Wei Qiu, Marcin Copik, Yun Wang, Alexandru Calotoiu, Torsten Hoefler

# User-guided Page Merging for Memory Deduplication in Serverless Systems

IEEE **BIG DATA** 2023
15–18 December • Sorrento, Italy

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

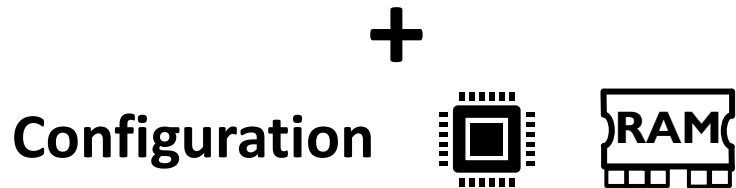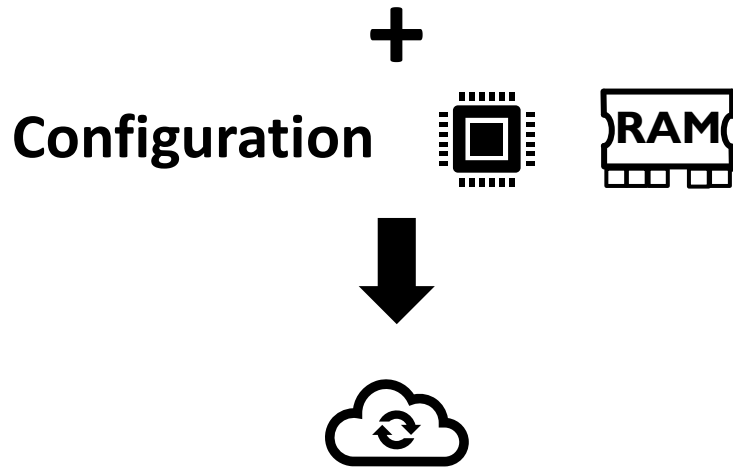# How does Function-as-a-Service (FaaS) work?

```python
def handler_function(req: dict, context: dict):

    model = cloud_storage.download_model()

    input = parse_input(req['payload'])

    output = model.inference(input)

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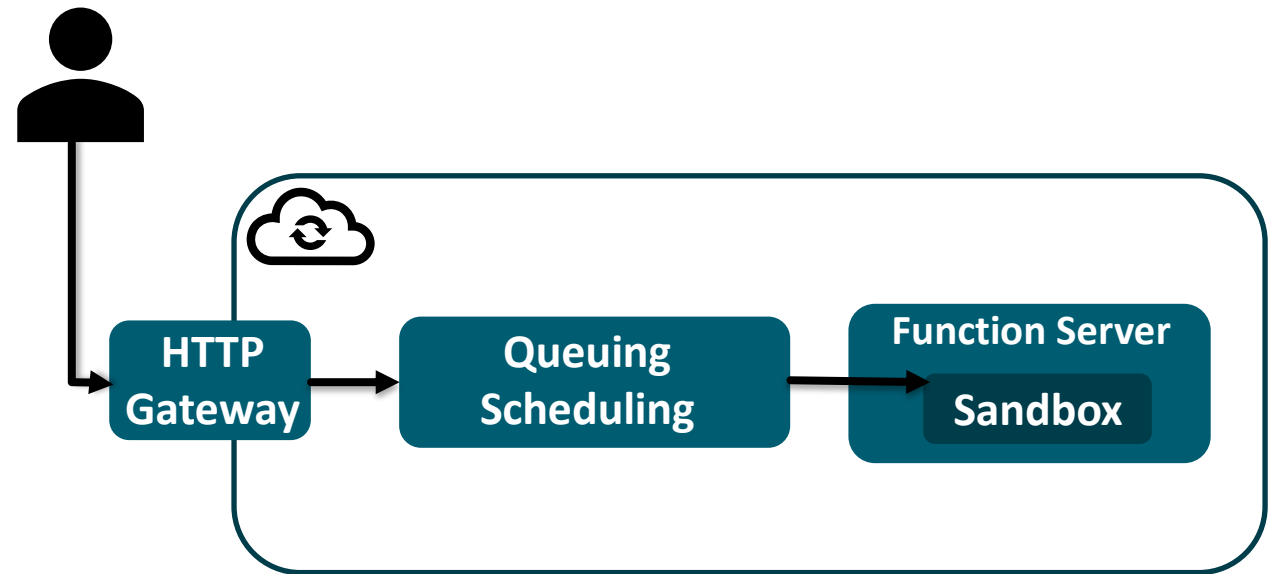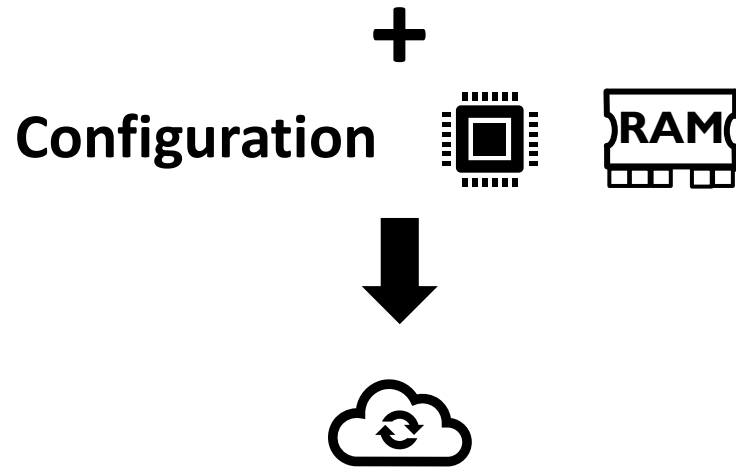
# How does Function-as-a-Service (FaaS) work?

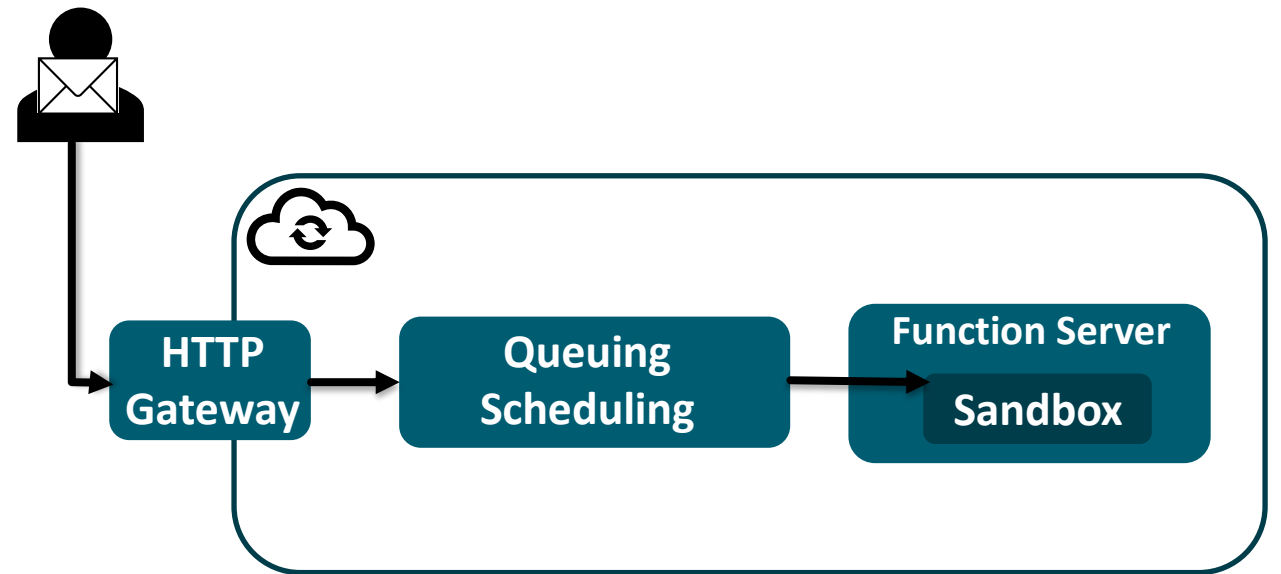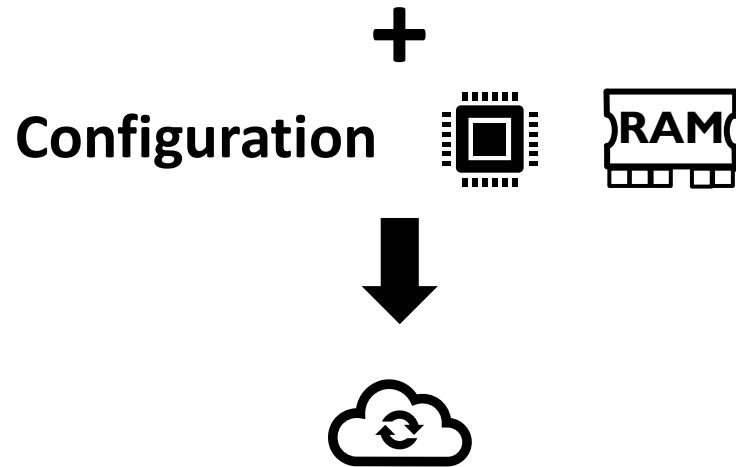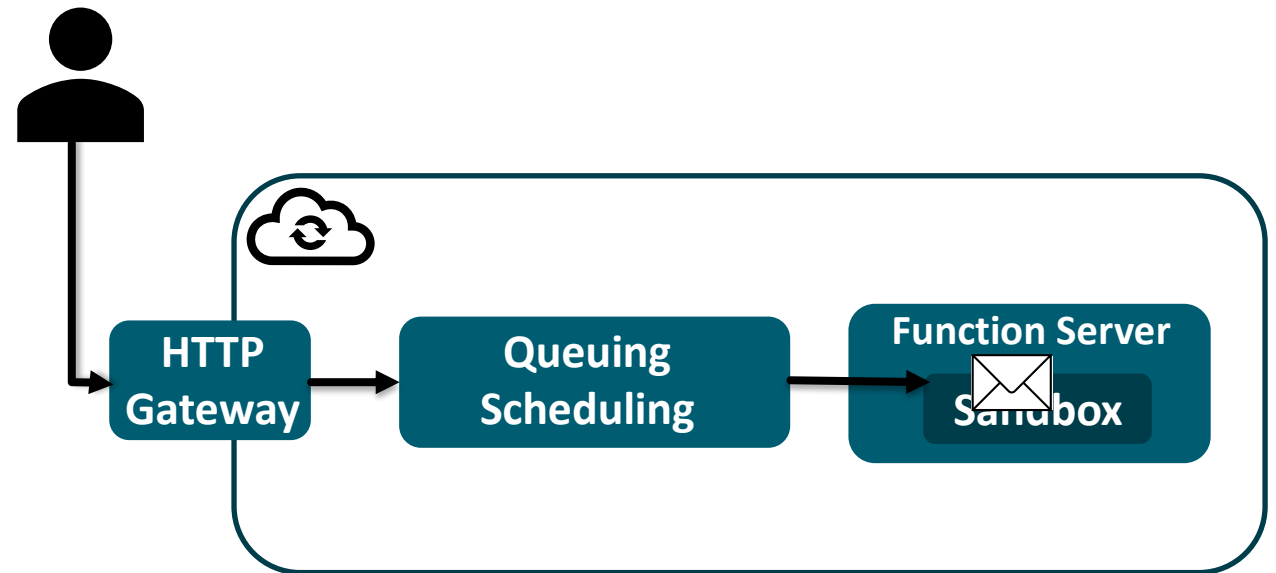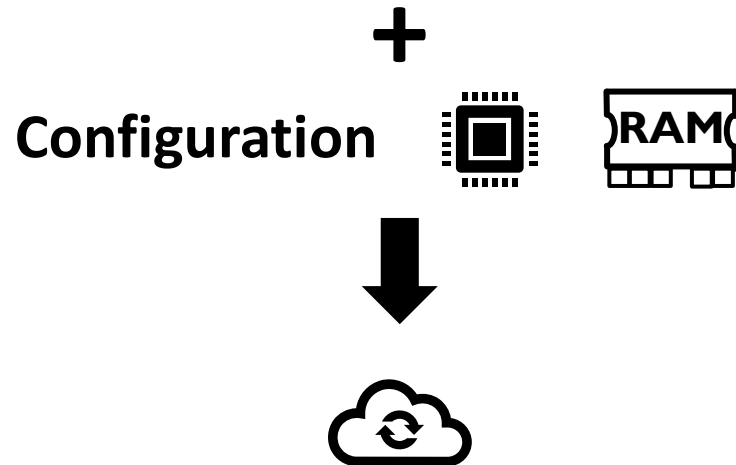```python
def handler_function(req: dict, context: dict):

    model = cloud_storage.download_model()

    input = parse_input(req['payload'])

    output = model.inference(input)

    return output
```

**+**

**Configuration**

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

```python
def handler_function(req: dict, context: dict):

    model = cloud_storage.download_model()

    input = parse_input(req['payload'])

    output = model.inference(input)

    return output
```

**+**

**Configuration**

3

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

```python
def handler_function(req: dict, context: dict):

    model = cloud_storage.download_model()

    input = parse_input(req['payload'])

    output = model.inference(input)

    return output
```

