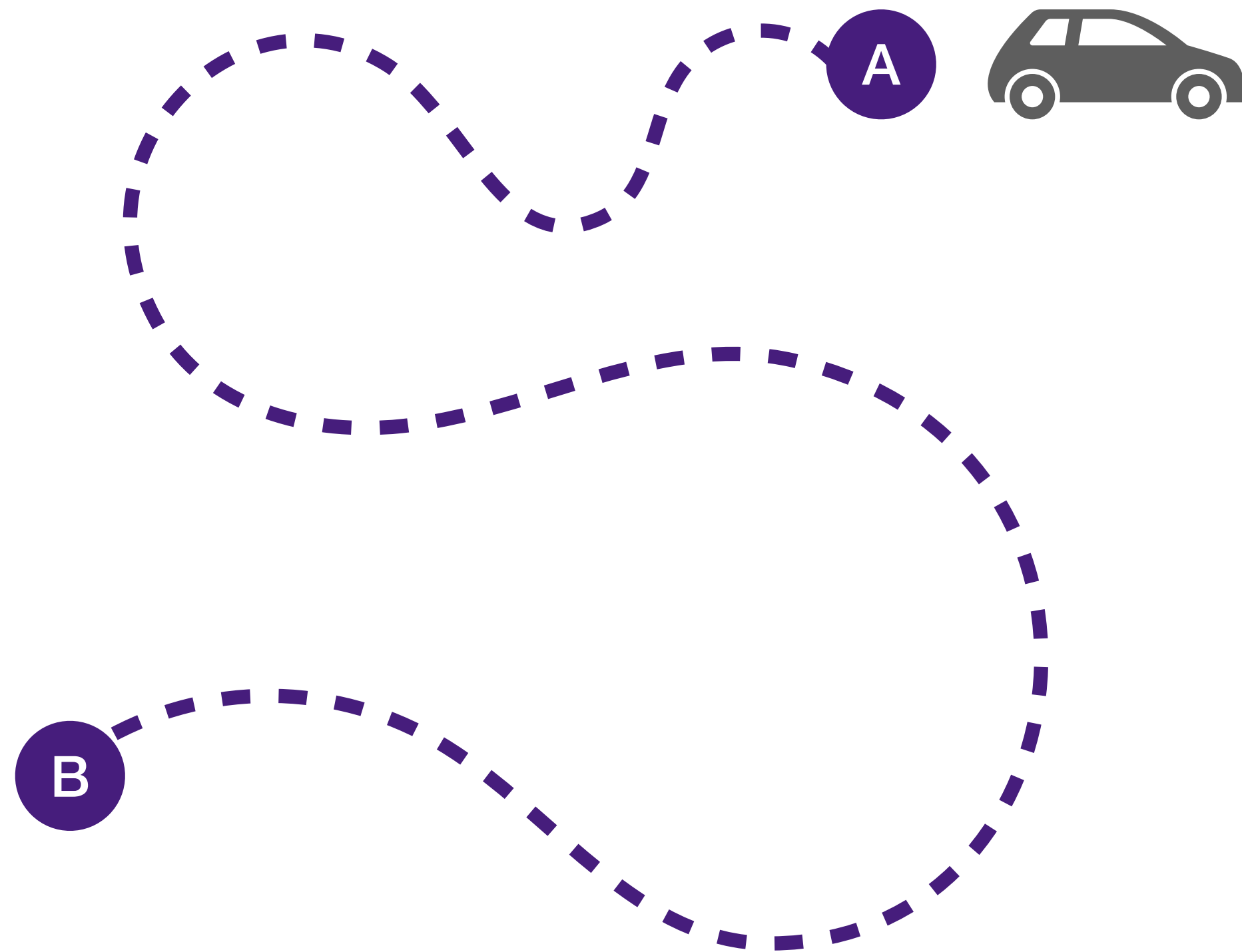# LIBRA: Harvesting Idle Resources Safely and Timely in Serverless Clusters

