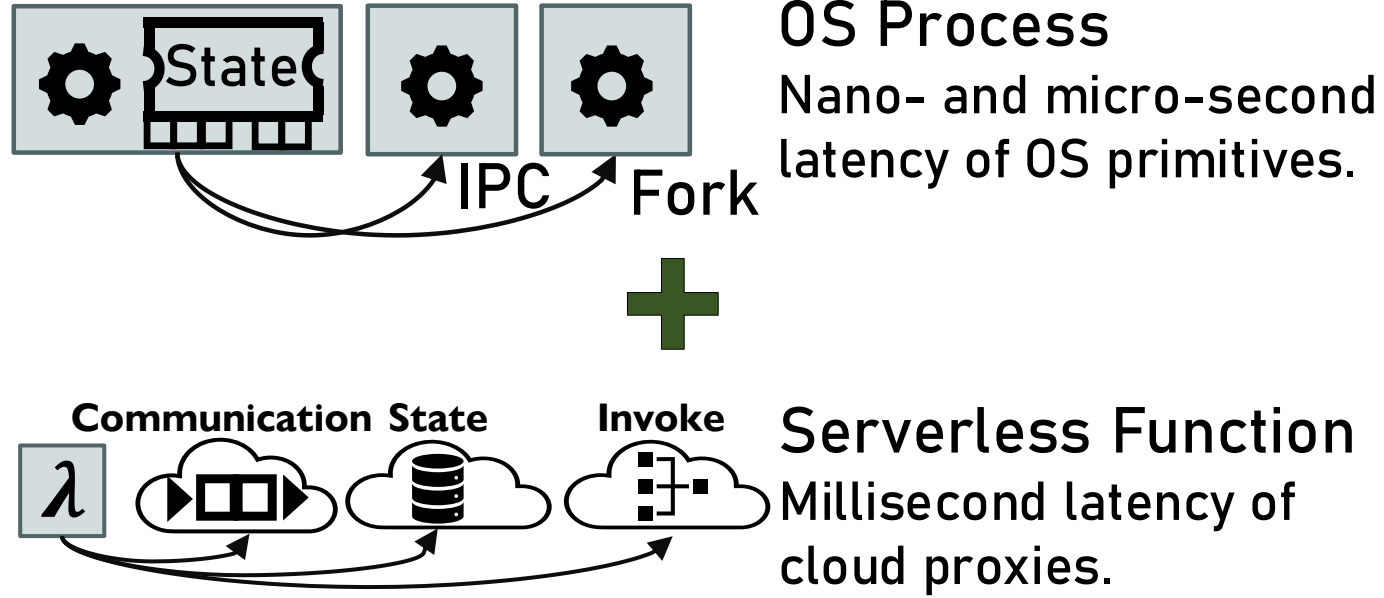


OS Process
Nano- and micro-second latency of OS primitives.

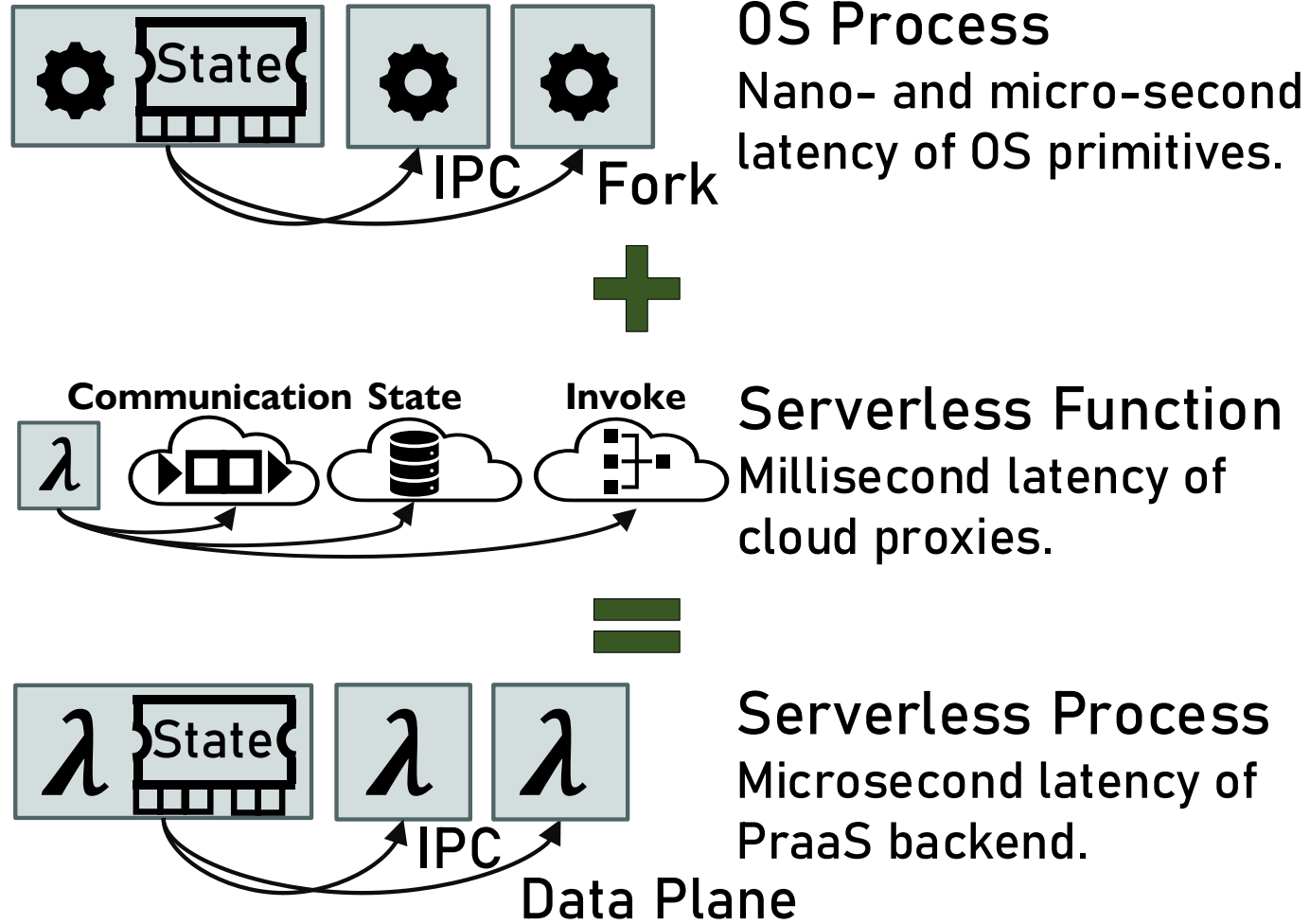


Serverless Function
Millisecond latency of cloud proxies.

Serverless Process



Serverless Process

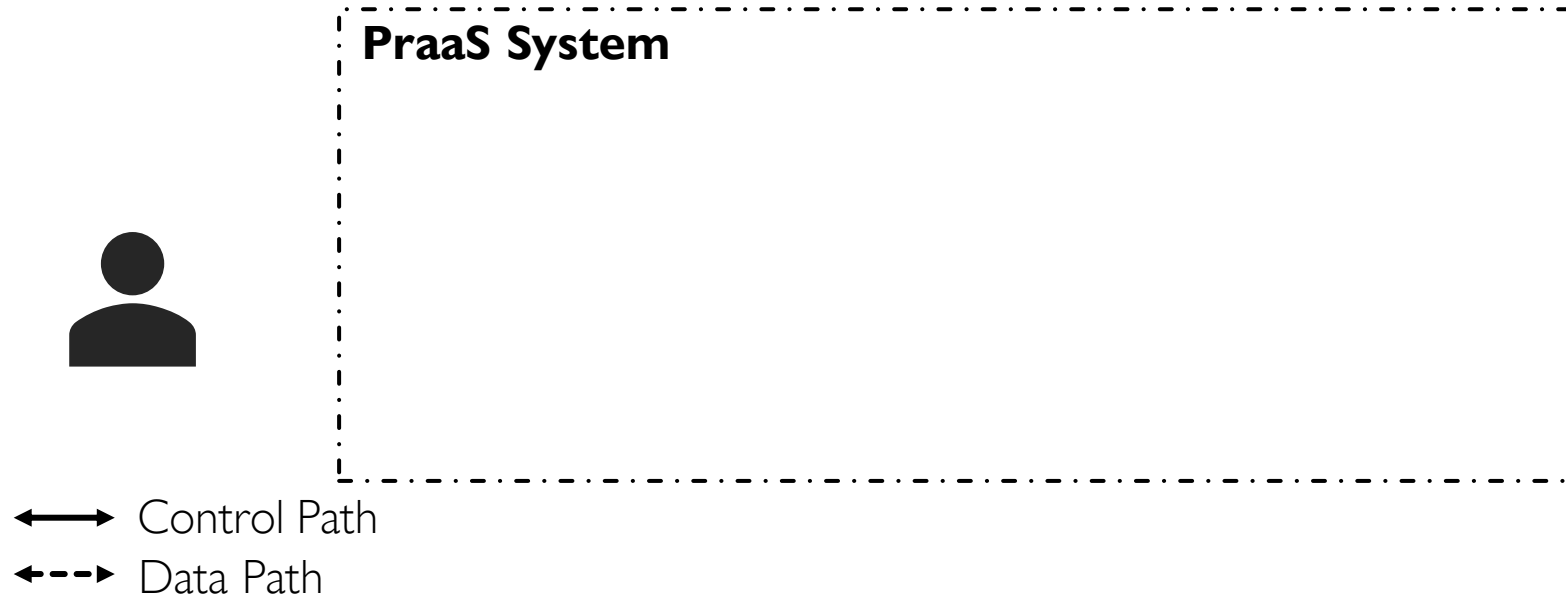


OS Process
Nano- and micro-second latency of OS primitives.