**+**

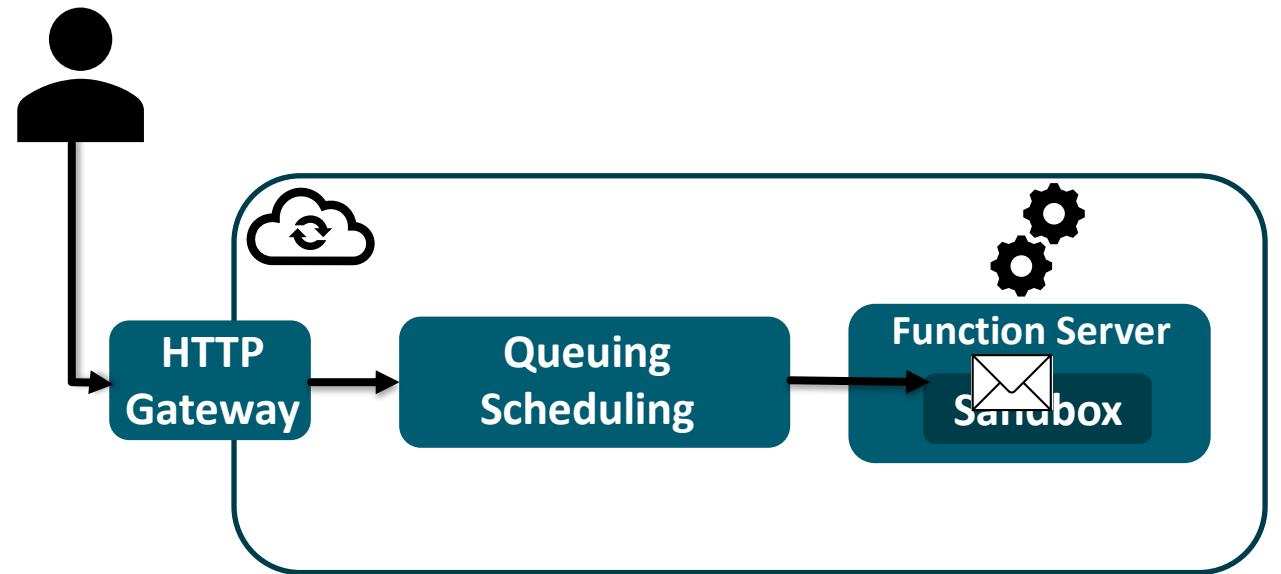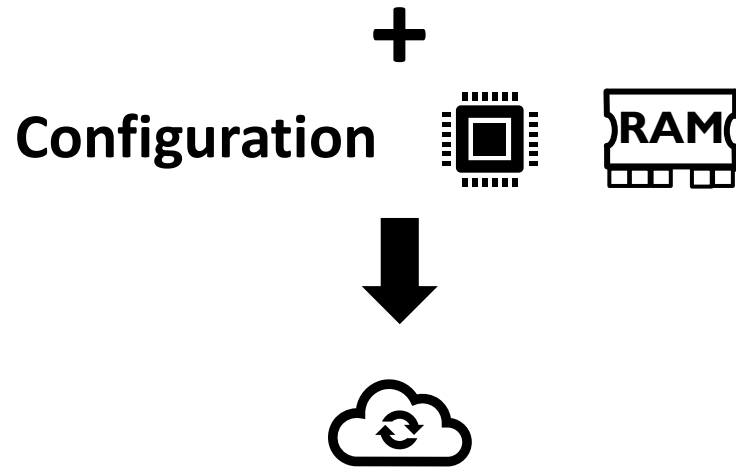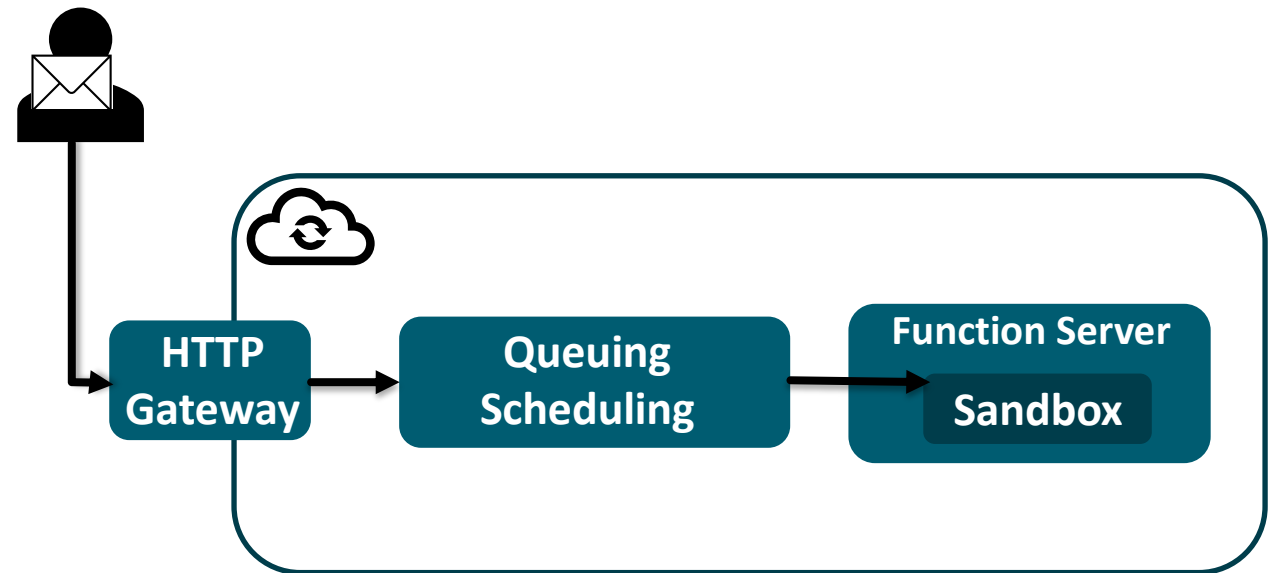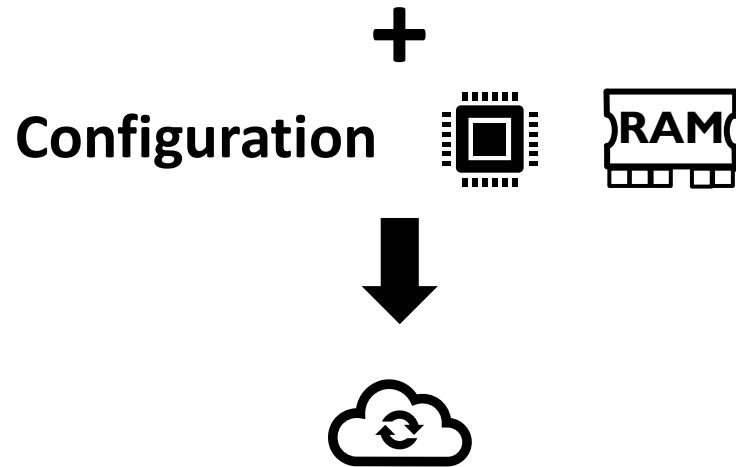**Configuration**