Hanfei Yu, Christian Fontenot, Hao Wang, Jian Li*, Xu Yuan†, Seung-Jong Park

Louisiana State University, SUNY-Binghamton University*, University of Louisiana at Lafayette†
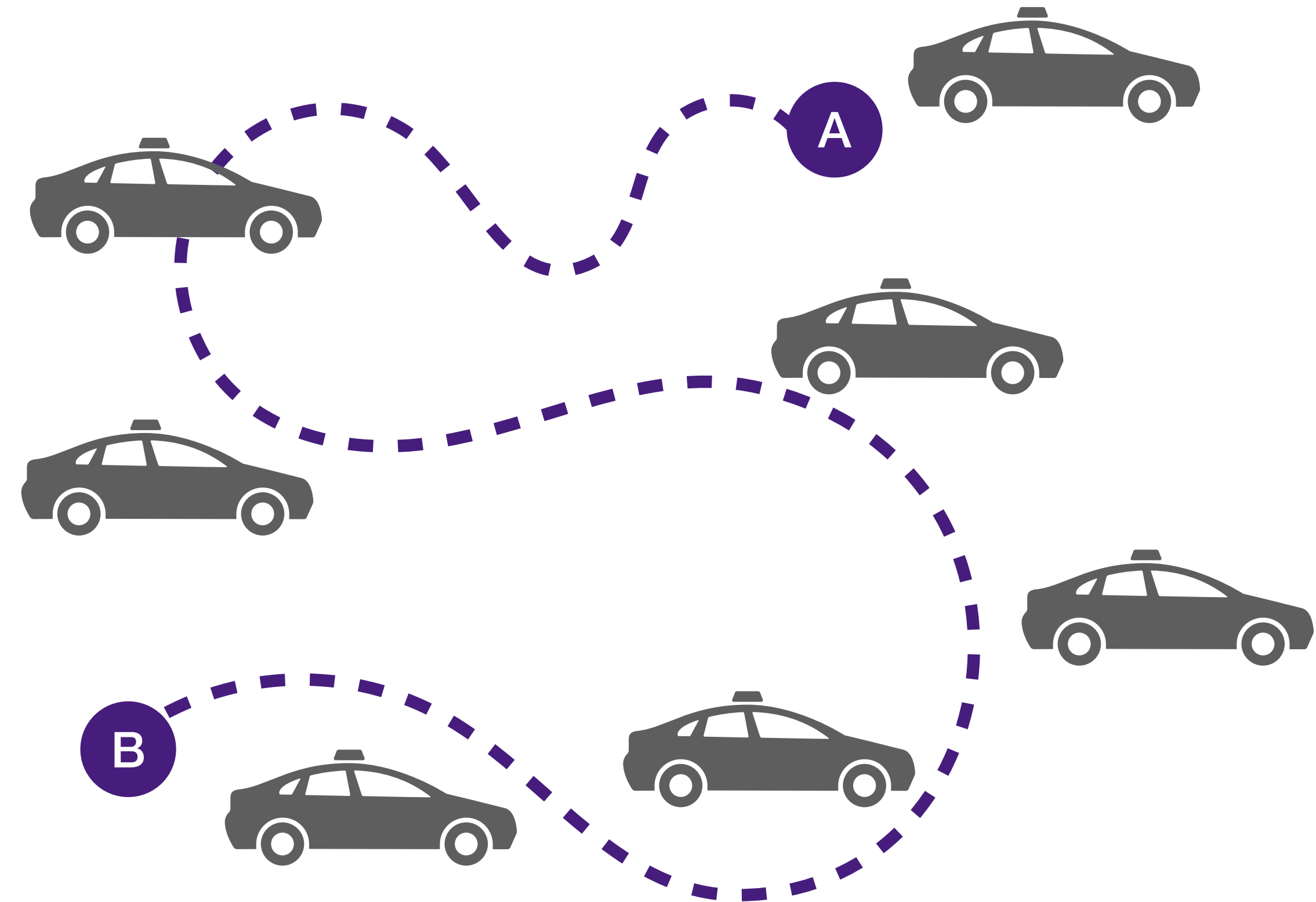
IntelliSys Lab

# Cloud / HPC

# Serverless



**Car rental**

**Uber/Lyft**

# *Idle* Resources

**"Which of the following factors contribute to avoidable cloud spend, also known as cloud waste, at your organization?"**

| Factor | Percentage |
|---|---|
| Idle or underused resources | 66% |
| Overprovisioning resources | 59% |
| Lack of needed skills | 47% |
| Manual containerization | 37% |
| We do not have any avoidable cloud spend | 6% |

FORRESTER

**Unlocking Multicloud's Operational Potential**

Centralize And Automate Your Cloud Platform To Increase Operational Efficiency And Security

A FORRESTER CONSULTING THOUGHT LEADERSHIP PAPER COMMISSIONED BY HASHICORP, AUGUST 2022

[1] Hashicorp-Forrester, "Unlocking Multicloud's Operational Potential," 2022. [Online].
    Available: https://www.datocms-assets.com/2885/1659554932-unlocking-multiclouds-operational-potential-forrester-hashicorp.pdf