Serverless Function
Millisecond latency of cloud proxies.

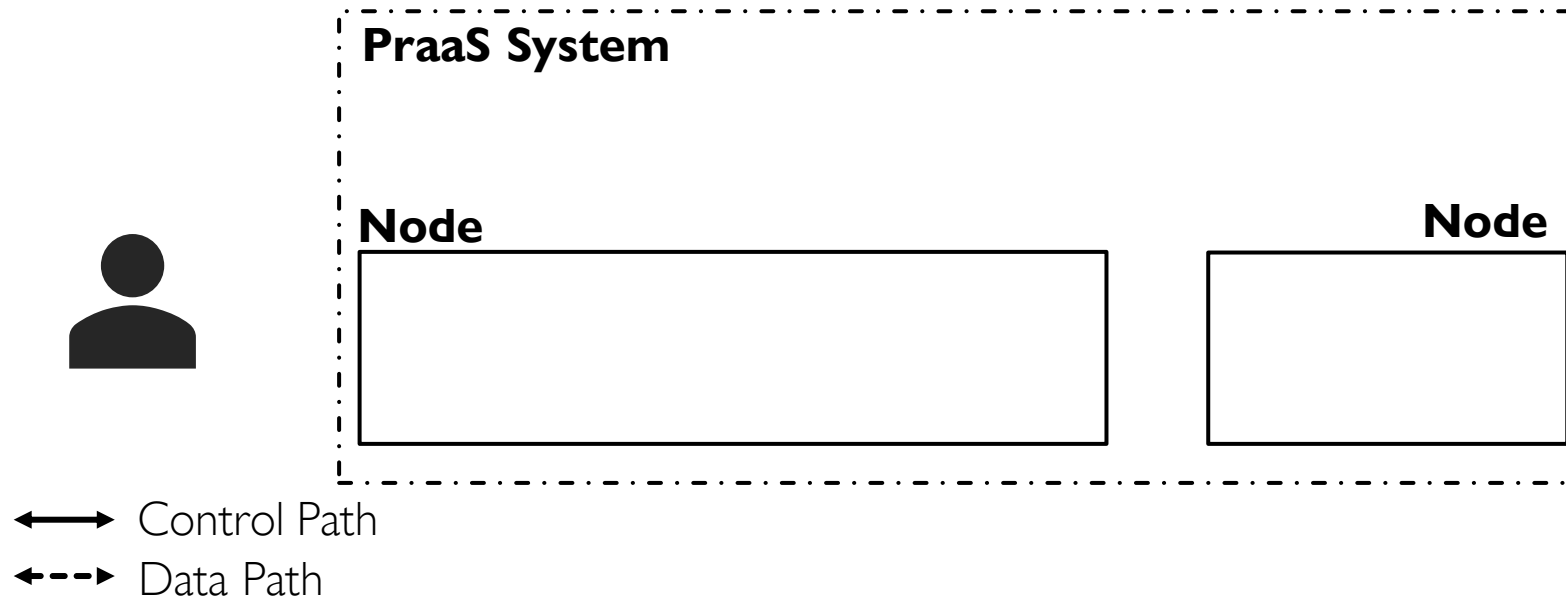
Serverless Process
Microsecond latency of PaaS backend.

PraaS: Process-as-Service



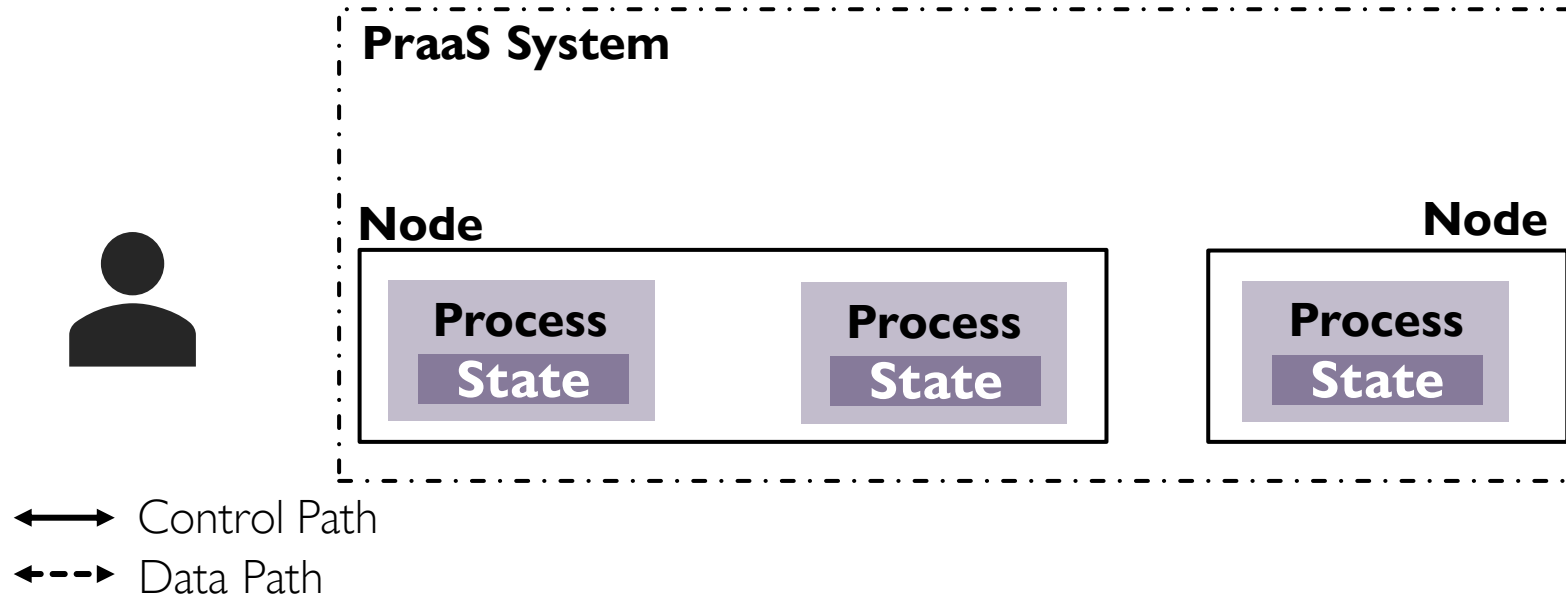
“Process-as-a-Service: FaaS Stateful Computing with Optimized Data Planes”, paper preprint.

PraaS: Process-as-Service



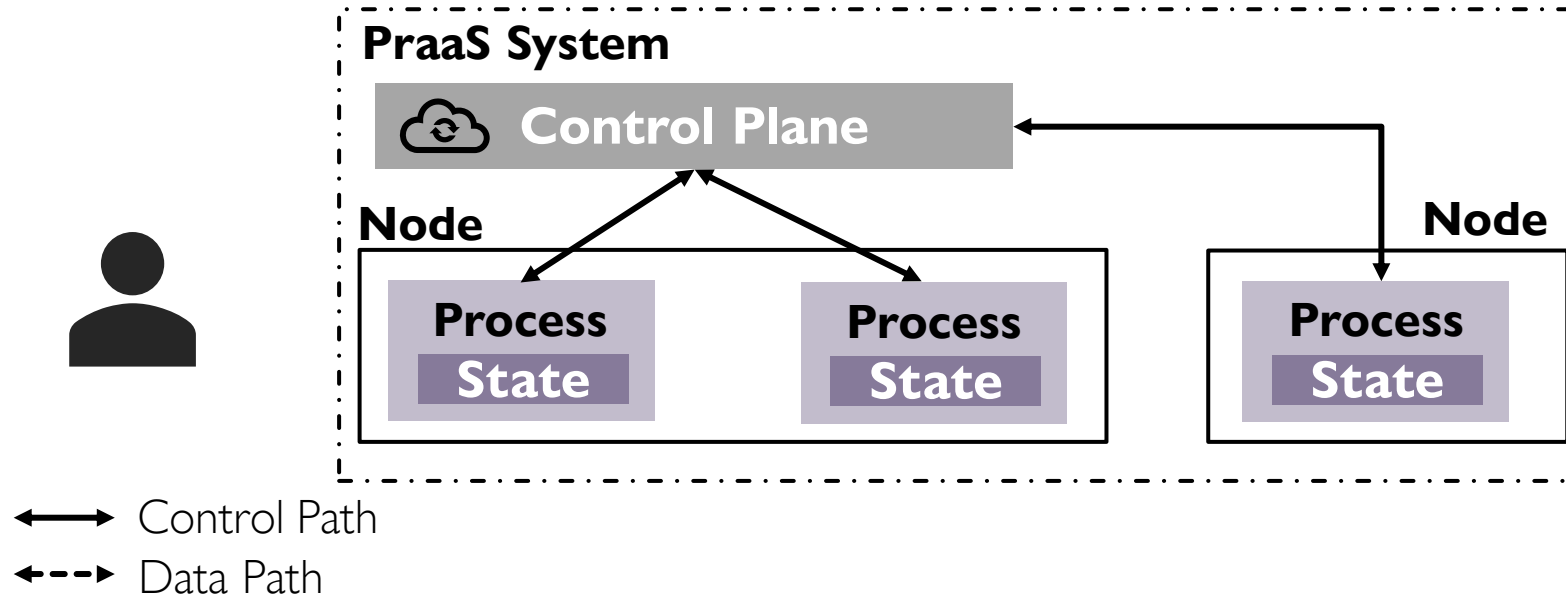
“Process-as-a-Service: FaaS Stateful Computing with Optimized Data Planes”, paper preprint.