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

```python
def handler_function(req: dict, context: dict):

    model = cloud_storage.download_model()

    input = parse_input(req['payload'])

    output = model.inference(input)

    return output
```
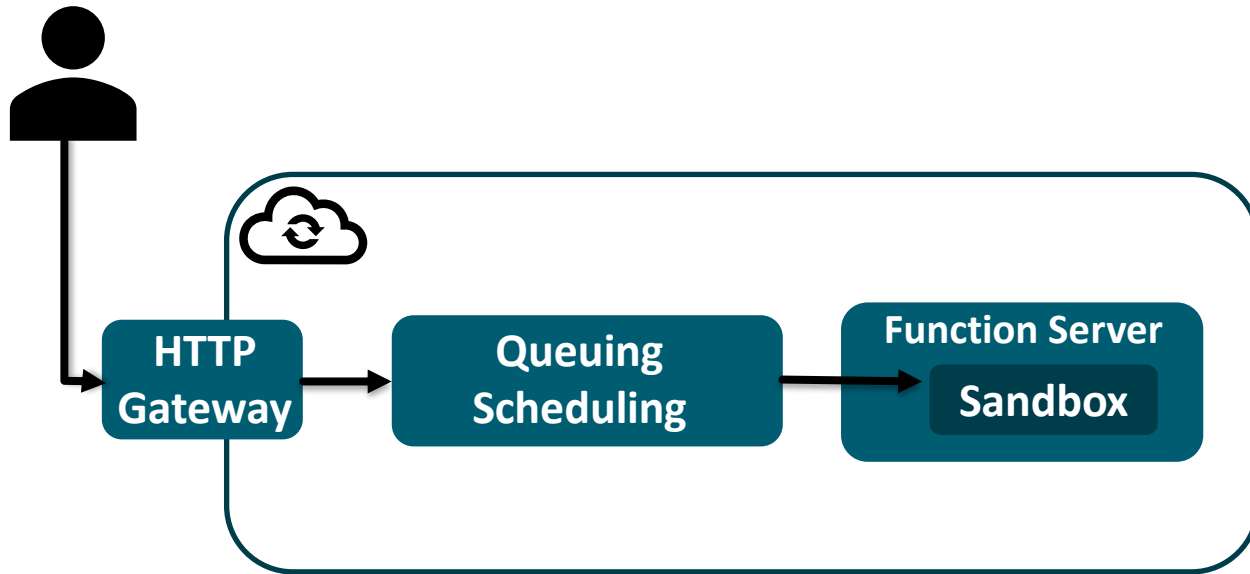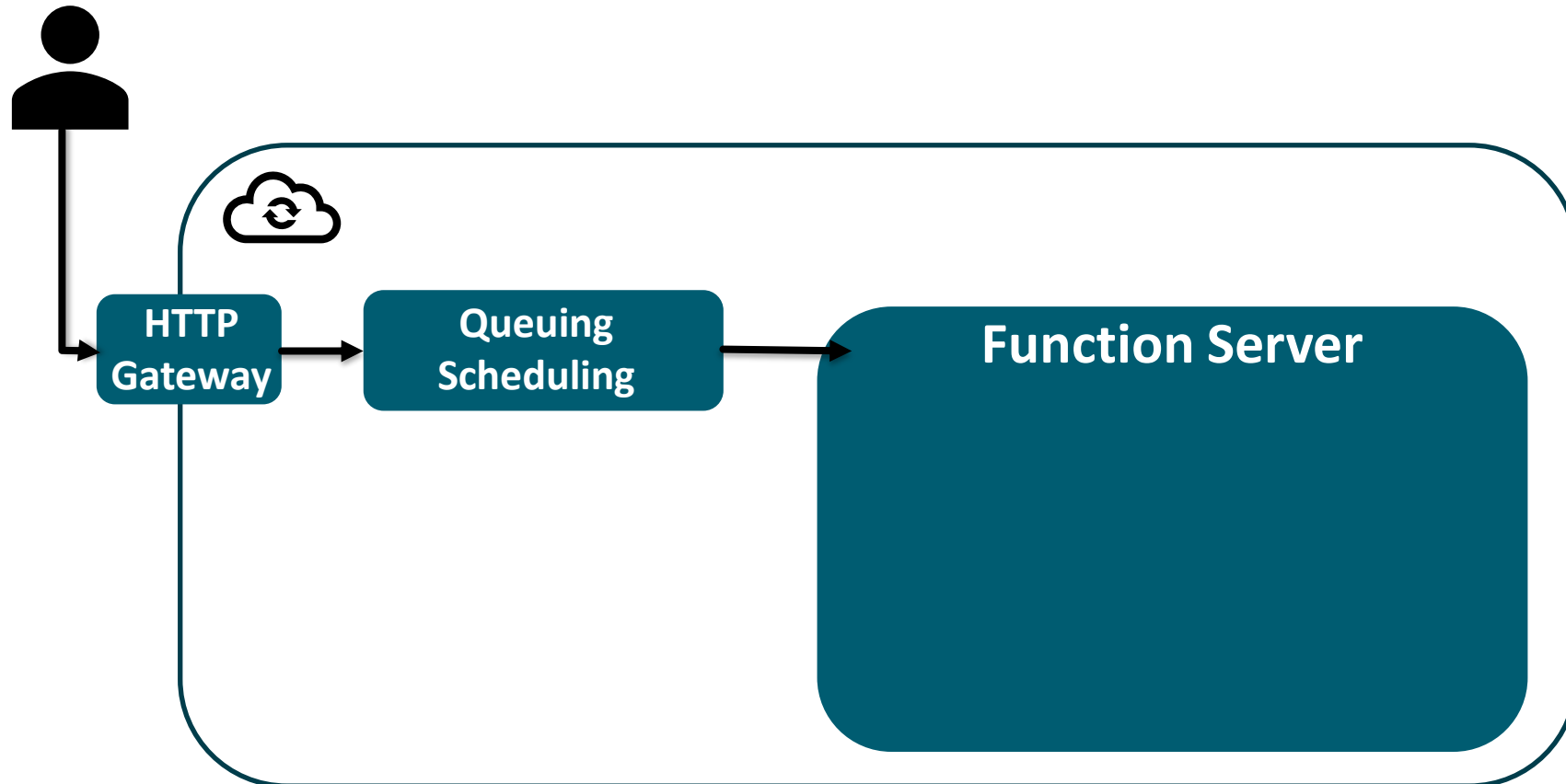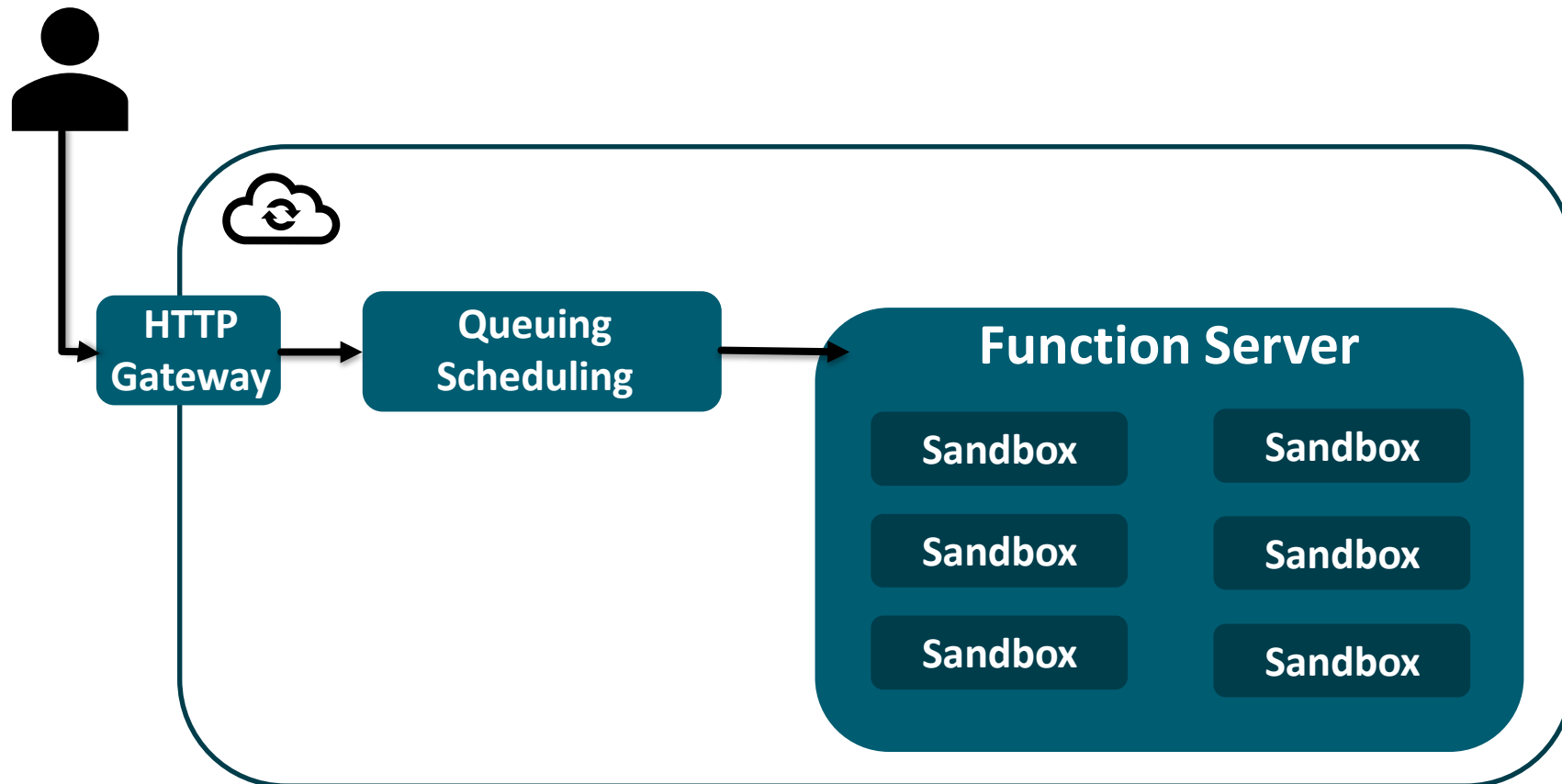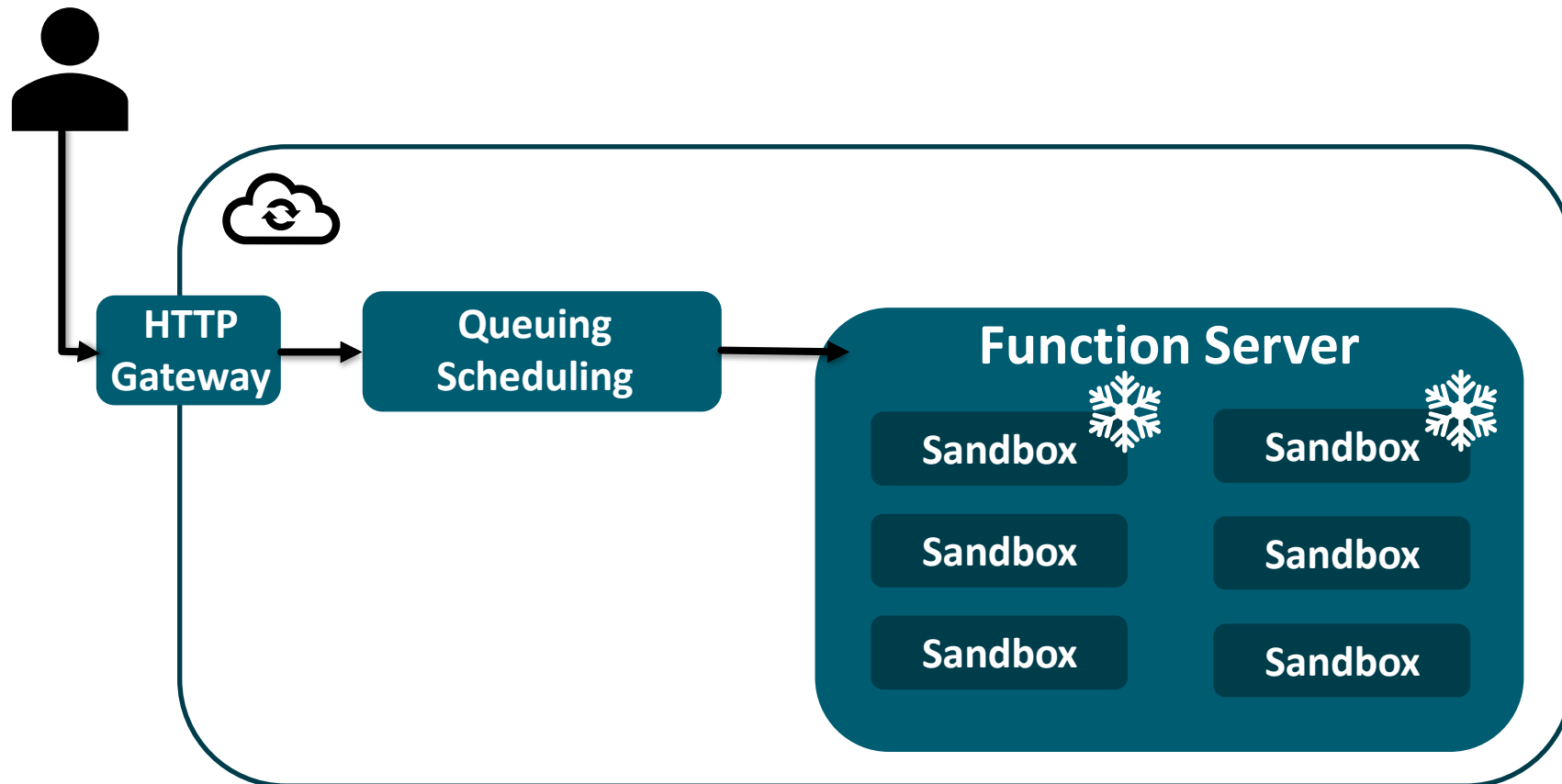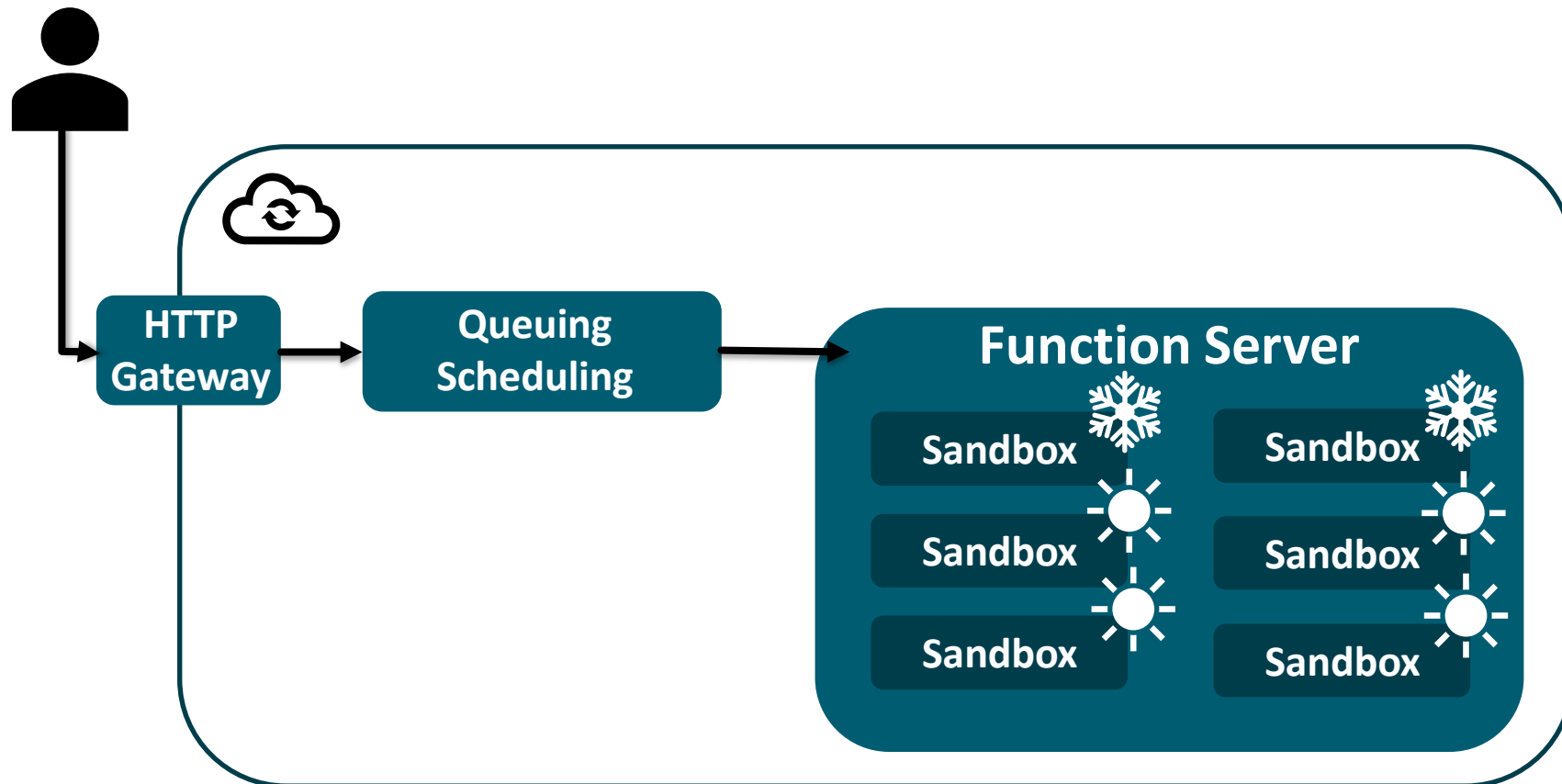
**+**

**Configuration**

# How does serverless systems really work?

# How does serverless systems really work?

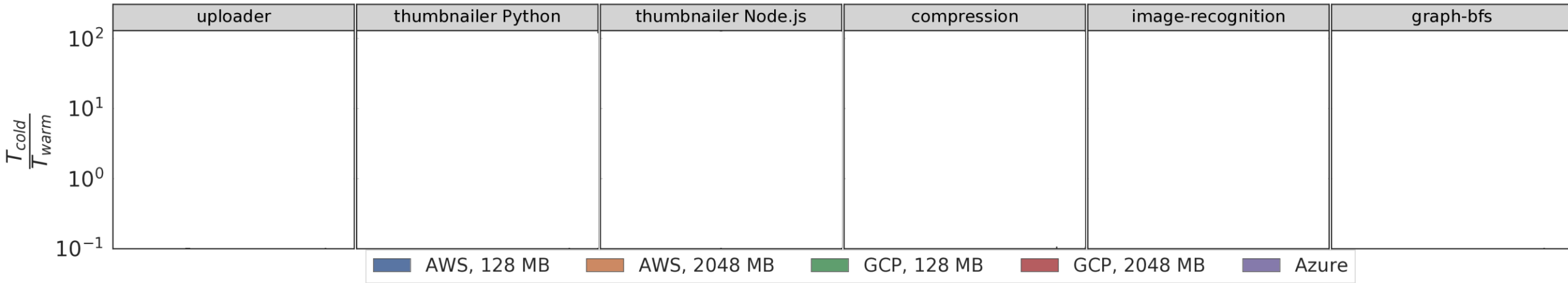# How does serverless systems really work?