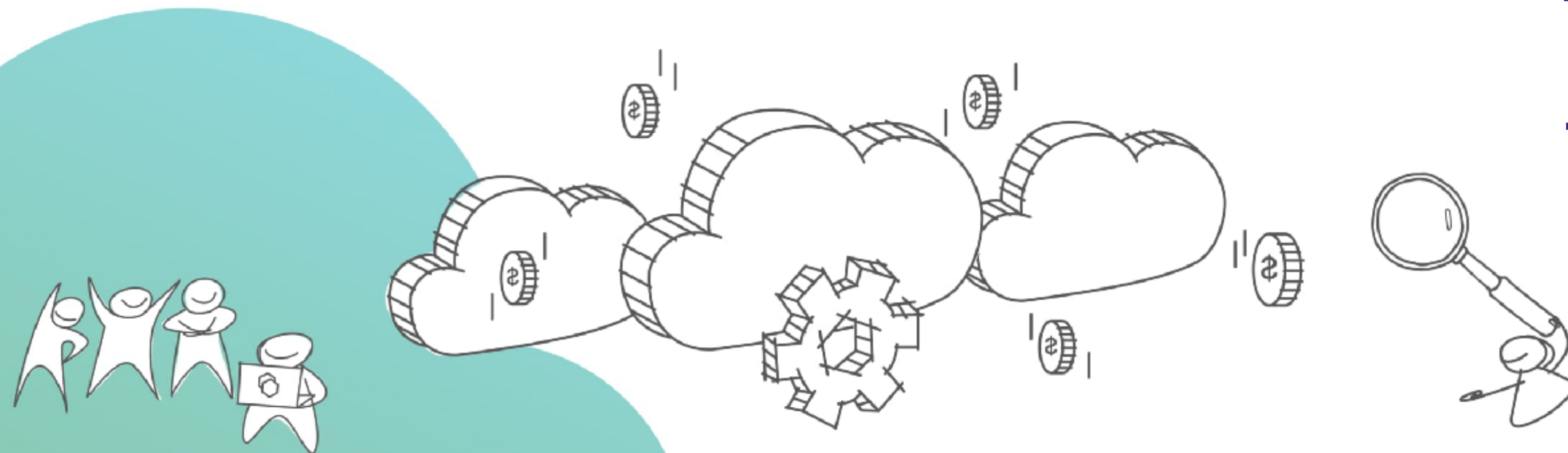**User over-provisioning**

**Computation patterns**

**Varying input data**

…

**Idle Resources**

**Low-priority VMs**
SmartHarvest [EuroSys 21]

**Memory harvesting VMs**
MHVM [ASPLOS 22]

**Harvest gSSDs**
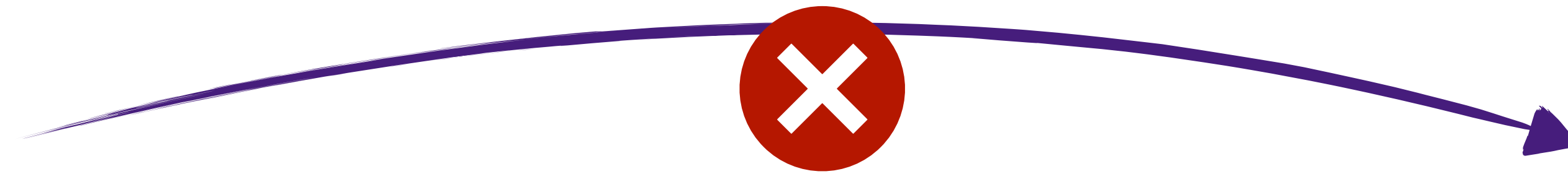BlockFlex [OSDI 22]

**Function acceleration**
Freyr [WWW 22]

…

**VM-based Cloud**

CherryPick [NSDI'17]

SmartHarvest [EuroSys'21]

Ernest [NSDI'16]

MHVM [ASPLOS'22]

**User-side** ← → **Provider-side**

OFC [EuroSys'21]

COSE [INFOCOM'20]

Freyr [WWW'22]

StepConf [INFOCOM'22]

⭐ Libra [HPDC'23]