PraaS: Process-as-Service



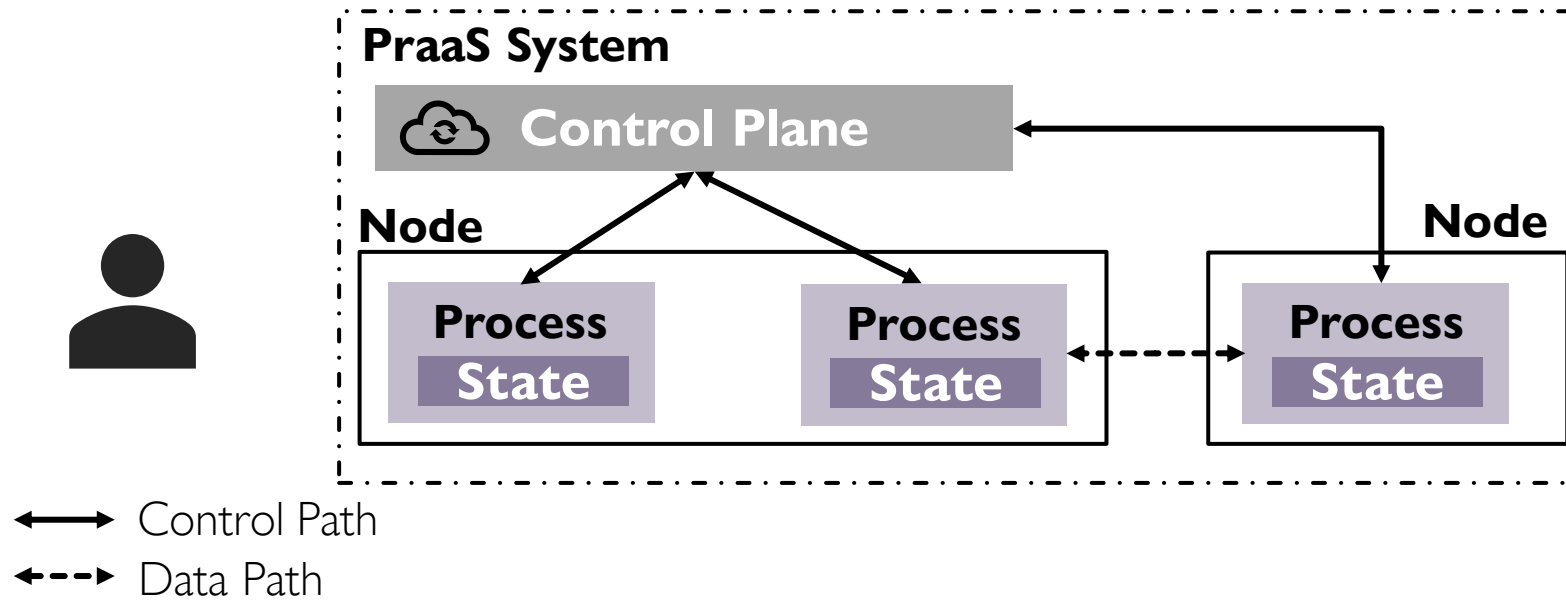
“Process-as-a-Service: FaaS Stateful Computing with Optimized Data Planes”, paper preprint.

PraaS: Process-as-Service



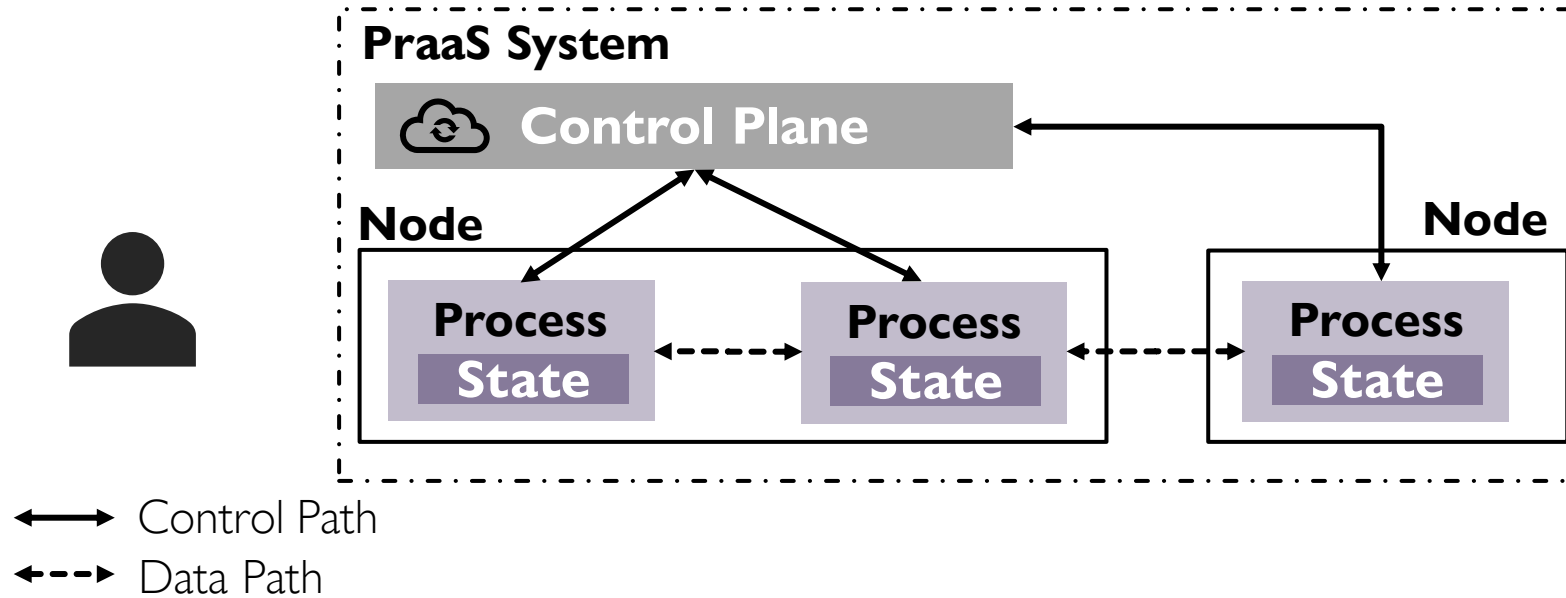
“Process-as-a-Service: FaaS Stateful Computing with Optimized Data Planes”, paper preprint.