# How does serverless systems really work?

# How does serverless systems really work?

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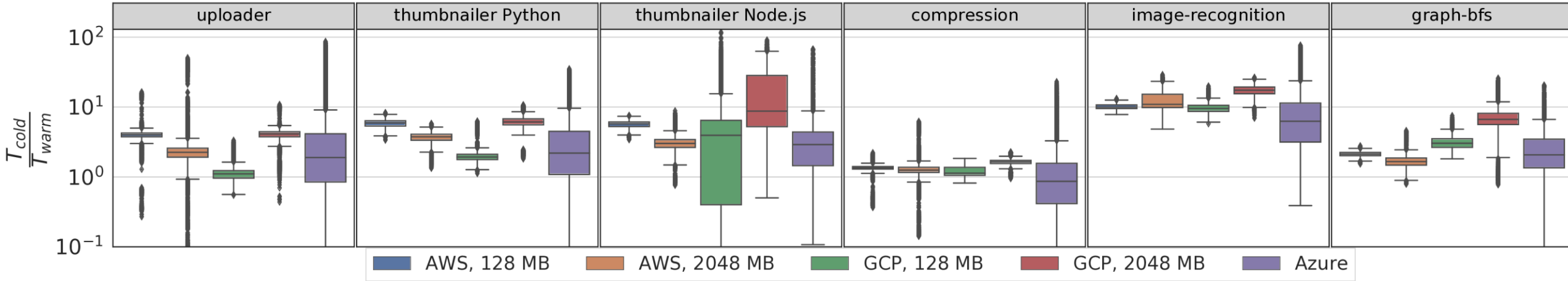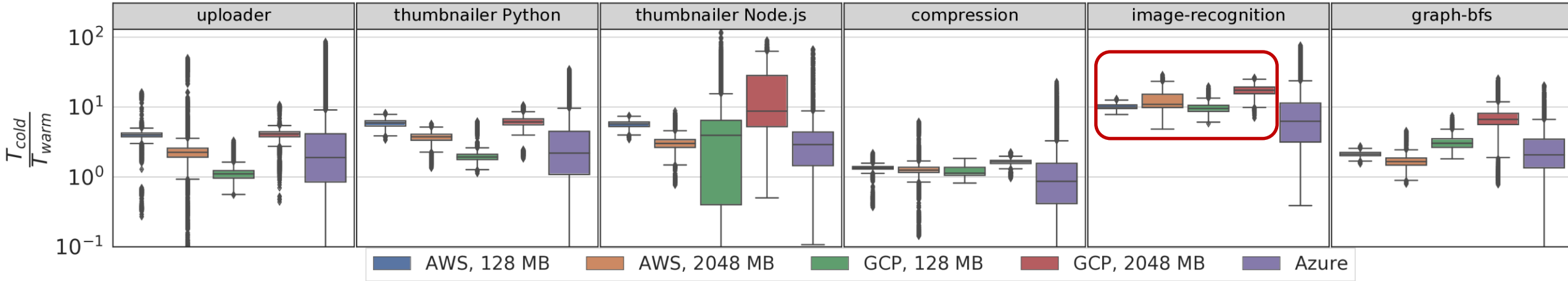
**"SeBS: a Serverless Benchmark Suite for Function-as-a-Service Computing", ACM/IFIP Middleware 2021**

# Cold Startup Overheads

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

# Cold Startup Overheads

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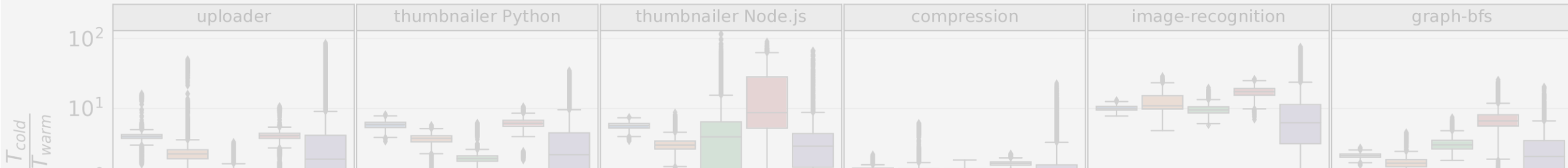
# Cold Startup Overheads



👍 **Faster invocations.**
👍 **Predictable invocation latency.**
👍 **Cacheable state.**

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

# Cold Startup Overheads



👍 **Faster invocations.**

👍 **Predictable invocation latency.**     👎 **Increased memory consumption.**

👍 **Cacheable state.**

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

# Cold Startup Overheads



$\frac{T_{cold}}{T_{warm}}$

| uploader | thumbnailer Python | thumbnailer Node.js | compression | image-recognition | graph-bfs |

Memory contributes 10% of capital and operational expenditures (MareNostrum, 2013) and 18% of peak power consumption in data center (2018).

👍 Faster invocations.
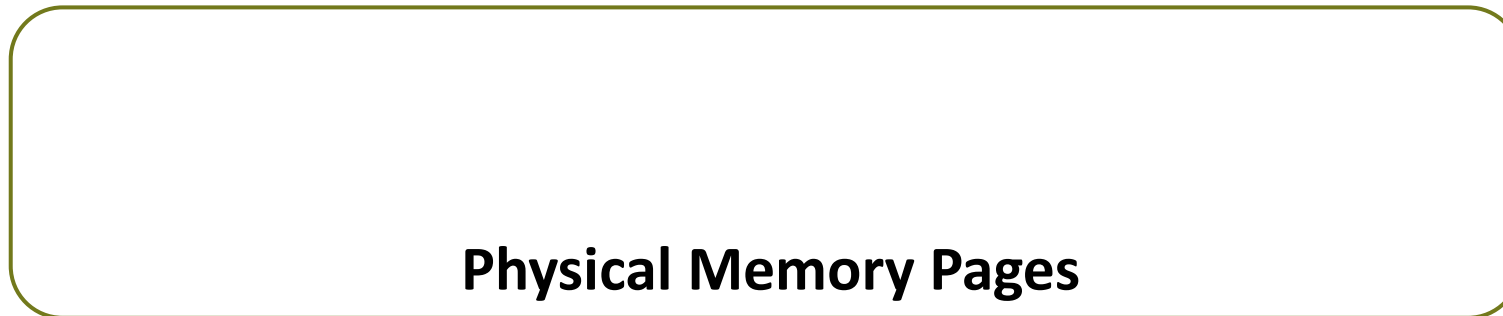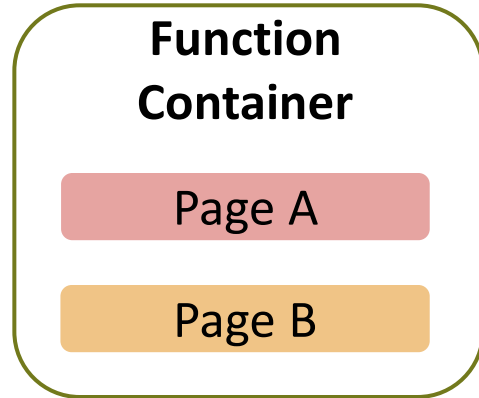👍 Predictable invocation latency.          👎 Increased memory consumption.
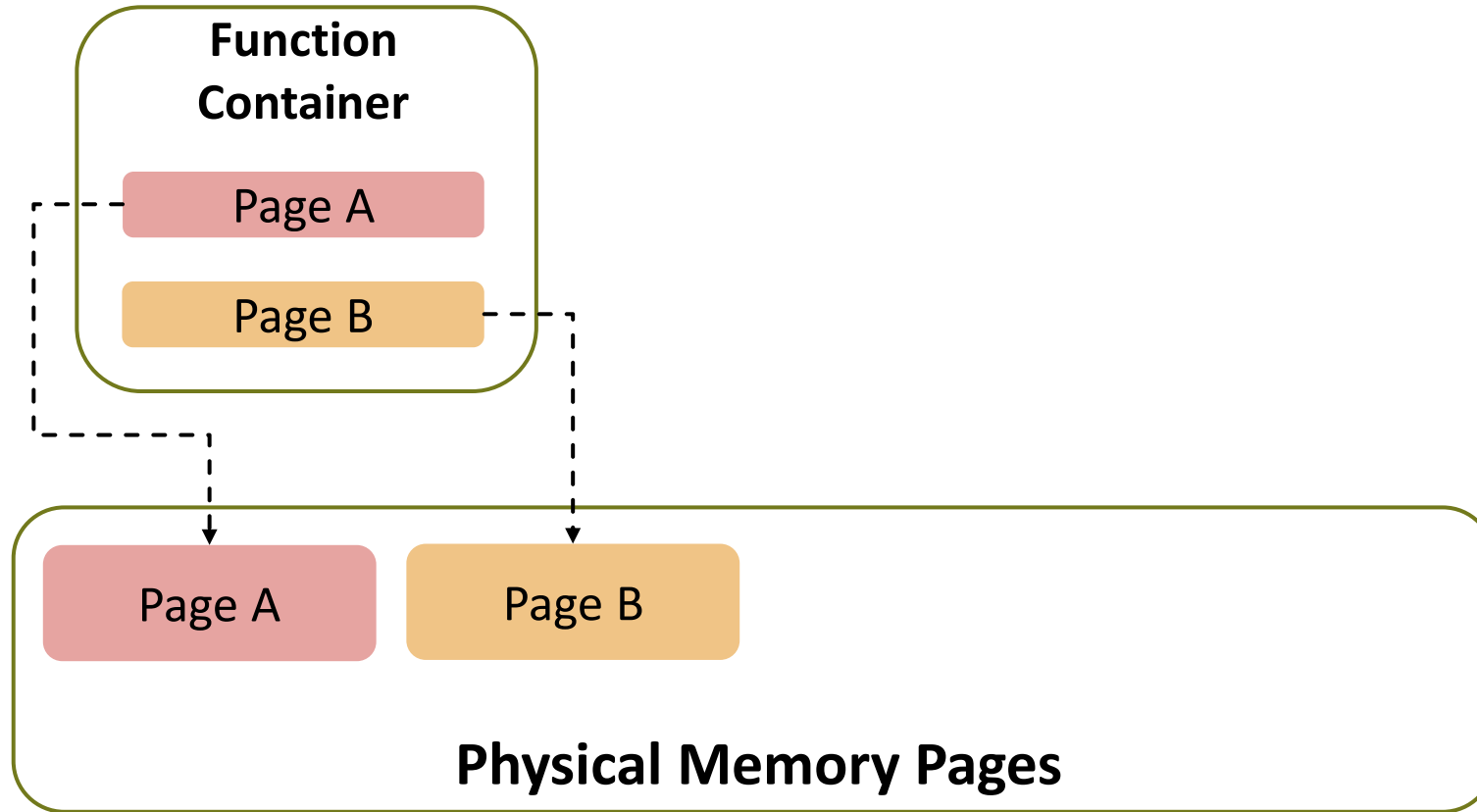👍 Cacheable state.