**Serverless**

**Harvesting idle resources timely and safely?**

**Harvesting for VM-based cloud / HPC**

Long-running

Always-on

Workload-agnostic

**Serverless computing**
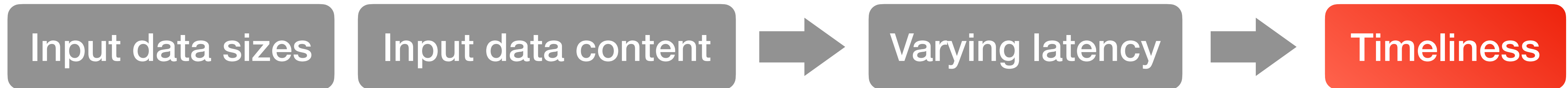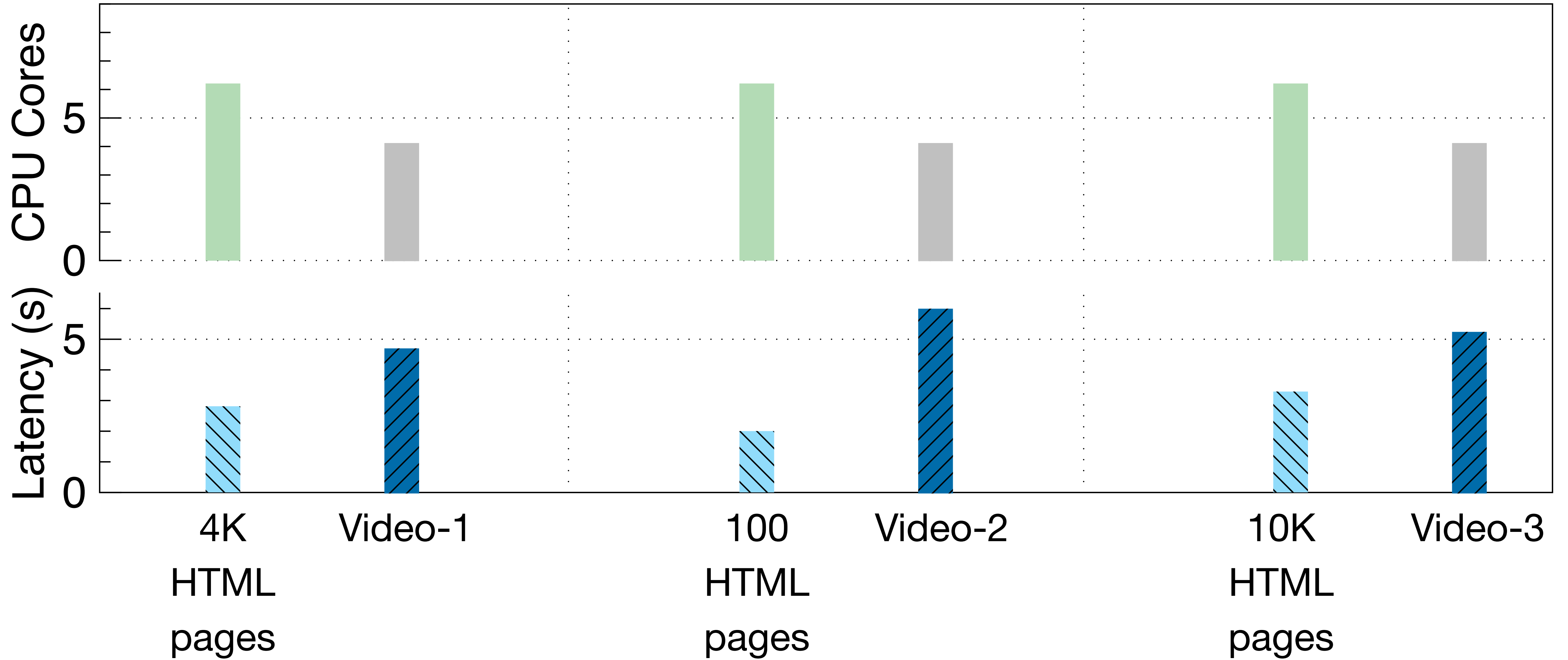
Short-living

Timeliness

Varying input sizes

# Limitations of Existing Work

- Provider-side works (OFC, Freyr)

  - Hard to generalize to *varying input data*

  - Ignorance of *resource timeliness*

- User-side works (COSE, StepConf)

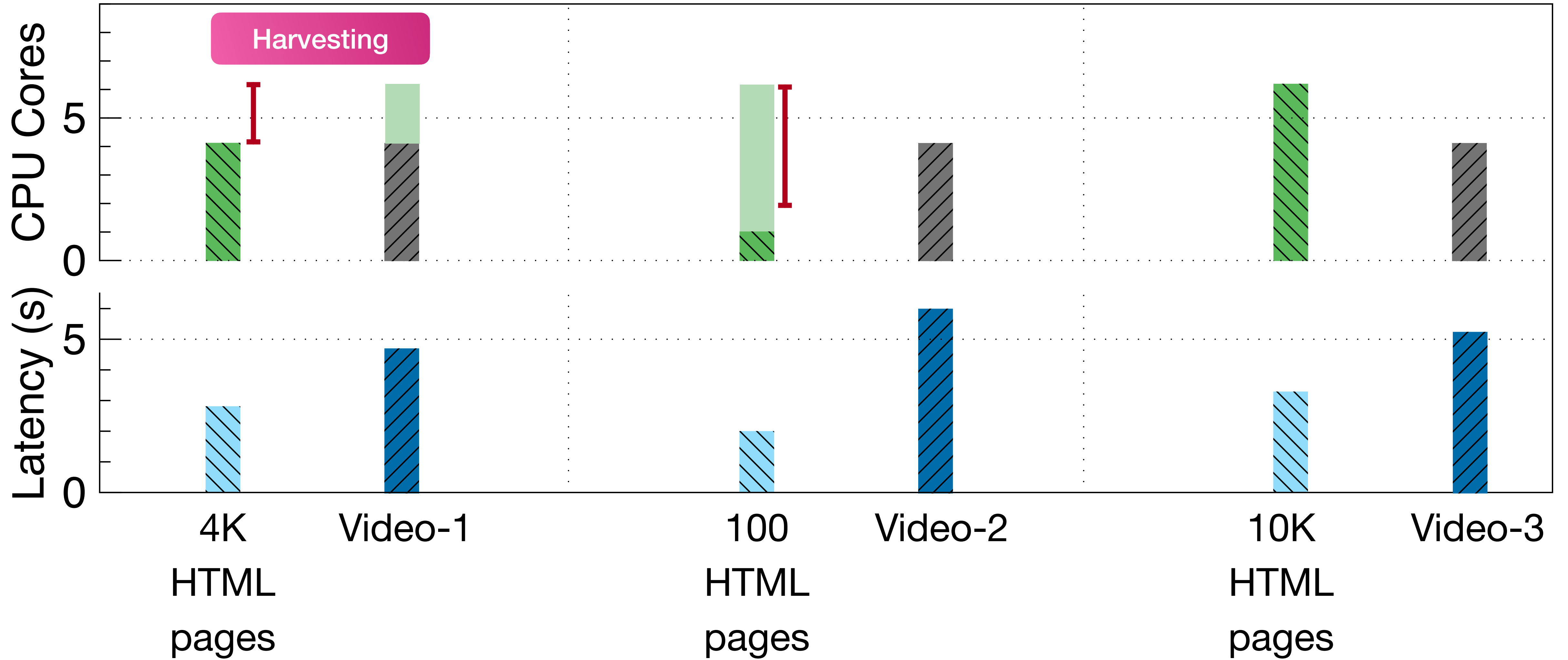  - Static configuration cannot satisfy dynamic resource demands