PraaS: Process-as-Service



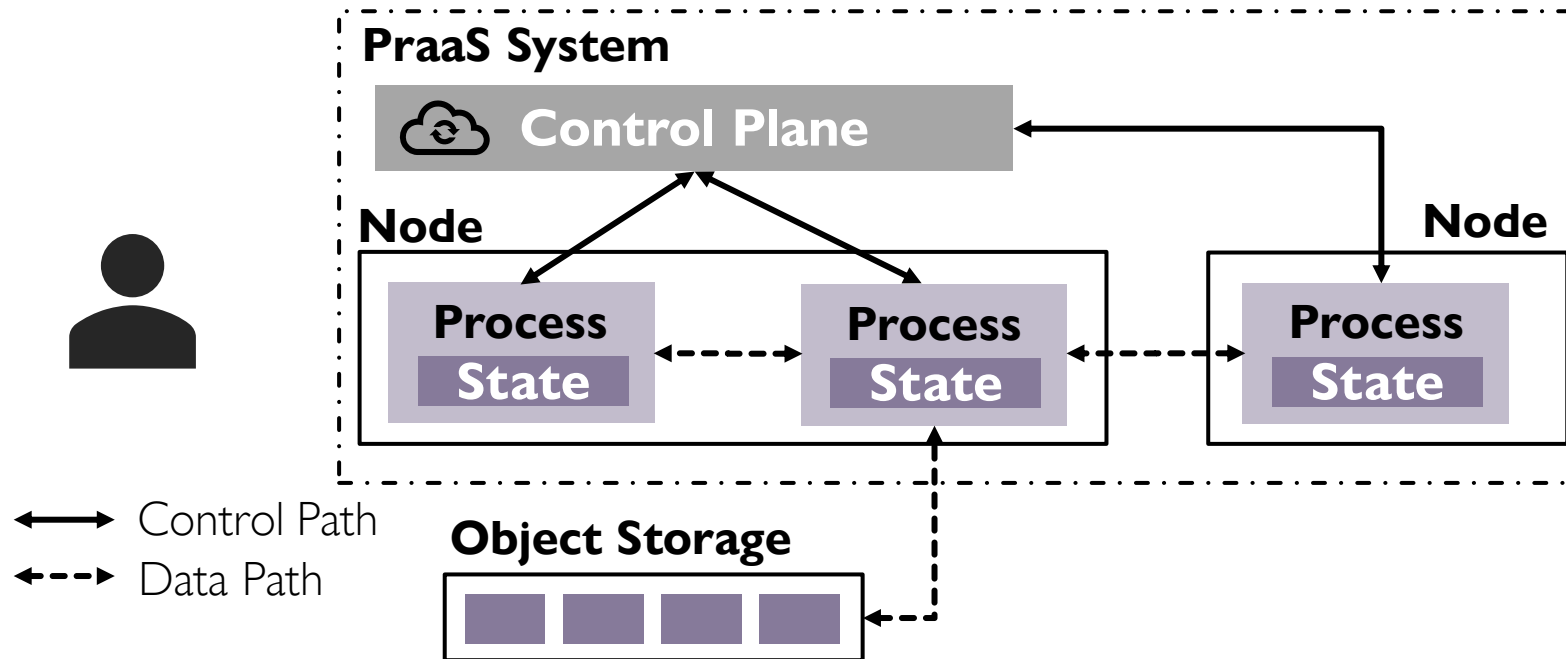
“Process-as-a-Service: FaaS Stateful Computing with Optimized Data Planes”, paper preprint.

PraaS: Process-as-Service



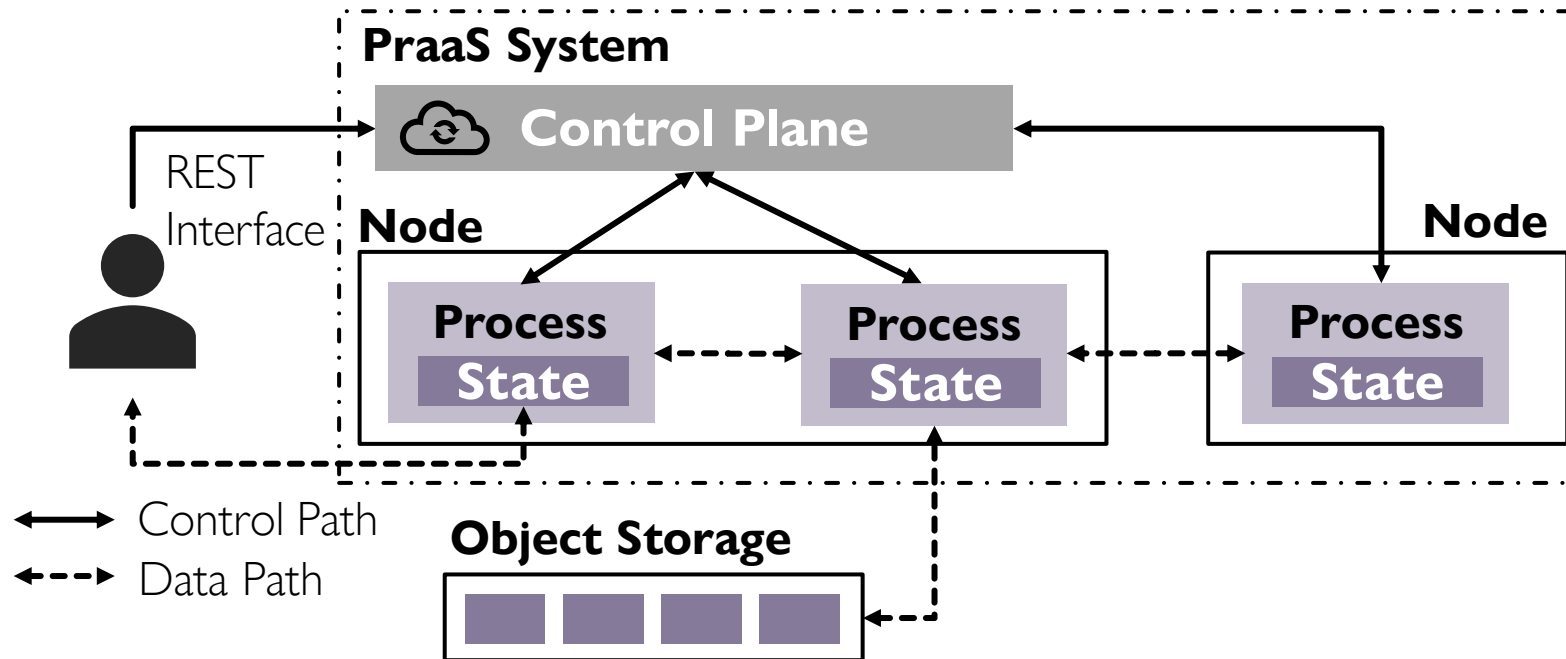
“Process-as-a-Service: FaaS Stateful Computing with Optimized Data Planes”, paper preprint.

PraaS: Process-as-Service

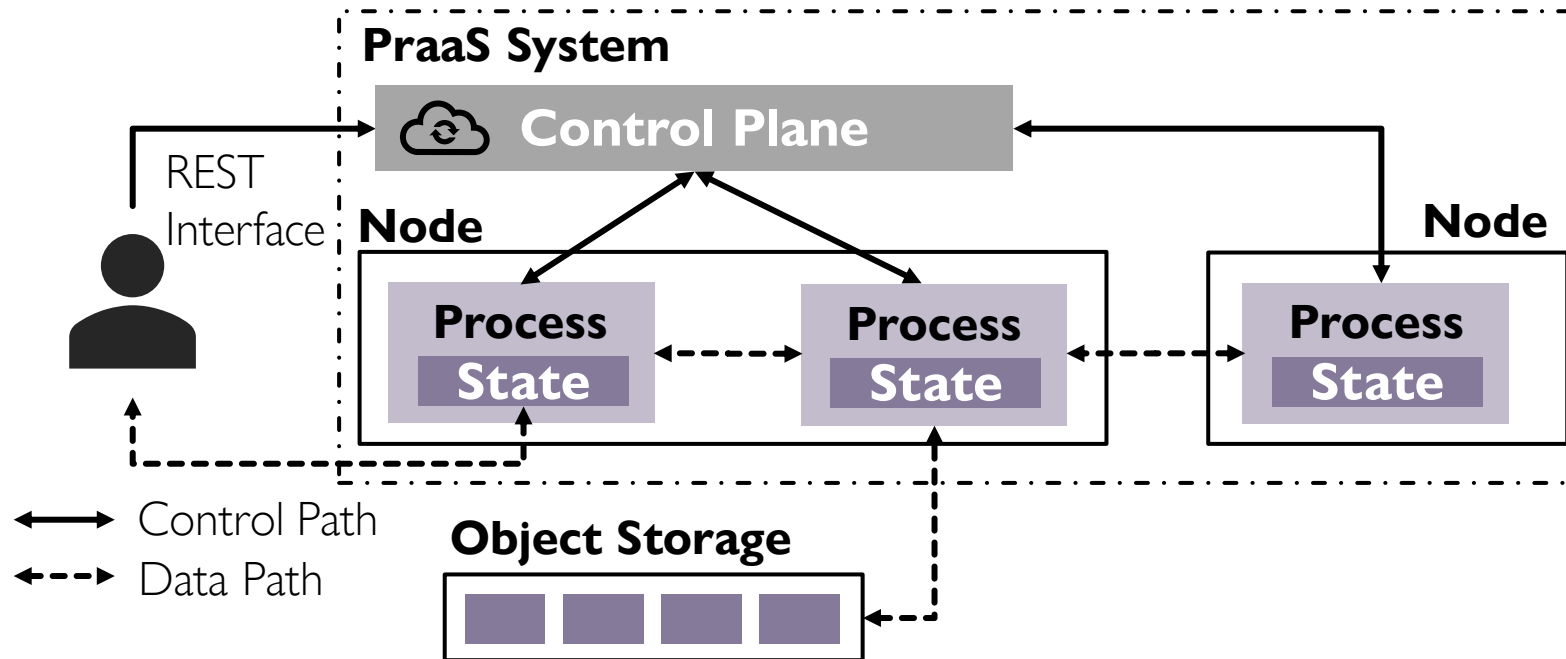


“Process-as-a-Service: FaaS Stateful Computing with Optimized Data Planes”, paper preprint.