# Memory Duplication in Serverless

**Physical Memory Pages**

# Memory Duplication in Serverless

**Function Container**
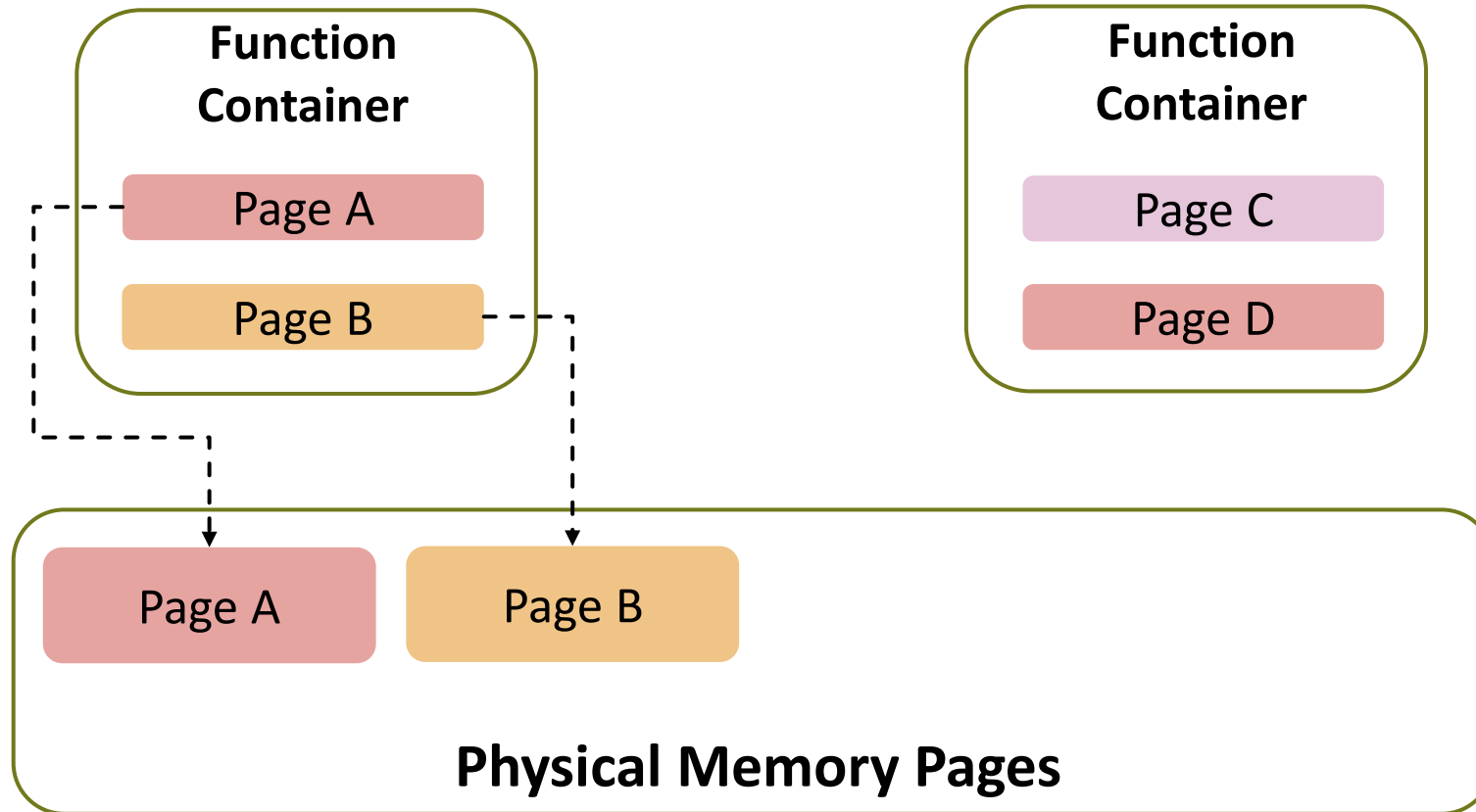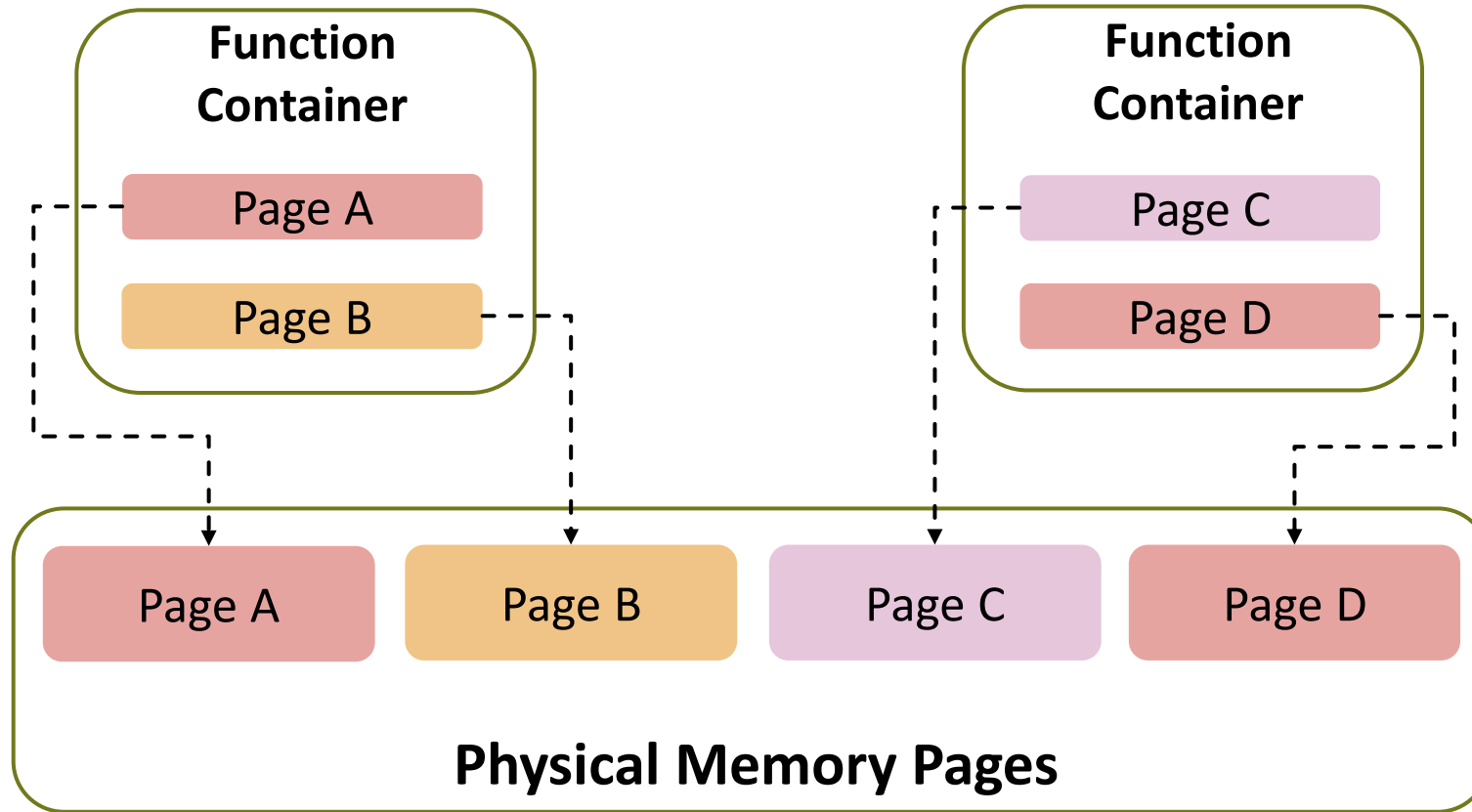
Page A

Page B

**Physical Memory Pages**

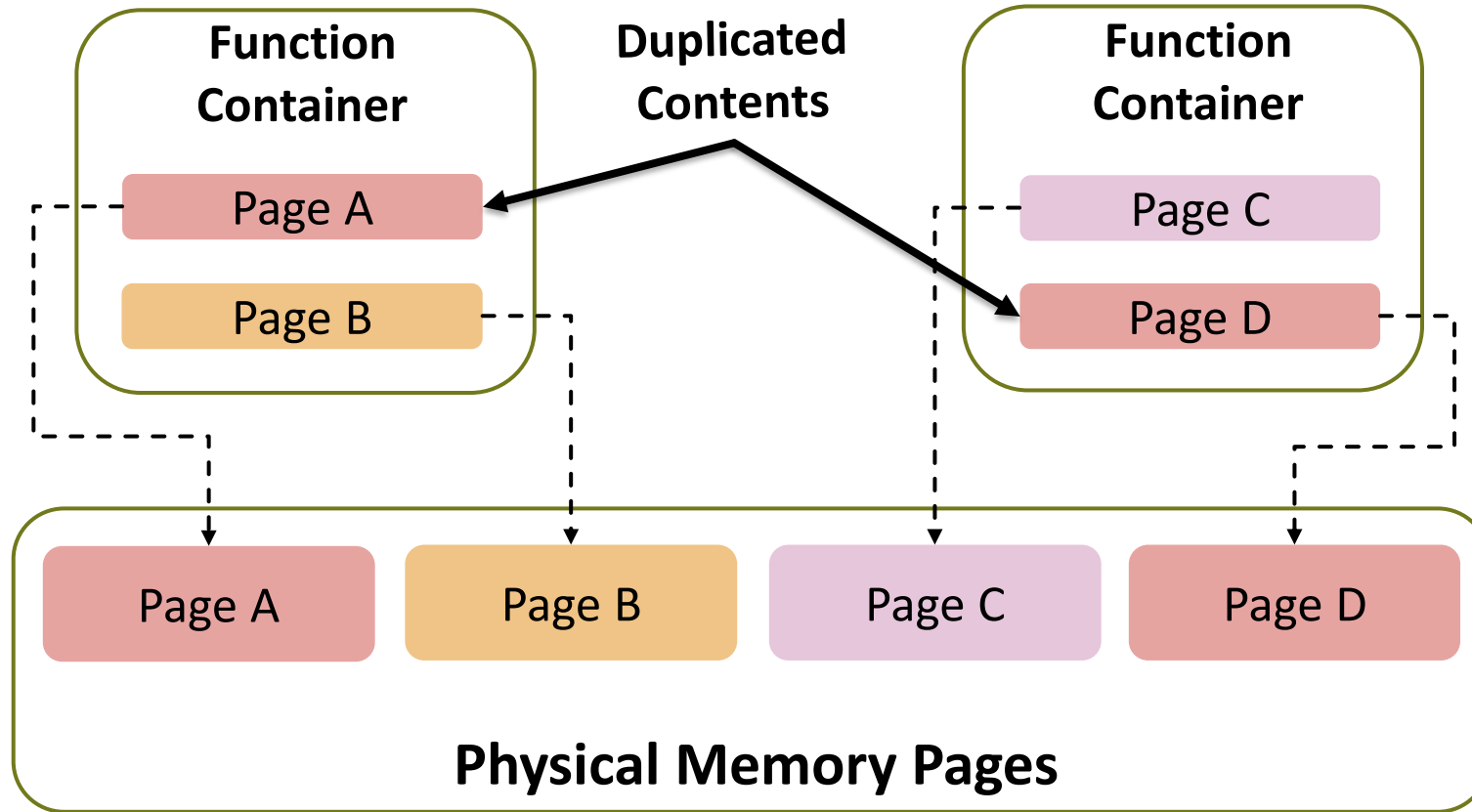# Memory Duplication in Serverless

# Memory Duplication in Serverless
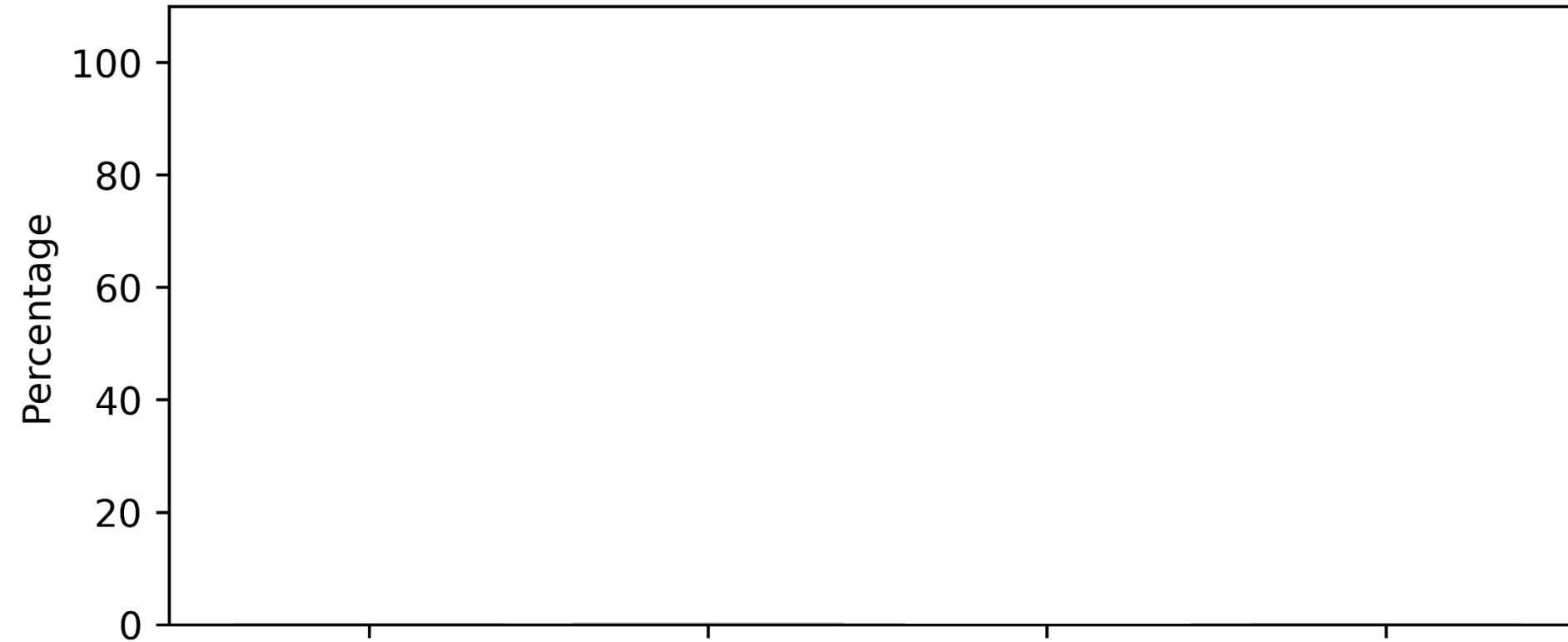
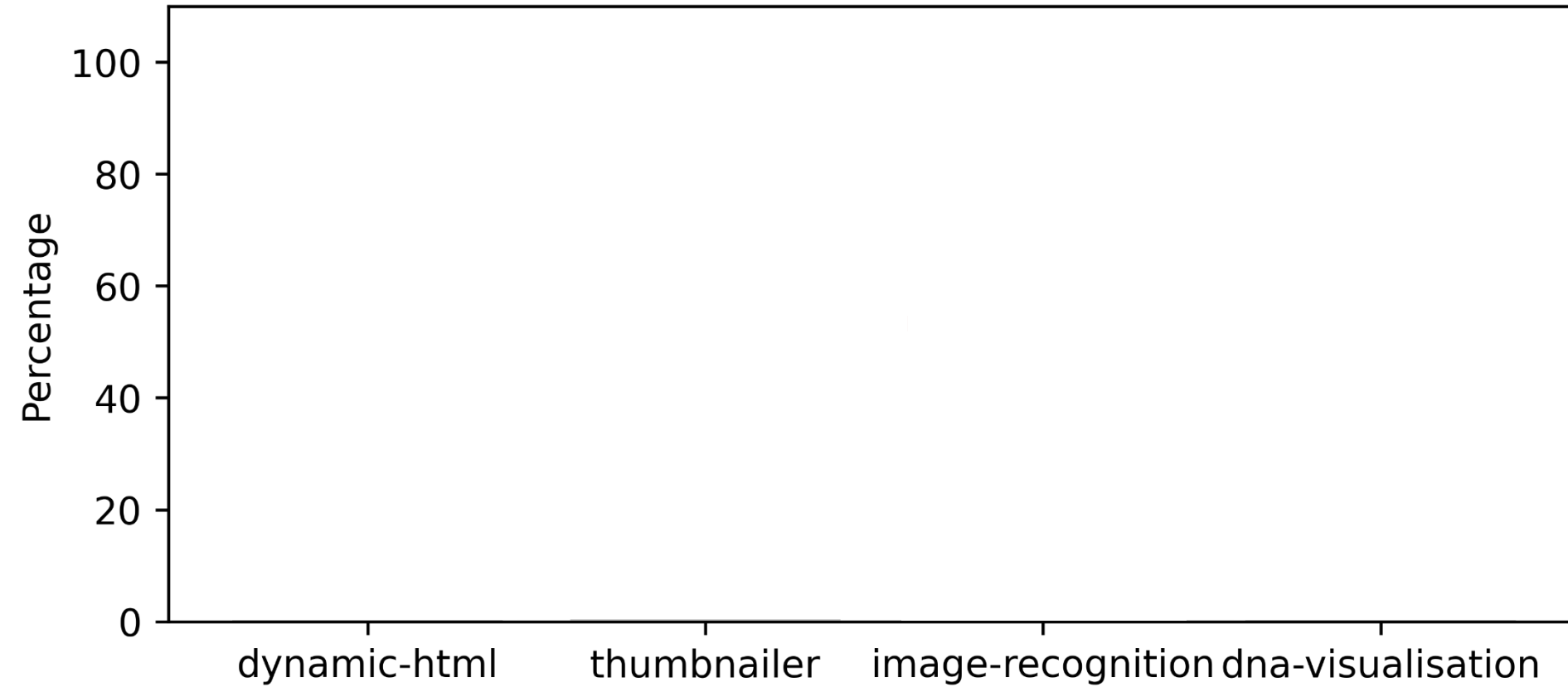# Memory Duplication in Serverless