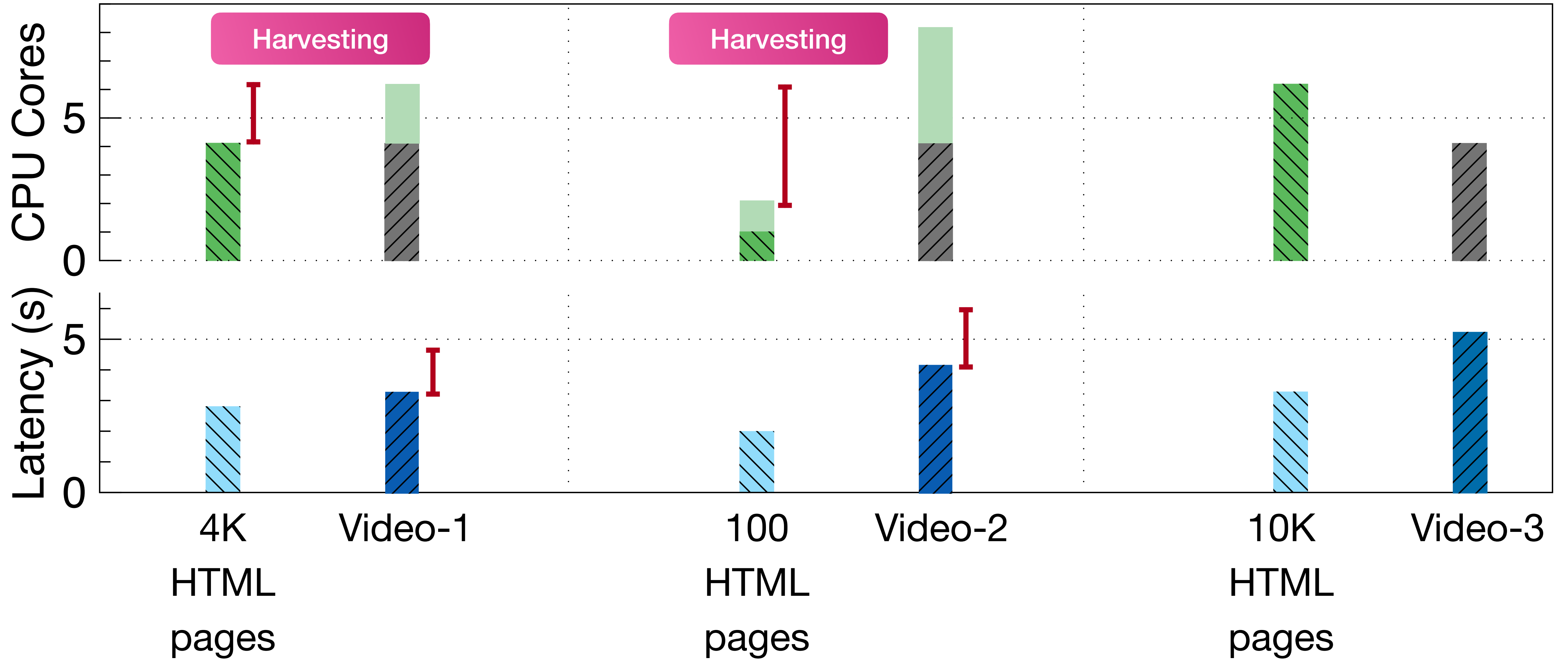DH: Dynamic HTML          VP: Video Processing