PraaS: Process-as-Service



PraaS: Process-as-Service

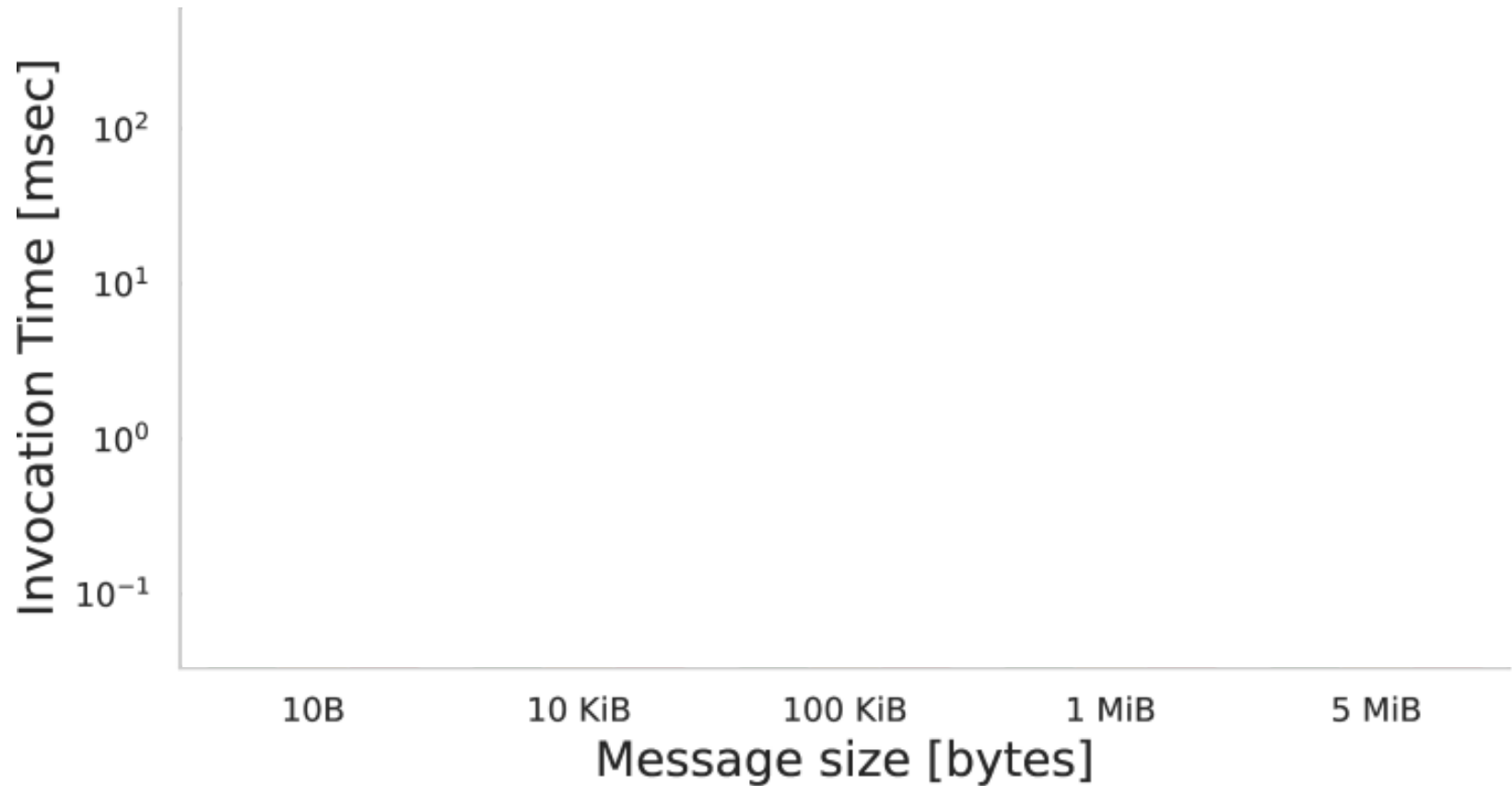


Works on AWS Fargate, Knative, Kubernetes.



“Process-as-a-Service: FaaS Stateful Computing with Optimized Data Planes”, paper preprint.

Serverless Process on Fargate vs AWS Lambda



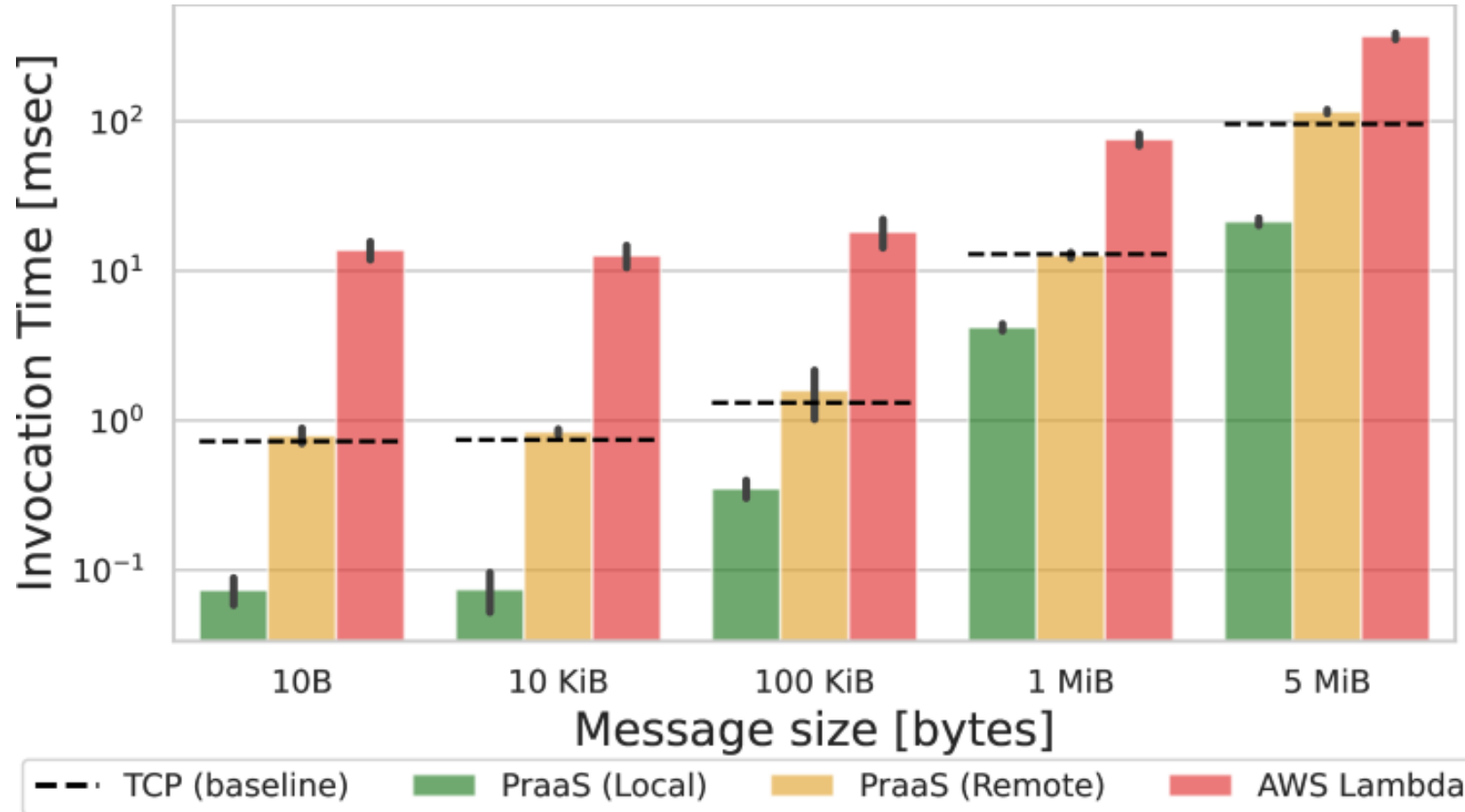
“Process-as-a-Service: FaaS Stateful Computing with Optimized Data Planes”, paper preprint.