# Memory Duplication in Serverless

# Sharing opportunities
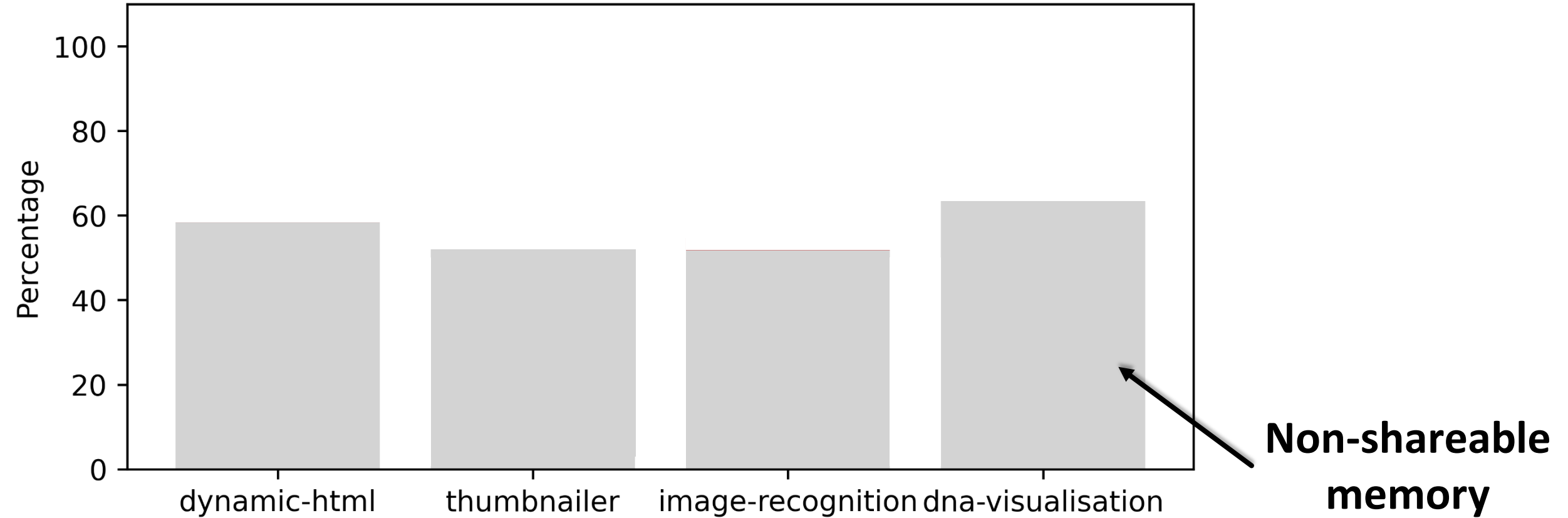
## Memory percentage over RSS

# Sharing opportunities

## Memory percentage over RSS

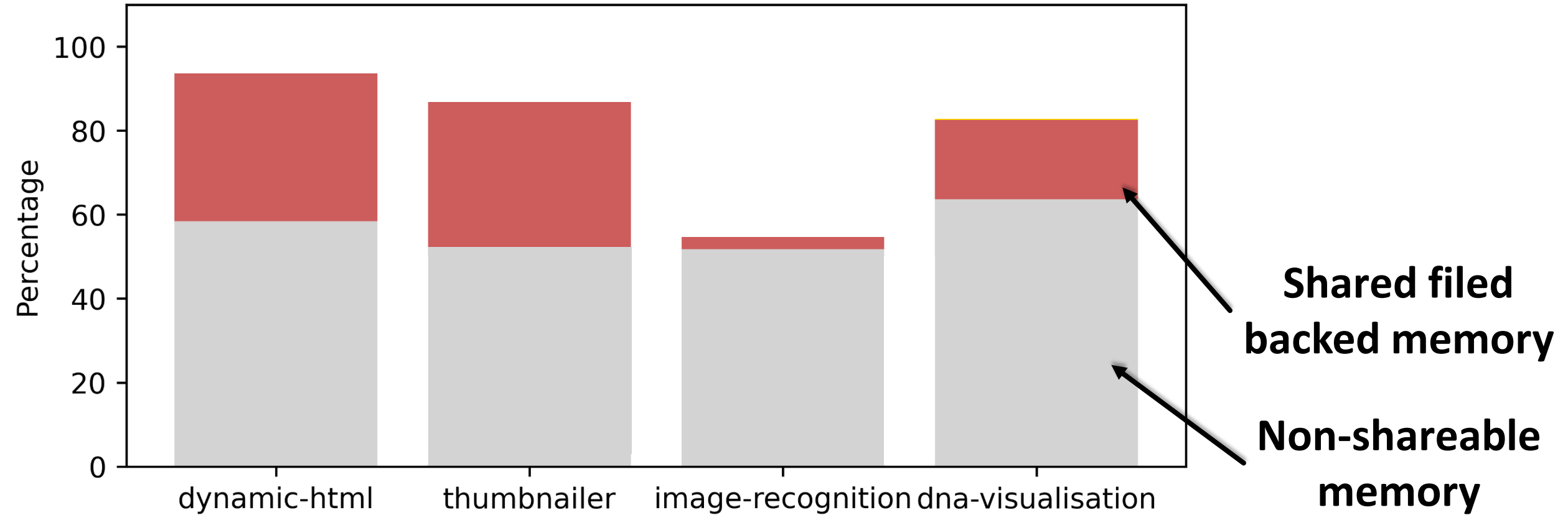**"SeBS: a Serverless Benchmark Suite for Function-as-a-Service Computing", ACM/IFIP Middleware 2021**

# Sharing opportunities



Memory percentage over RSS

**Non-shareable memory**

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

# Sharing opportunities



Memory percentage over RSS

**Shared filed backed memory**

**Non-shareable memory**

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

# Sharing opportunities



Memory percentage over RSS

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

# Sharing opportunities



Memory percentage over RSS

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

# Sharing opportunities



Memory percentage over RSS

Shareable file backed memory

Shareable memory

Shared filed backed memory

Non-shareable memory

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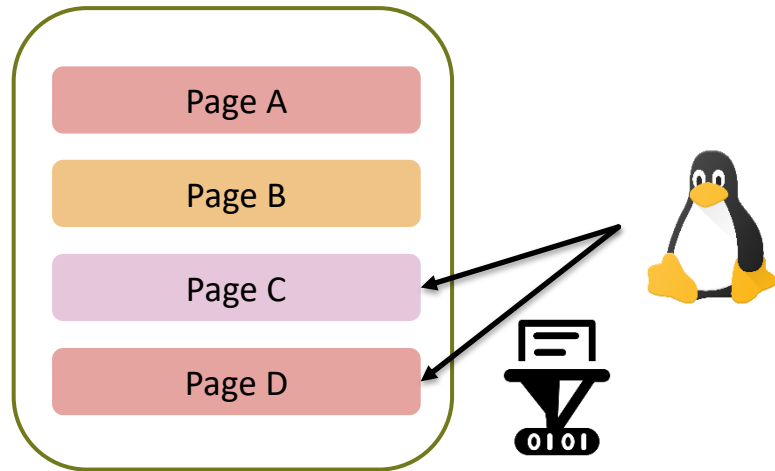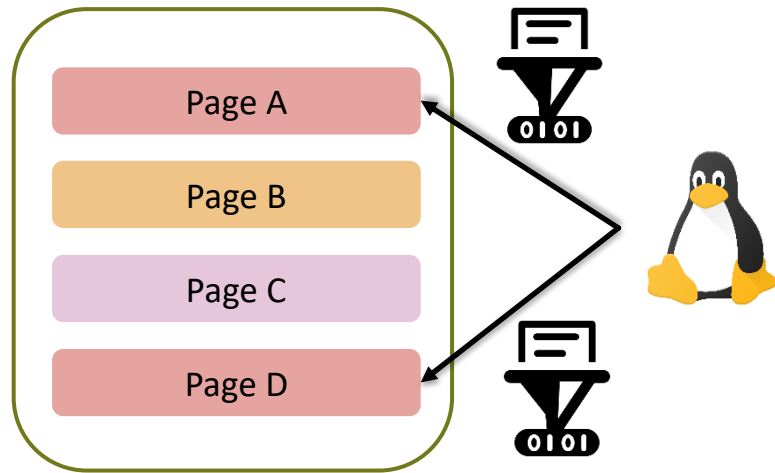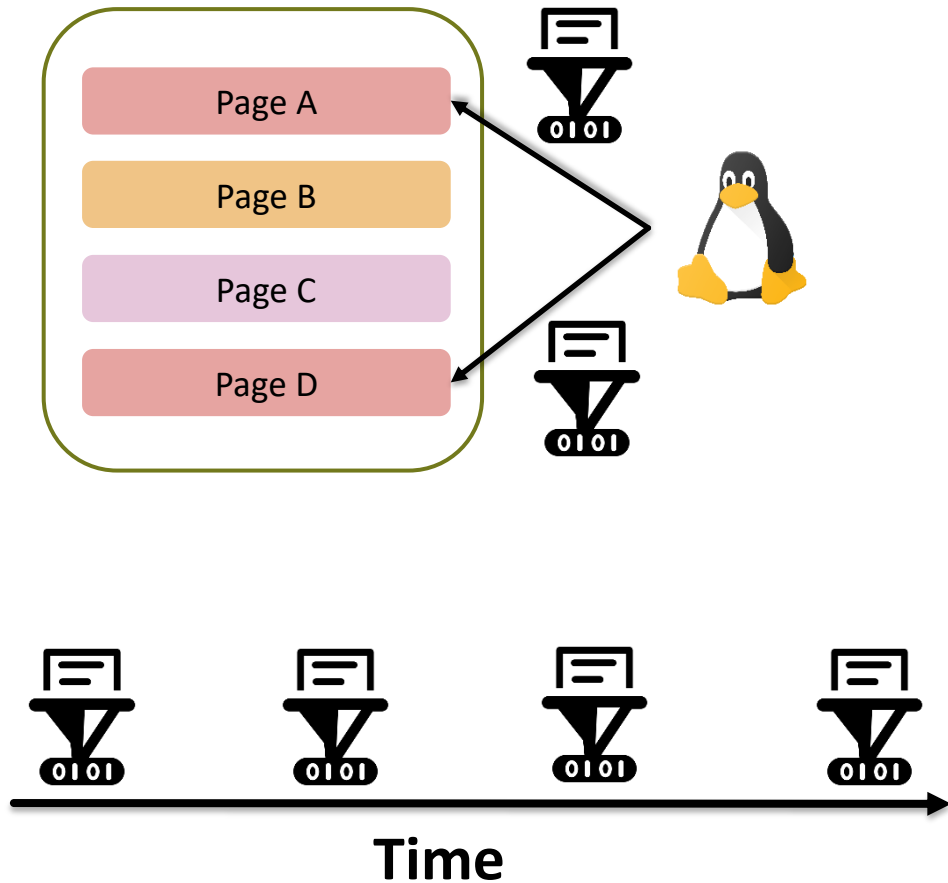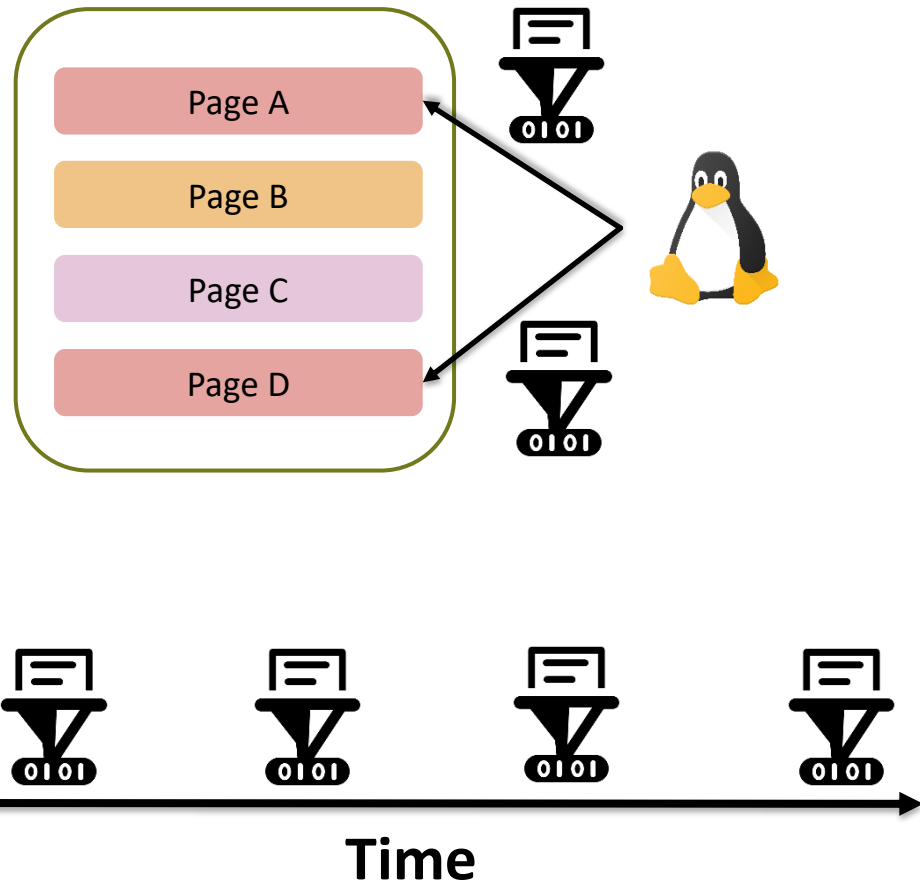
# Existing Memory Deduplication Techniques

**Example**: Kernel Samepage Merging (KSM)

# Existing Memory Deduplication Techniques

**Example**: Kernel Samepage Merging (KSM)

# Existing Memory Deduplication Techniques

**Example**: Kernel Samepage Merging (KSM)

# Existing Memory Deduplication Techniques

**Example**: Kernel Samepage Merging (KSM)

# Existing Memory Deduplication Techniques

**Example**: Kernel Samepage Merging (KSM)

# Existing Memory Deduplication Techniques

**Example**: Kernel Samepage Merging (KSM)



**Time**

# Existing Memory Deduplication Techniques

**Example**: Kernel Samepage Merging (KSM)



Page A

Page B

Page C

Page D

**Satori: Enlightened page sharing**

Grzegorz Miłoś, Derek G. Murray, Steven Hand
University of Cambridge Computer Laboratory
Cambridge, United Kingdom
First.Last@cl.cam.ac.uk