Input data sizes    Input data content  ➡  Varying latency  ➡  Timeliness

DH: Dynamic HTML          VP: Video Processing
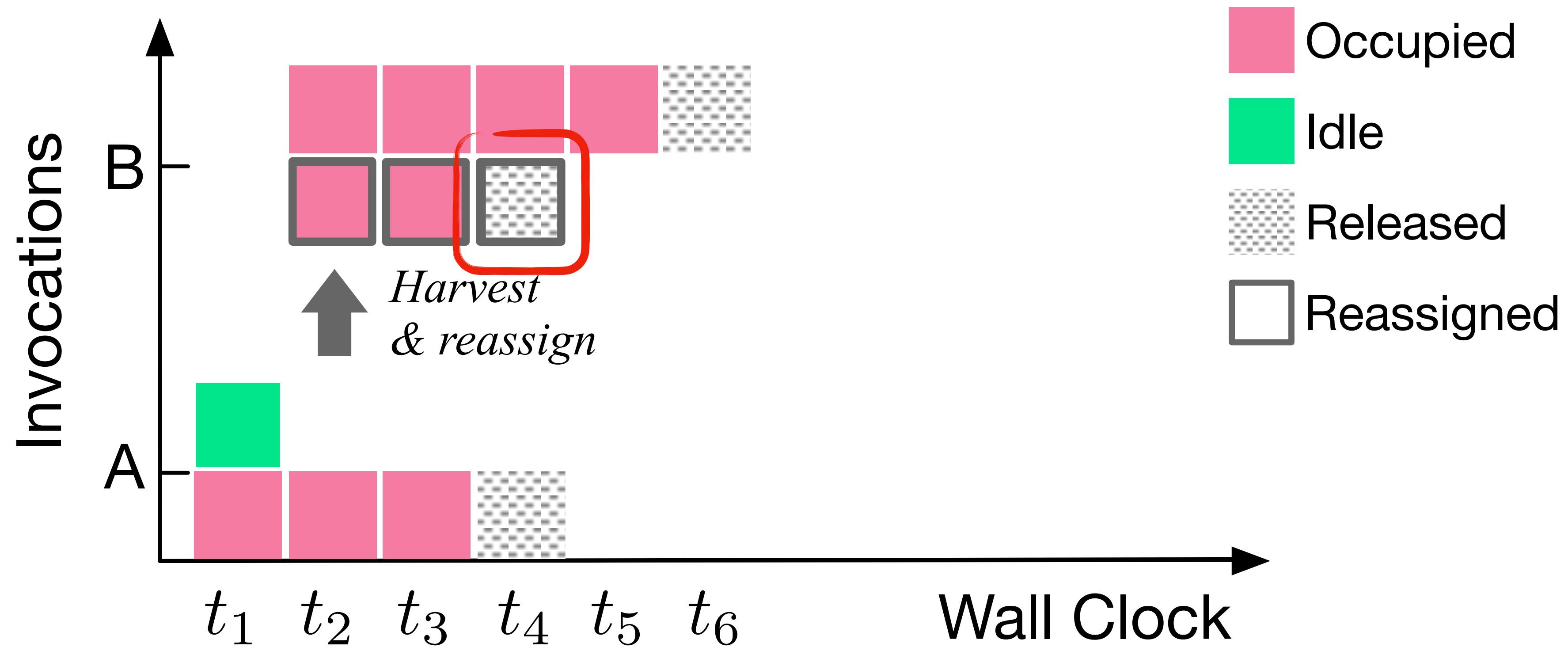
DH: Dynamic HTML          VP: Video Processing
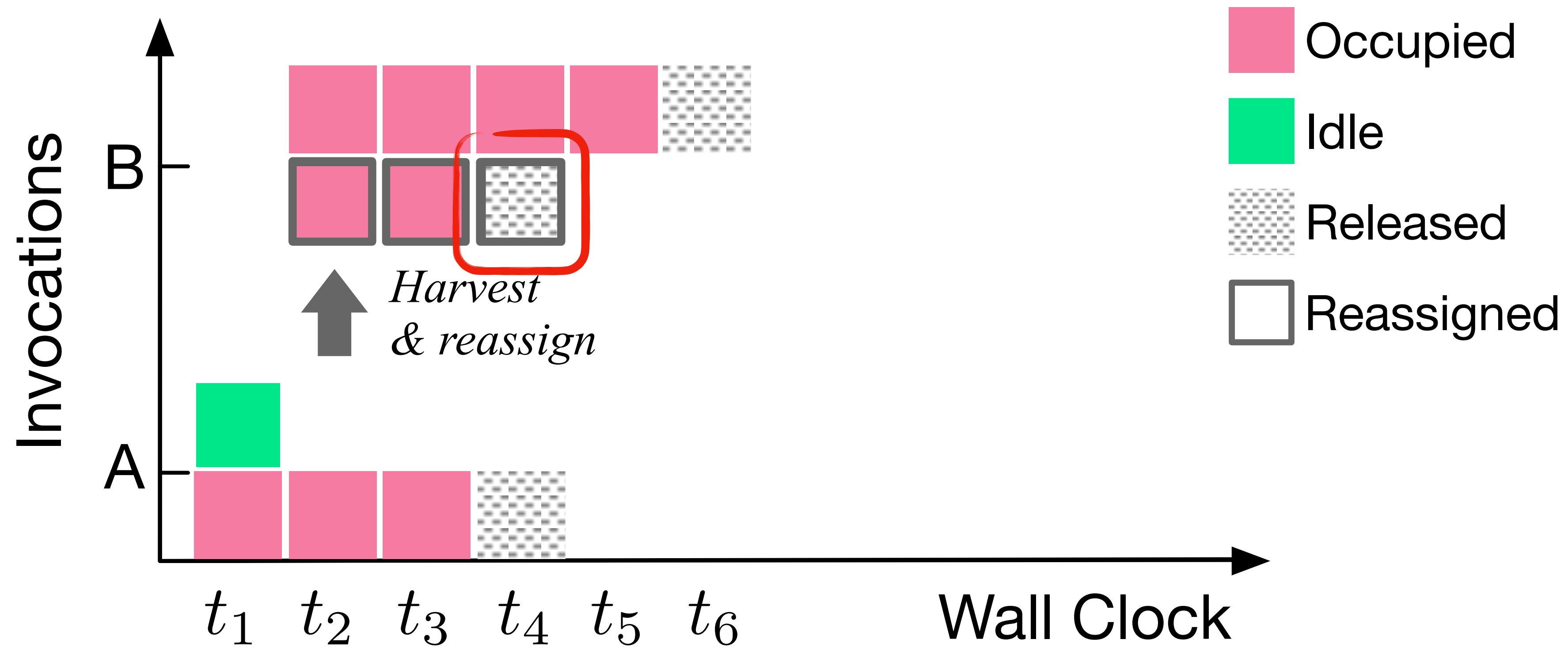
**Harvesting idle resources to accelerate functions**