Serverless Process on Fargate vs AWS Lambda



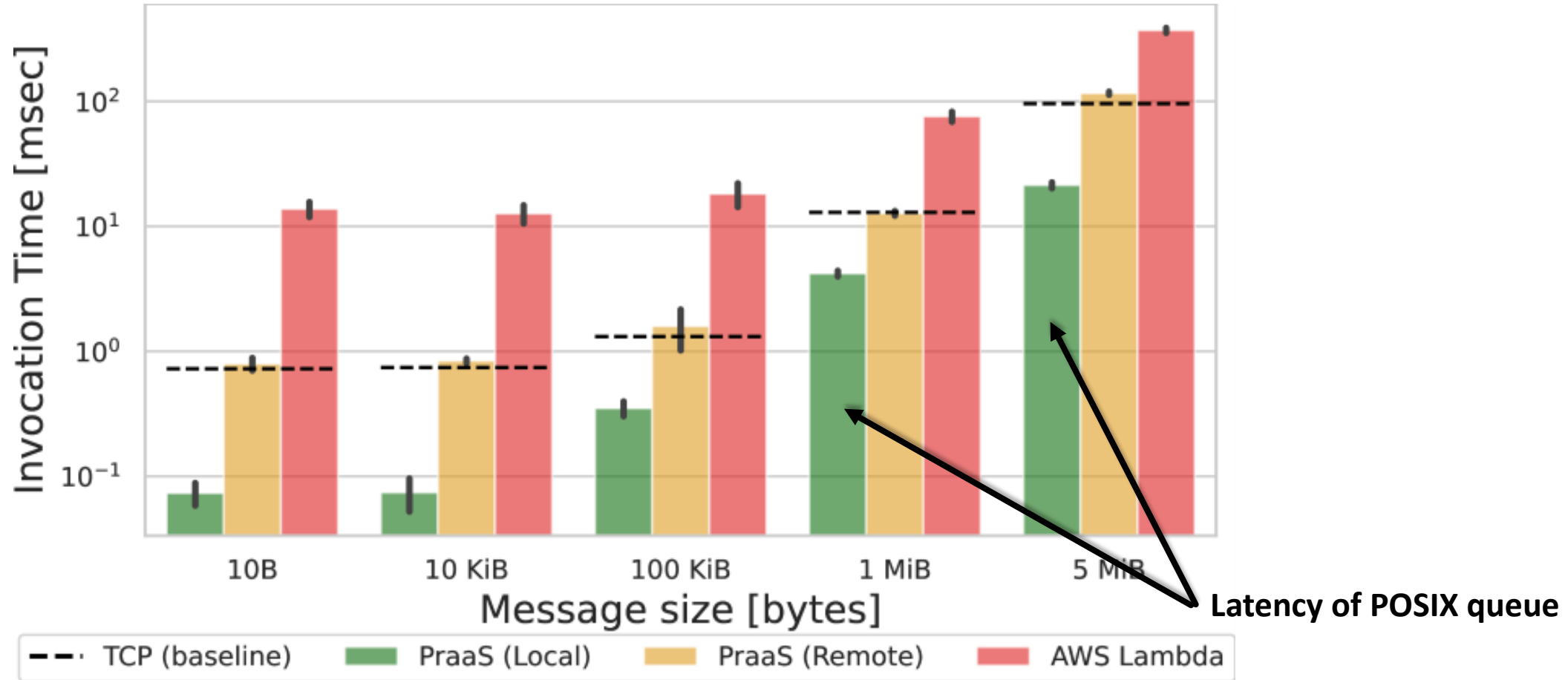
“Process-as-a-Service: FaaS Stateful Computing with Optimized Data Planes”, paper preprint.

Serverless Process on Fargate vs AWS Lambda



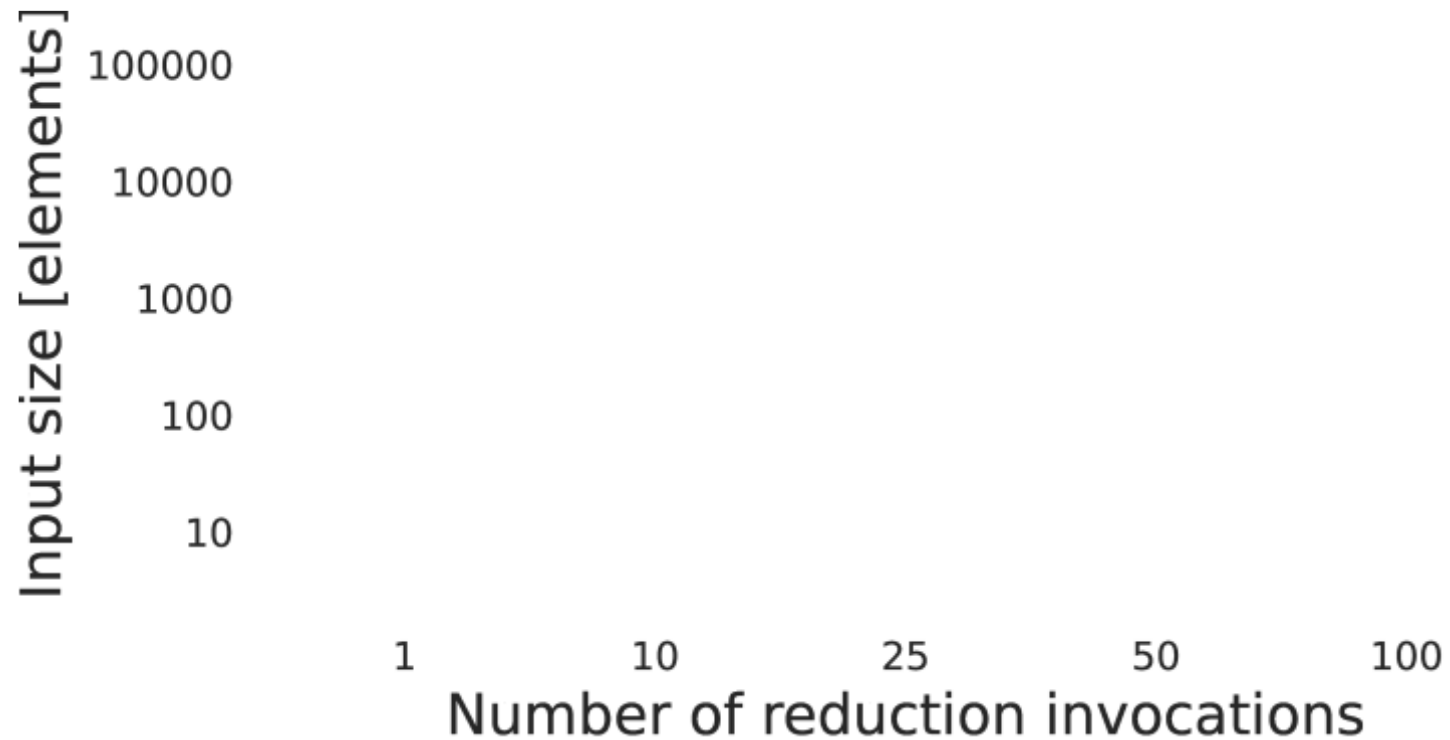
“Process-as-a-Service: FaaS Stateful Computing with Optimized Data Planes”, paper preprint.

Serverless Process on Fargate vs AWS Lambda



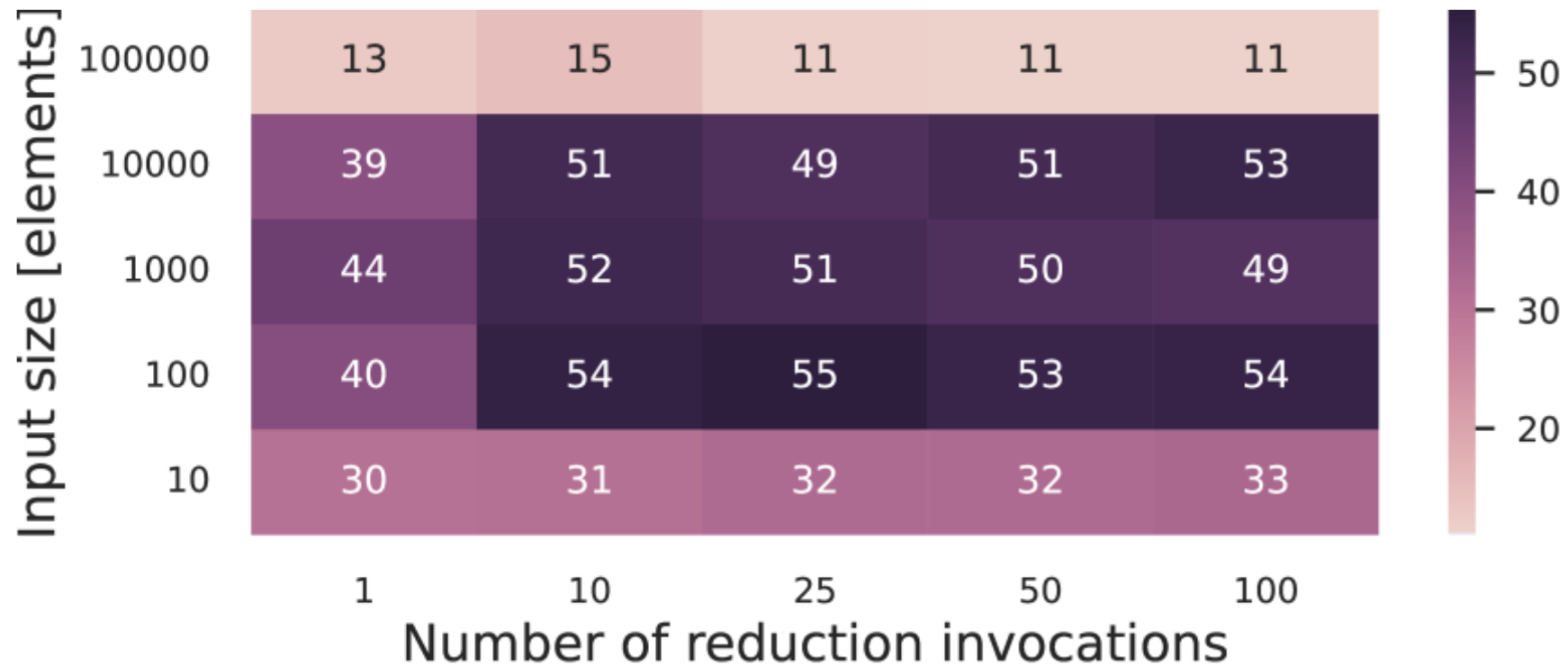
“Process-as-a-Service: FaaS Stateful Computing with Optimized Data Planes”, paper preprint.