Michael A. Fetterman
NVIDIA Corporation
Bedford, Massachusetts, USA
mafetter@nvidia.com

*"For example in VMware ESX Server the default memory scan frequency is set to once an hour, with a maximum of six times per hour. Therefore, the theoretical **mean duplicate discovery time for the default setting is 40min**, which means that short-lived sharing opportunities will be missed."*

**Time**

# Existing Memory Deduplication Techniques

## Satori: Enlightened page sharing

Grzegorz Miłoś, Derek G. Murray, Steven Hand
*University of Cambridge Computer Laboratory*
*Cambridge, United Kingdom*
First.Last@cl.cam.ac.uk

Michael A. Fetterman
*NVIDIA Corporation*
*Bedford, Massachusetts, USA*
mafetter@nvidia.com

*"For example in VMware ESX Server the default memory scan frequency is set to once an hour, with a maximum of six times per hour. Therefore, the theoretical **mean duplicate discovery time for the default setting is 40min**, which means that short-lived sharing opportunities will be missed."*

# Existing Memory Deduplication Techniques

**Serverless in the Wild: Characterizing and Optimizing the Serverless Workload at a Large Cloud Provider**

Mohammad Shahrad, Rodrigo Fonseca, Íñigo Goiri, Gohar Chaudhry, Paul Batum, Jason Cooke, Eduardo Laureano, Colby Tresness, Mark Russinovich, and Ricardo Bianchini, *Microsoft Azure and Microsoft Research*

*Satori: Enlightened page sharing*

Grzegorz Miłoś, Derek G. Murray, Steven Hand          Michael A. Fetterman
*University of Cambridge Computer Laboratory*               *NVIDIA Corporation*
*Cambridge, United Kingdom*                                   *Bedford, Massachusetts, USA*
First.Last@cl.cam.ac.uk                               mafetter@nvidia.com

*"We observe that 50% of the functions execute for less than 1s on average, and 50% of the functions have maximum execution time shorter than ~3s; 90% of the functions take at most 60s, and **96% of functions take less than 60s on average.**"*

*"For example in VMware ESX Server the default memory scan frequency is set to once an hour, with a maximum of six times per hour. Therefore, the theoretical **mean duplicate discovery time for the default setting is 40min**, which means that short-lived sharing opportunities will be missed."*

# Existing Memory Deduplication Techniques

## Serverless in the Wild: Characterizing and Optimizing the Serverless Workload at a Large Cloud Provider

Mohammad Shahrad, Rodrigo Fonseca, Íñigo Goiri, Gohar Chaudhry, Paul Batum, Jason Cooke, Eduardo Laureano, Colby Tresness, Mark Russinovich, and Ricardo Bianchini, *Microsoft Azure and Microsoft Research*

## SeBS: A Serverless Benchmark Suite for Function-as-a-Service Computing

Marcin Copik
marcin.copik@inf.ethz.ch
ETH Zürich
Switzerland

Grzegorz Kwaśniewski
ETH Zürich
Switzerland

Maciej Besta
ETH Zürich
Switzerland

Michał Podstawski
Future Processing SA
Poland

Torsten Hoefler
ETH Zürich
Switzerland

*"We observe that 50% of the functions execute for less than 1s on average, and 50% of the functions have maximum execution time shorter than  ~3s; 90% of the functions take at most 60s, and **96% of functions take less than 60s on average.**"*

Number of idle function containers is **reduced every 380 seconds** on a major commercial serverless provider.

# Memory Deduplication in Serverless

# Memory Deduplication in Serverless

**Speed**
Deduplication in seconds,
not minutes.

# Memory Deduplication in Serverless

**Speed**
Deduplication in seconds,
not minutes.

**Compatibility**
No changes to existing runtimes.

# Memory Deduplication in Serverless

**Speed**
Deduplication in seconds,
not minutes.

**Compatibility**
No changes to existing runtimes.

**Opt-in**
Zero overhead when unused and
no sharing for sensitive data.

# Memory Deduplication in Serverless

**Speed**
Deduplication in seconds,
not minutes.

**Compatibility**
No changes to existing runtimes.

**Opt-in**
Zero overhead when unused and
no sharing for sensitive data.

**Concurrency**
Support many containers
operating concurrently.

# UPM: User-Guided Page Merging

# UPM: User-Guided Page Merging

# UPM: User-Guided Page Merging

# UPM: User-Guided Page Merging

# UPM: User-Guided Page Merging

# UPM: User-Guided Page Merging

# UPM: User-Guided Page Merging

# UPM: User-Guided Page Merging

# UPM: User-Guided Page Merging

# UPM: User-Guided Page Merging

# UPM: User-Guided Page Merging

# UPM: User-Guided Page Merging

# UPM: User-Guided Page Merging

# UPM: User-Guided Page Merging

# UPM Algorithm

# UPM Algorithm



madvise
a memory region.

For each page
in region.

UPM, Kernel space

# UPM Algorithm

madvise
a memory region.

For each page
in region.

Calculate hash
value

Find page with
the same hash

UPM, Kernel space

# UPM Algorithm

# UPM Algorithm

# UPM Algorithm

# UPM Algorithm

# UPM Algorithm

# UPM Algorithm

# Implementation

❖ **Modified Linux kernel 4.15.18**

❖ **UPM is a new built-in kernel module**

❖ **Reuse concepts from Kernel Samepage Merging (KSM)**

# Evaluation

## Bare-Metal Server

# Evaluation

# Evaluation

# Evaluation

# Evaluation



Bare-Metal Server

Hypervisor

Virtual Machine

MINIO

Function

Function

Function

Function

Does UPM decrease memory footprint of functions? **?**

How much system memory can be saved? **?**

# Evaluation

**Bare-Metal Server**

Hypervisor

**Virtual Machine**

MINIO

Function    Function

Function    Function

Does UPM decrease memory footprint of functions? **?**

How much system memory can be saved? **?**

How much overhead does UPM add to serverless? **?**

# Function Memory Footprint – ResNet 50

4x Intel Xeon X7550 @ 2.00GHz,
64 cores total. 1 TB memory.



FaaS without UPM.

FaaS with UPM.

RSS = Private + Shared

$$PSS = Private + \frac{Shared}{\#Processes}$$

19

# Function Memory Footprint – ResNet 50

FaaS without UPM.

FaaS with UPM.

Legend:
- RSS
- PSS
- Private

$RSS = Private + Shared$

$PSS = Private + \dfrac{Shared}{\#Processes}$

# Function Memory Footprint – ResNet 50

FaaS without UPM.

FaaS with UPM.

RSS = Private + Shared

$$PSS = Private + \frac{Shared}{\#Processes}$$

# Function Memory Footprint – ResNet 50

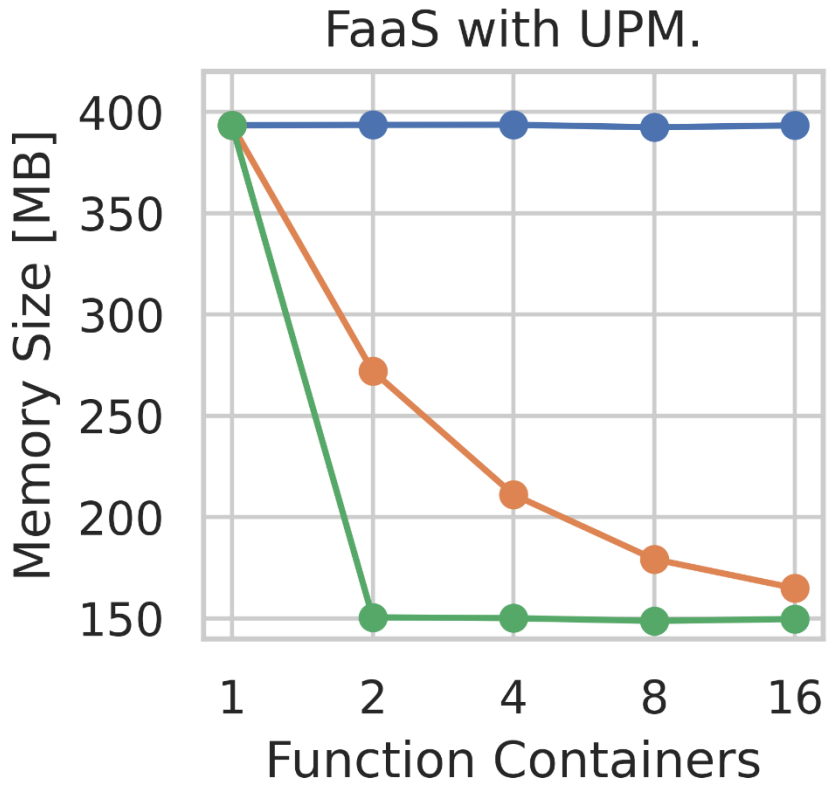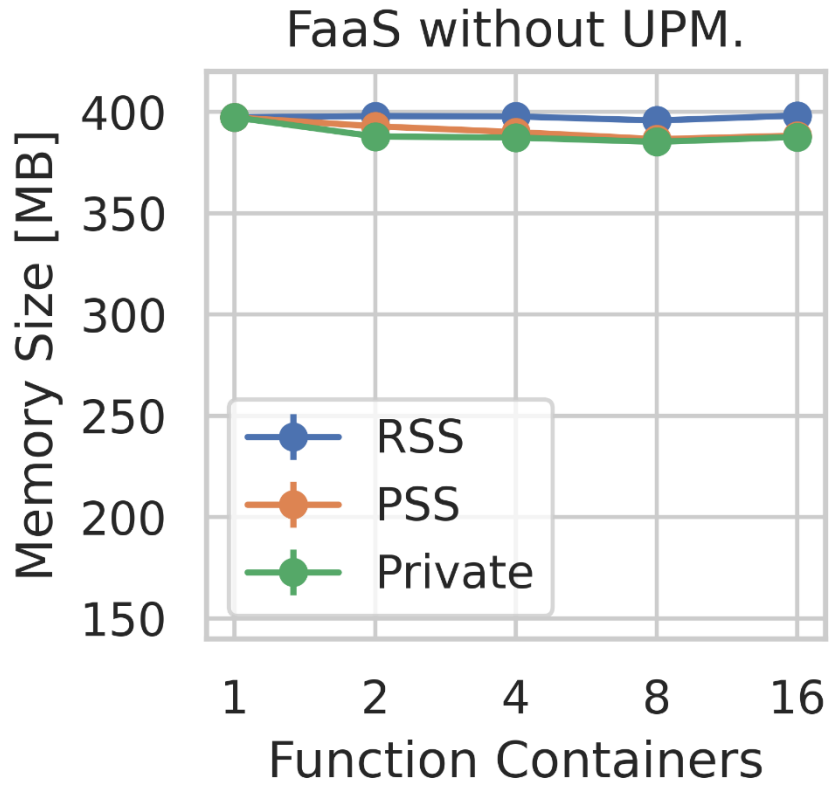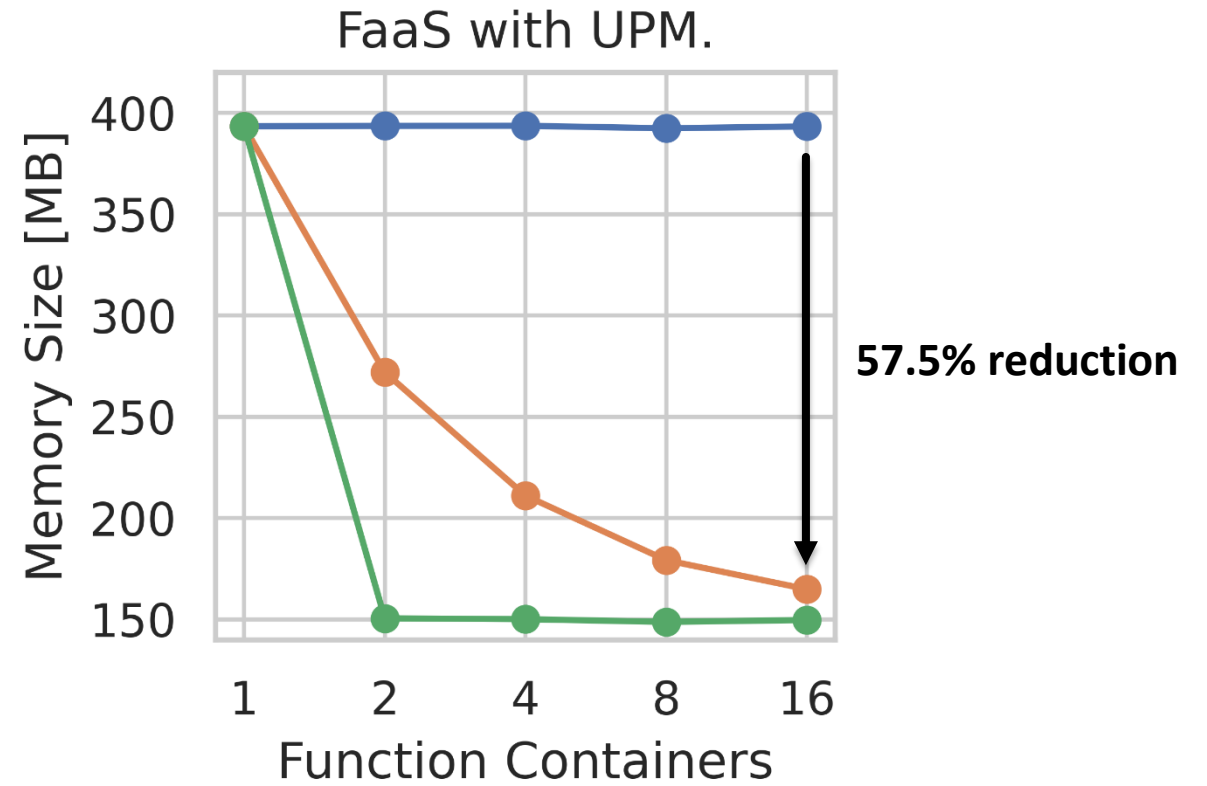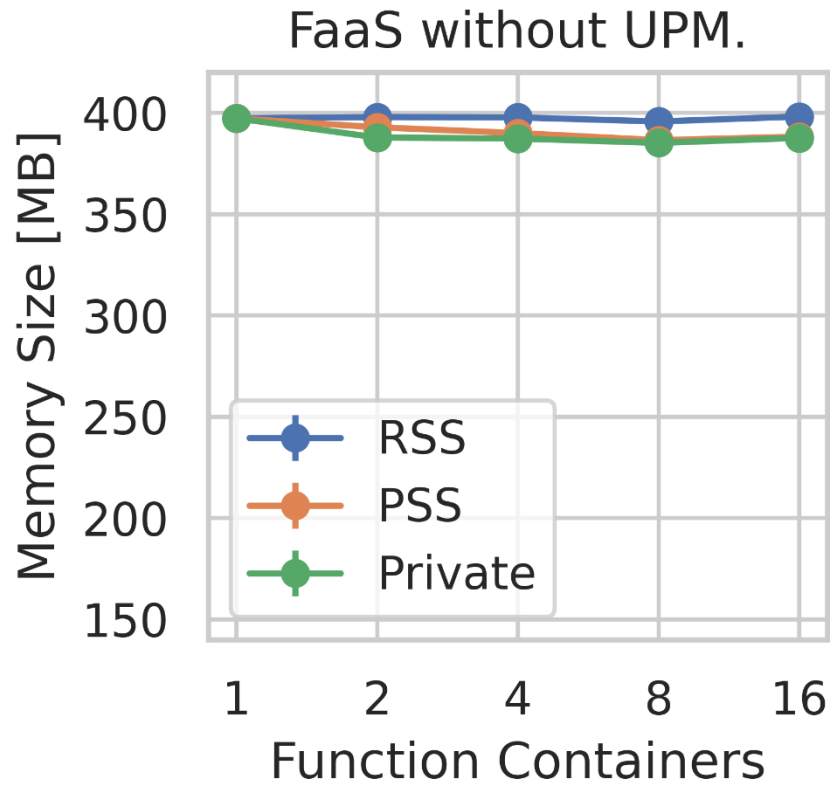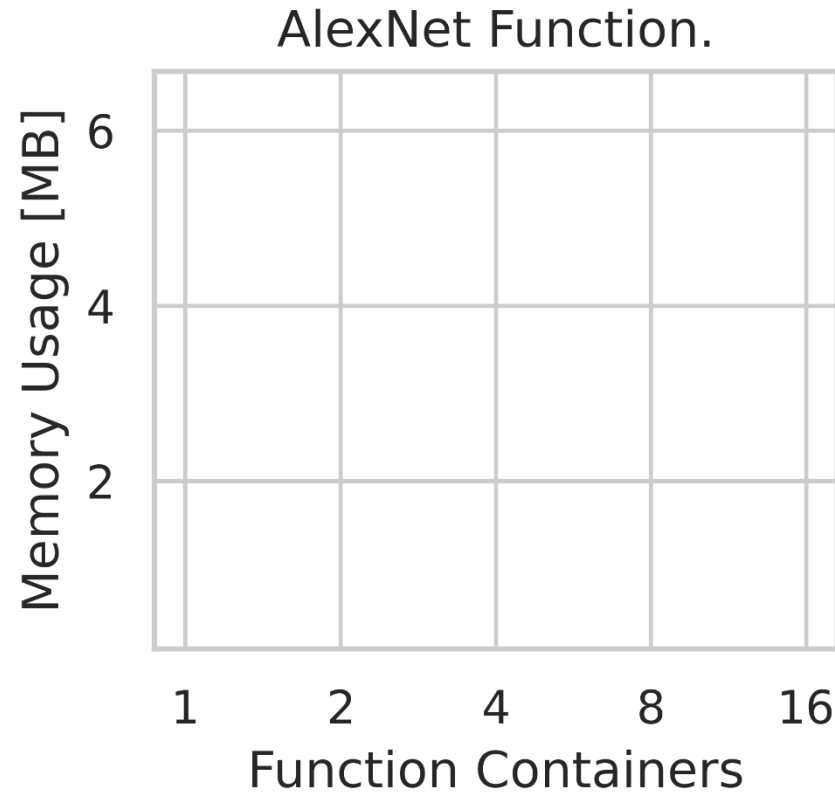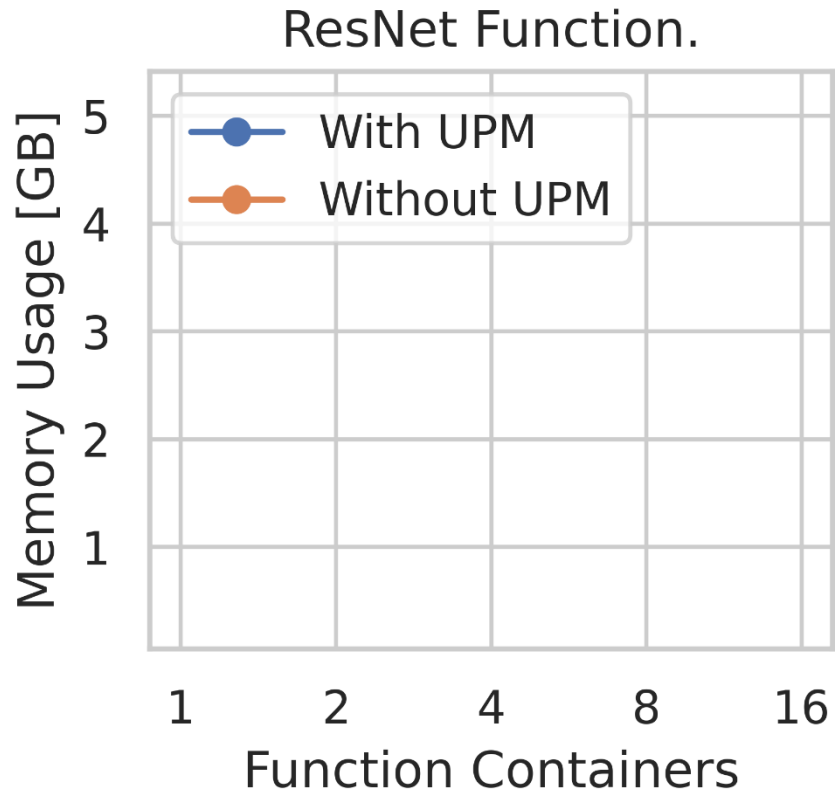4x Intel Xeon X7550 @ 2.00GHz,
64 cores total. 1 TB memory.