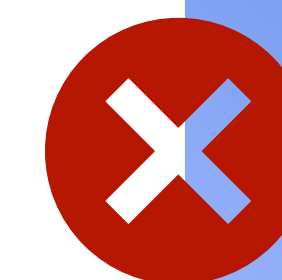# Timeliness of Resources



Invocation A: Over-provisioned, t1—t3

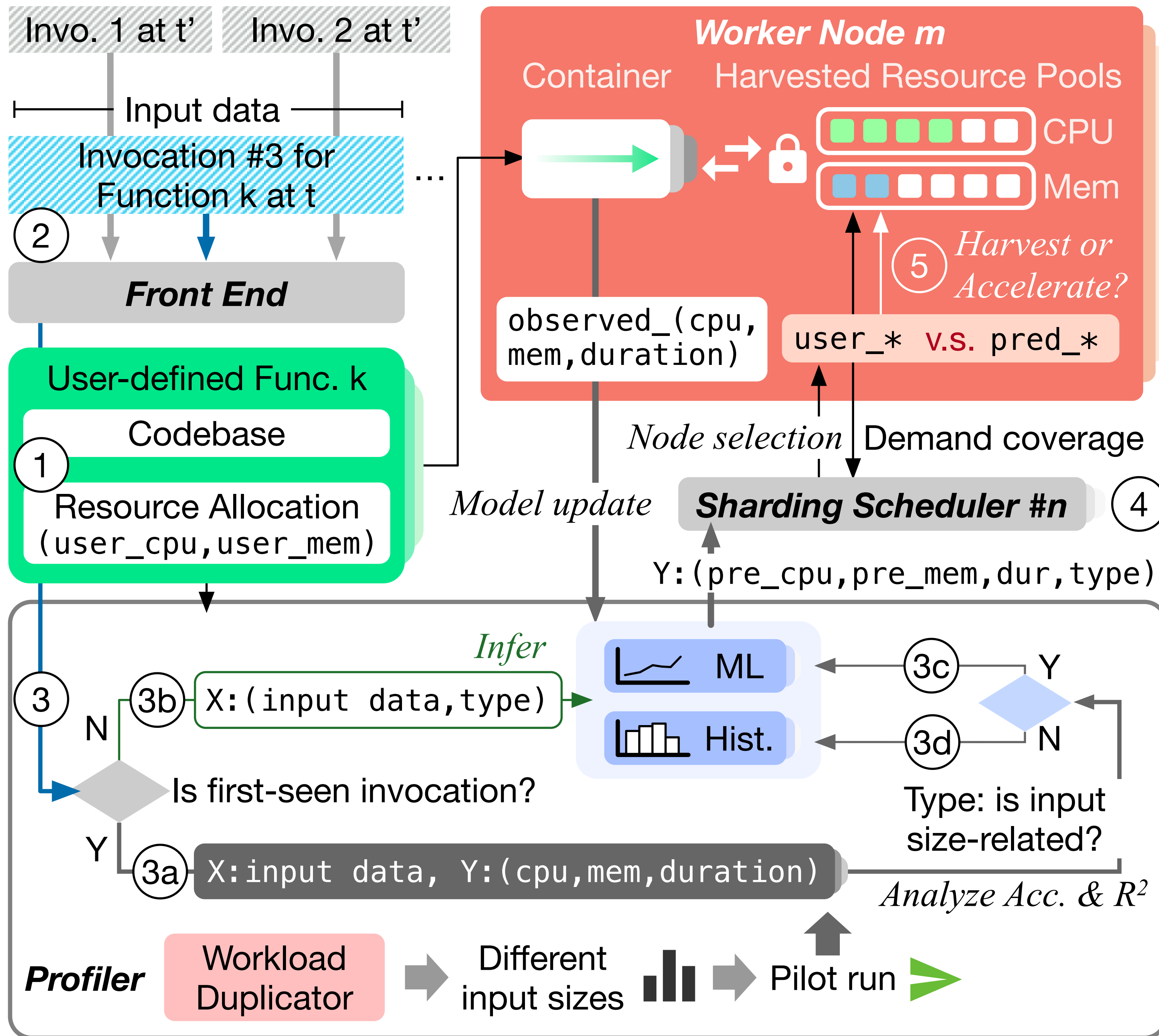Invocation B: Under-provisioned, t2—t5

# Timeliness of Resources



Invocation A: Over-provisioned, t1—t3

Invocation B: Under-provisioned, t2—t5

Unreserved Resources