Reduction Benchmark: Process State vs S3



“Process-as-a-Service: FaaS Stateful Computing with Optimized Data Planes”, paper preprint.

Reduction Benchmark: Process State vs S3



Serverless Solutions for HPC

Serverless Solutions for HPC



[spcl/serverless-benchmarks](https://github.com/spcl/serverless-benchmarks)

Serverless Solutions for HPC



[spcl/serverless-benchmarks](#)



[spcl/fmi](#)

Serverless Solutions for HPC



[spcl/serverless-benchmarks](#)



[spcl/fmi](#)



[spcl/rFaaS](#)

Serverless Solutions for HPC



[spcl/serverless-benchmarks](#)



[spcl/fmi](#)



[spcl/rFaaS](#)



[spcl/PraaS](#)

Serverless challenges in HPC

Serverless challenges in HPC

Poor vertical integration

Serverless challenges in HPC

Poor vertical integration

Expensive computing

Serverless challenges in HPC

Poor vertical integration

Expensive computing

Lack of heterogeneity

Serverless challenges in HPC

Poor vertical integration

Expensive computing

Lack of heterogeneity

Restricted environments

Serverless challenges in HPC

Poor vertical integration

Expensive computing

How to integrate
functions?

Lack of heterogeneity

Restricted environments


Conclusions



More of SPCL's research:

 youtube.com/@spcl **150+ Talks**

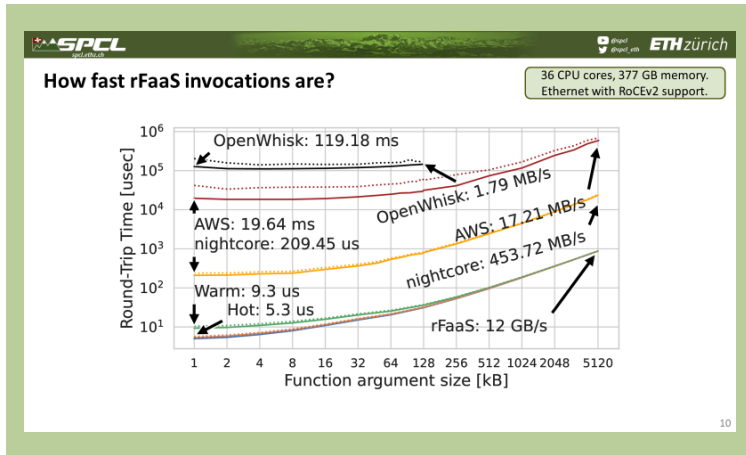
 twitter.com/spcl_eth **1.2K+ Followers**

 github.com/spcl **2K+ Stars**

... or spcl.ethz.ch



Conclusions



More of SPCL's research:

youtube.com/@spcl **150+ Talks**

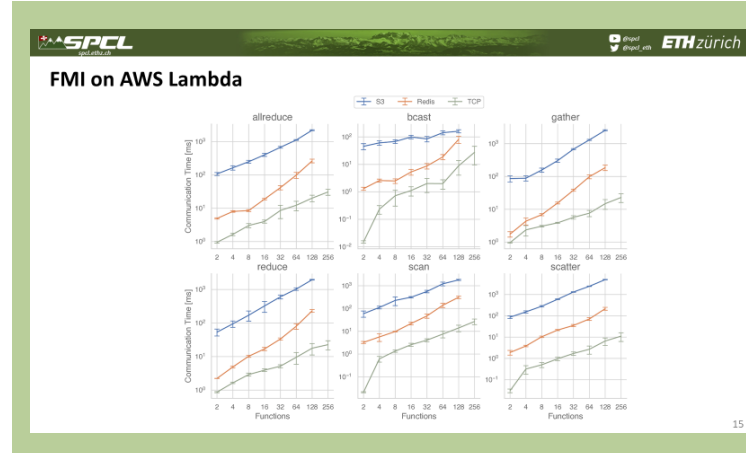
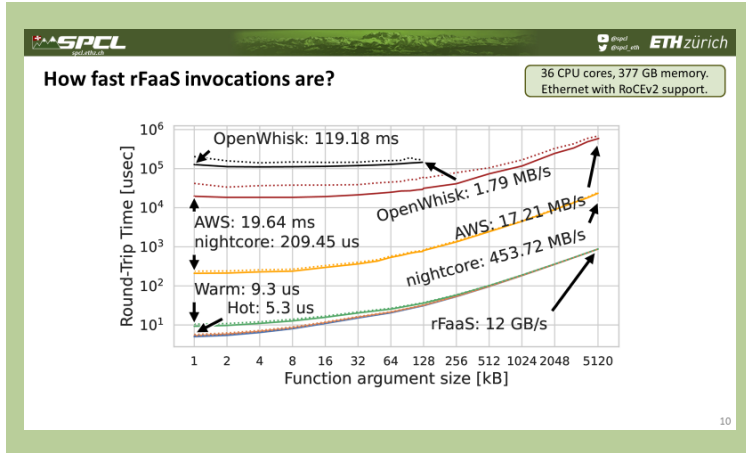
twitter.com/spcl_eth **1.2K+ Followers**

github.com/spcl **2K+ Stars**

... or spcl.ethz.ch



Conclusions



More of SPCL's research:

[youtube.com/@spcl](https://www.youtube.com/@spcl) **150+ Talks**

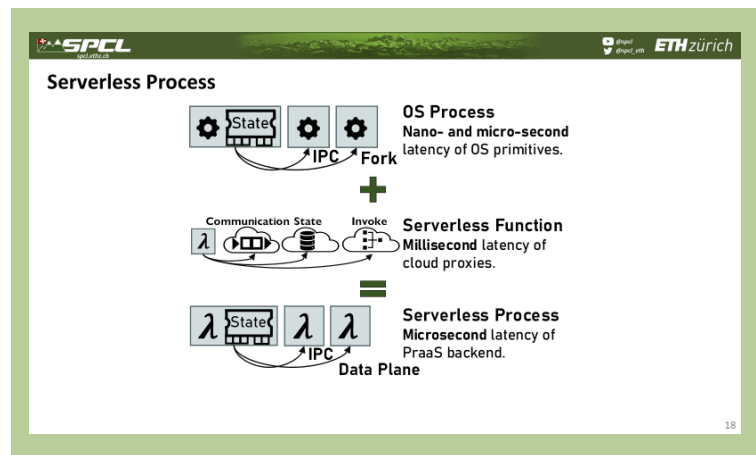
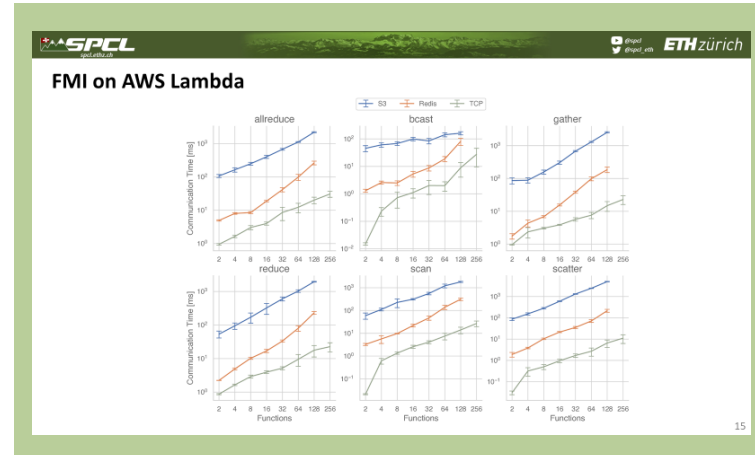
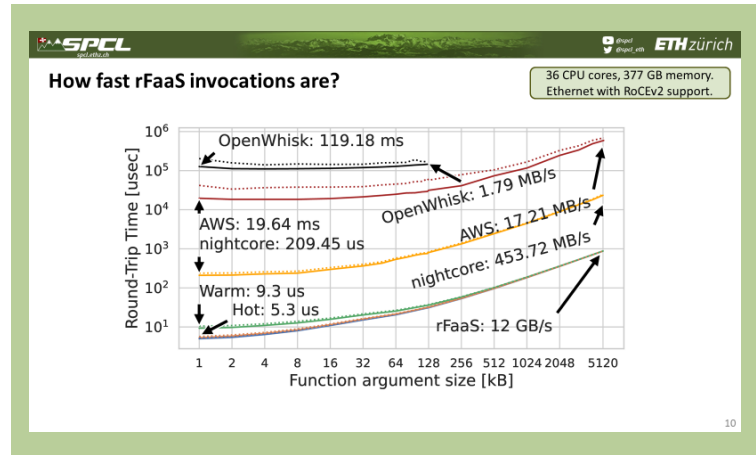
twitter.com/spcl_eth **1.2K+ Followers**