FaaS without UPM.

FaaS with UPM.



RSS = Private + Shared

$$PSS = Private + \frac{Shared}{\#Processes}$$

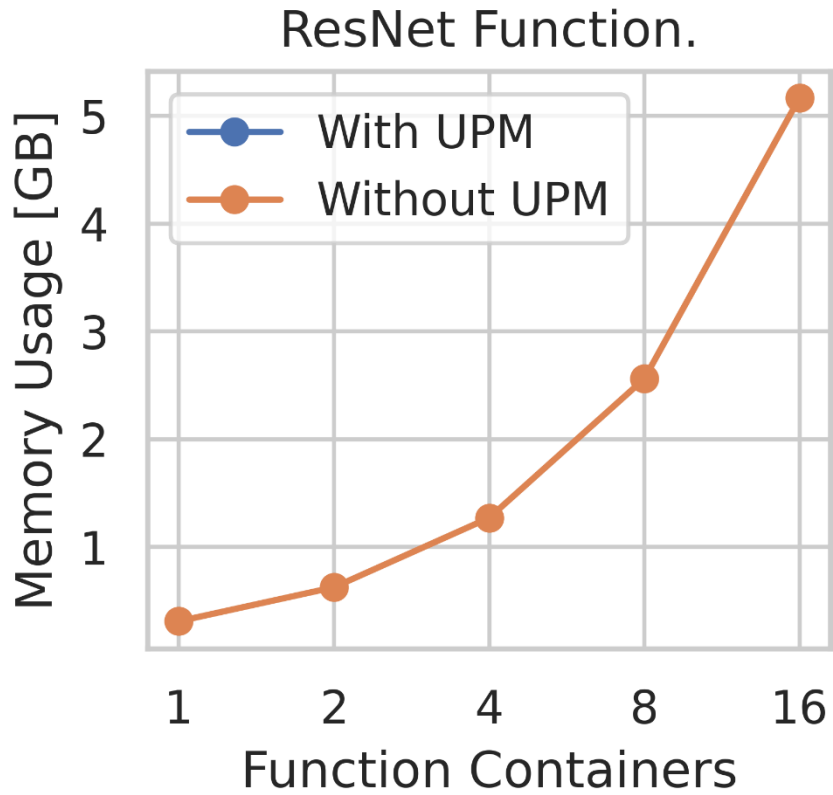# Function Memory Footprint – ResNet 50

4x Intel Xeon X7550 @ 2.00GHz, 64 cores total. 1 TB memory.



FaaS without UPM.

FaaS with UPM.

**23.5% reduction**

$RSS = Private + Shared$

$PSS = Private + \dfrac{Shared}{\#Processes}$

# Function Memory Footprint – AlexNet

## FaaS without UPM.



Memory Size [MB]: 400, 350, 300, 250, 200, 150

Function Containers: 1, 2, 4, 8, 16

Legend: RSS, PSS, Private

## FaaS with UPM.



Memory Size [MB]: 400, 350, 300, 250, 200, 150

Function Containers: 1, 2, 4, 8, 16

$$RSS = Private + Shared$$

$$PSS = Private + \frac{Shared}{\#Processes}$$

# Function Memory Footprint – AlexNet

FaaS without UPM.

FaaS with UPM.

RSS = Private + Shared

$$PSS = Private + \frac{Shared}{\#Processes}$$

# Function Memory Footprint – AlexNet

FaaS without UPM.

FaaS with UPM.

Legend:
- RSS
- PSS
- Private

$RSS = Private + Shared$

$PSS = Private + \dfrac{Shared}{\#Processes}$

# Function Memory Footprint – AlexNet

4x Intel Xeon X7550 @ 2.00GHz,
64 cores total. 1 TB memory.



FaaS without UPM.

FaaS with UPM.

RSS = Private + Shared

PSS = Private + $\frac{\text{Shared}}{\#\text{Processes}}$

# Function Memory Footprint – AlexNet

4x Intel Xeon X7550 @ 2.00GHz, 64 cores total. 1 TB memory.



**57.5% reduction**

RSS = Private + Shared

PSS = Private + $\frac{\text{Shared}}{\text{\#Processes}}$

# System Memory Consumption

**4x Intel Xeon X7550 @ 2.00GHz, 64 cores total. 1 TB memory.**



ResNet Function.

AlexNet Function.

# System Memory Consumption

4x Intel Xeon X7550 @ 2.00GHz, 64 cores total. 1 TB memory.



ResNet Function.

AlexNet Function.

# System Memory Consumption

4x Intel Xeon X7550 @ 2.00GHz, 64 cores total. 1 TB memory.



ResNet Function. / AlexNet Function.

# System Memory Consumption

4x Intel Xeon X7550 @ 2.00GHz, 64 cores total. 1 TB memory.



ResNet Function.

AlexNet Function.

21% reduction

55% reduction

21

# System Memory Consumption

4x Intel Xeon X7550 @ 2.00GHz, 64 cores total. 1 TB memory.



**ResNet Function.**

With UPM
Without UPM

**21% reduction**

**5 more containers**

**AlexNet Function.**

**55% reduction**

**21 more containers**

# Cold Startup Overhead – ResNet 50

2x Intel Xeon 4110 @ 2.10GHz, 16 cores total. 125 GB memory.

# Cold Startup Overhead – ResNet 50

2x Intel Xeon 4110 @ 2.10GHz, 16 cores total. 125 GB memory.



**Function**

# Cold Startup Overhead – ResNet 50

2x Intel Xeon 4110 @ 2.10GHz, 16 cores total. 125 GB memory.



22

# Cold Startup Overhead – ResNet 50

2x Intel Xeon 4110 @ 2.10GHz, 16 cores total. 125 GB memory.



Total

Function

madvise

# Cold Startup Overhead – AlexNet

# Cold Startup Overhead – AlexNet

**Function**

# Cold Startup Overhead – AlexNet

2x Intel Xeon 4110 @ 2.10GHz, 16 cores total. 125 GB memory.



**Function madvise**

23

# Cold Startup Overhead – AlexNet

2x Intel Xeon 4110 @ 2.10GHz, 16 cores total. 125 GB memory.



**Total**

**Function madvise**

23

# UPM at GitHub

# UPM at GitHub

spcl/UPM

spcl/UPM-artifact-data

# High-Performance Serverless Solutions

# High-Performance Serverless Solutions

 spcl/serverless-benchmarks

 spcl/rFaaS

 spcl/FMI

 spcl/PraaS

 spcl/FaaSKeeper

# Conclusions

# Conclusions

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

```
def handler_function(req: dict, context: dict):
    model = cloud_storage.download_model()
    input = parse_input(req['payload'])
    output = model.inference(input)
    return output
```

Configuration

HTTP Gateway → Queuing Scheduling → Function Server / Sandbox

# Conclusions

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



Sharing opportunities

# Conclusions

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



Sharing opportunities



UPM: User-Guided Page Merging

# Conclusions

# Conclusions

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



Sharing opportunities



UPM: User-Guided Page Merging



System Memory Consumption

Paper preprint     Other projects

CSCS
Centro Svizzero di Calcolo Scientifico
Swiss National Supercomputing Centre