github.com/spcl **2K+ Stars**




... or spcl.ethz.ch



Conclusions



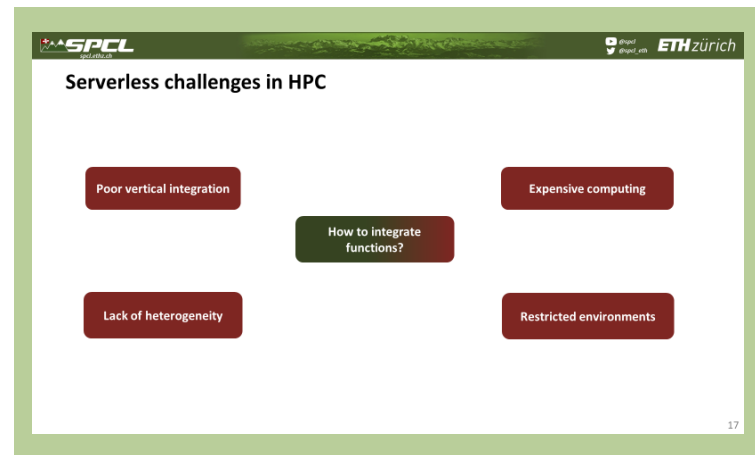
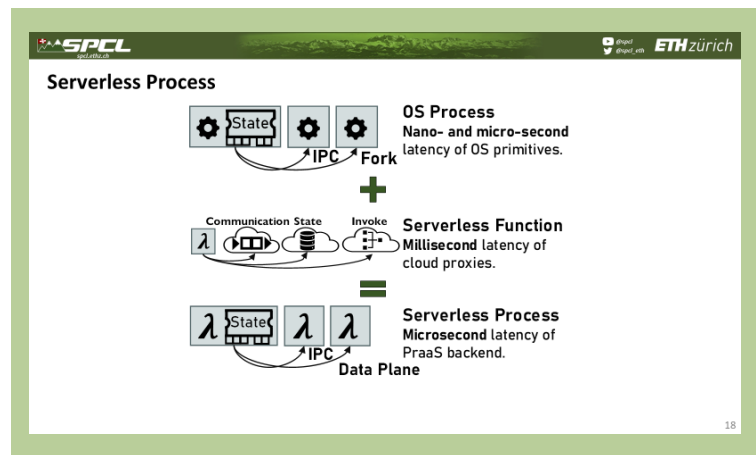
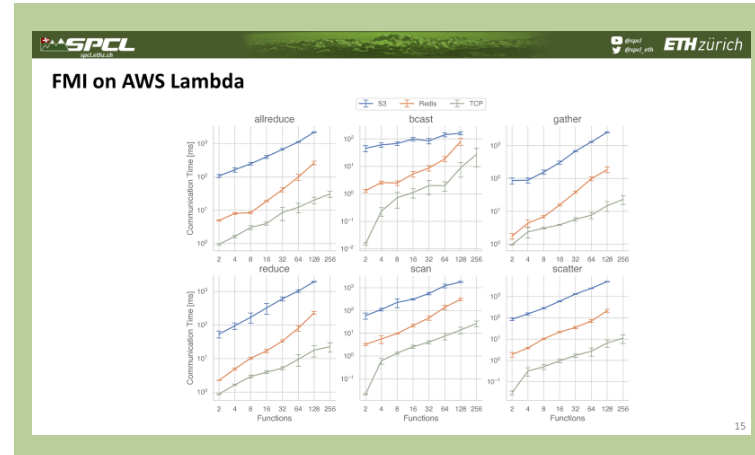
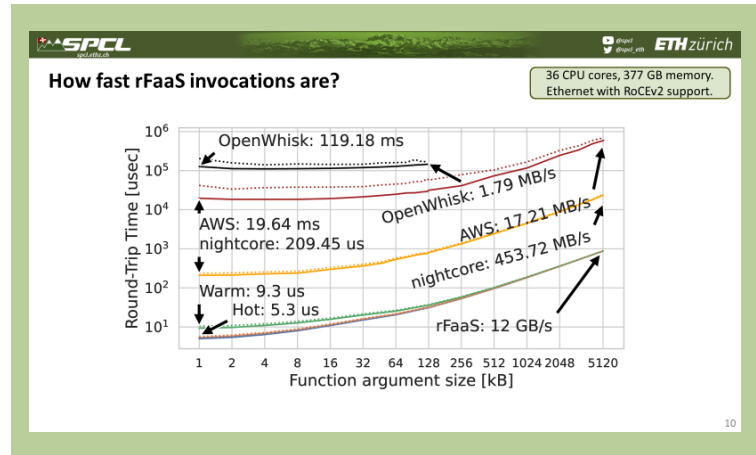
More of SPCL's research:

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

... or spcl.ethz.ch



Conclusions



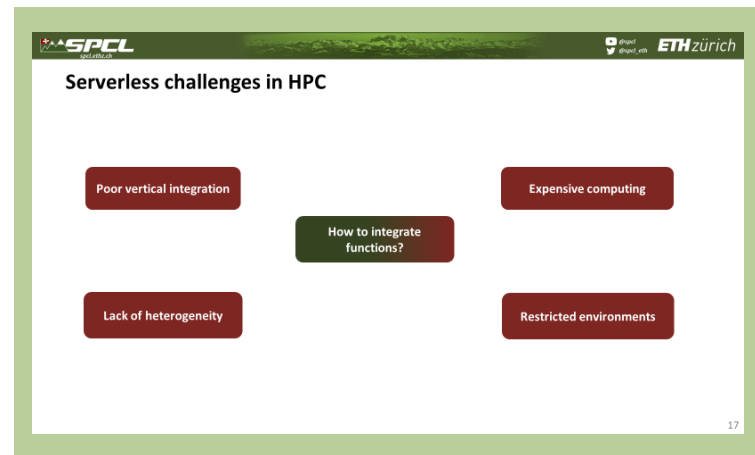
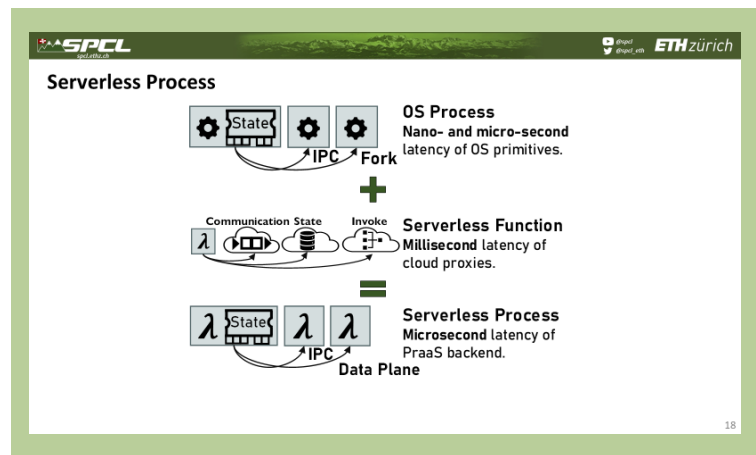
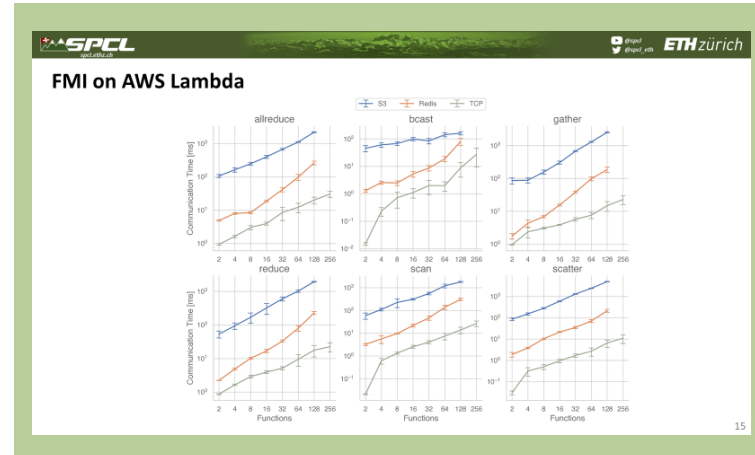
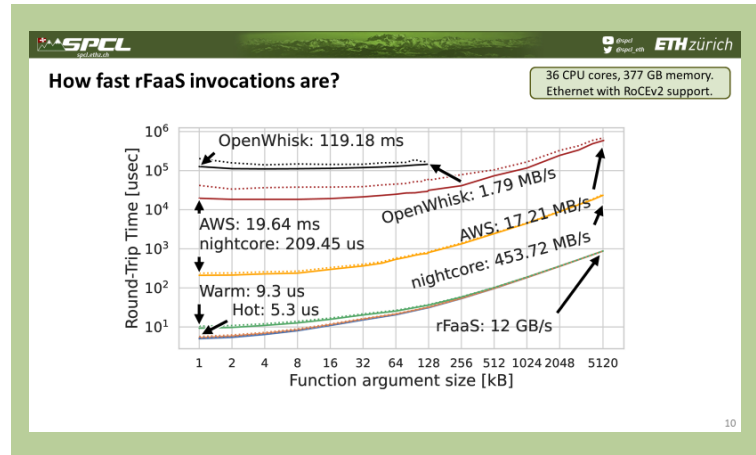
More of SPCL's research:

- 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](https://www.youtube.com/@spcl) **150+ Talks**

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

github.com/spcl **2K+ Stars**

... or spcl.ethz.ch



Poster. **Personal website.**