# Libra Overview

- Profiling
- Scheduling
- Harvesting & Acceleration
- Safeguarding

## User-defined Func. k

Codebase

(1)

Resource Allocation
`(user_cpu,user_mem)`

# Step 1: Function Deployment

- Deploy code
- User-defined resource allocations

Invo. 1 at t'    Invo. 2 at t'

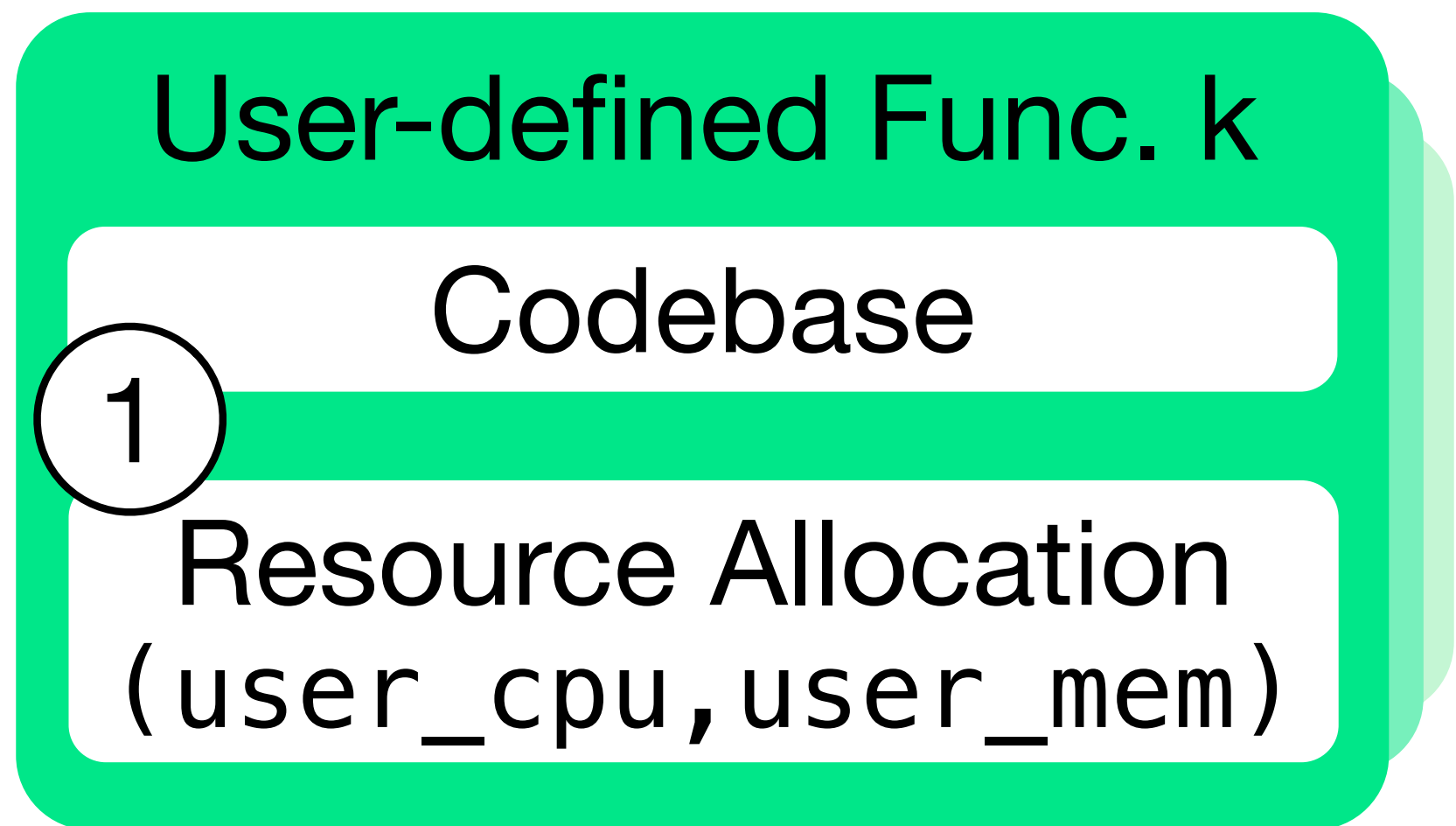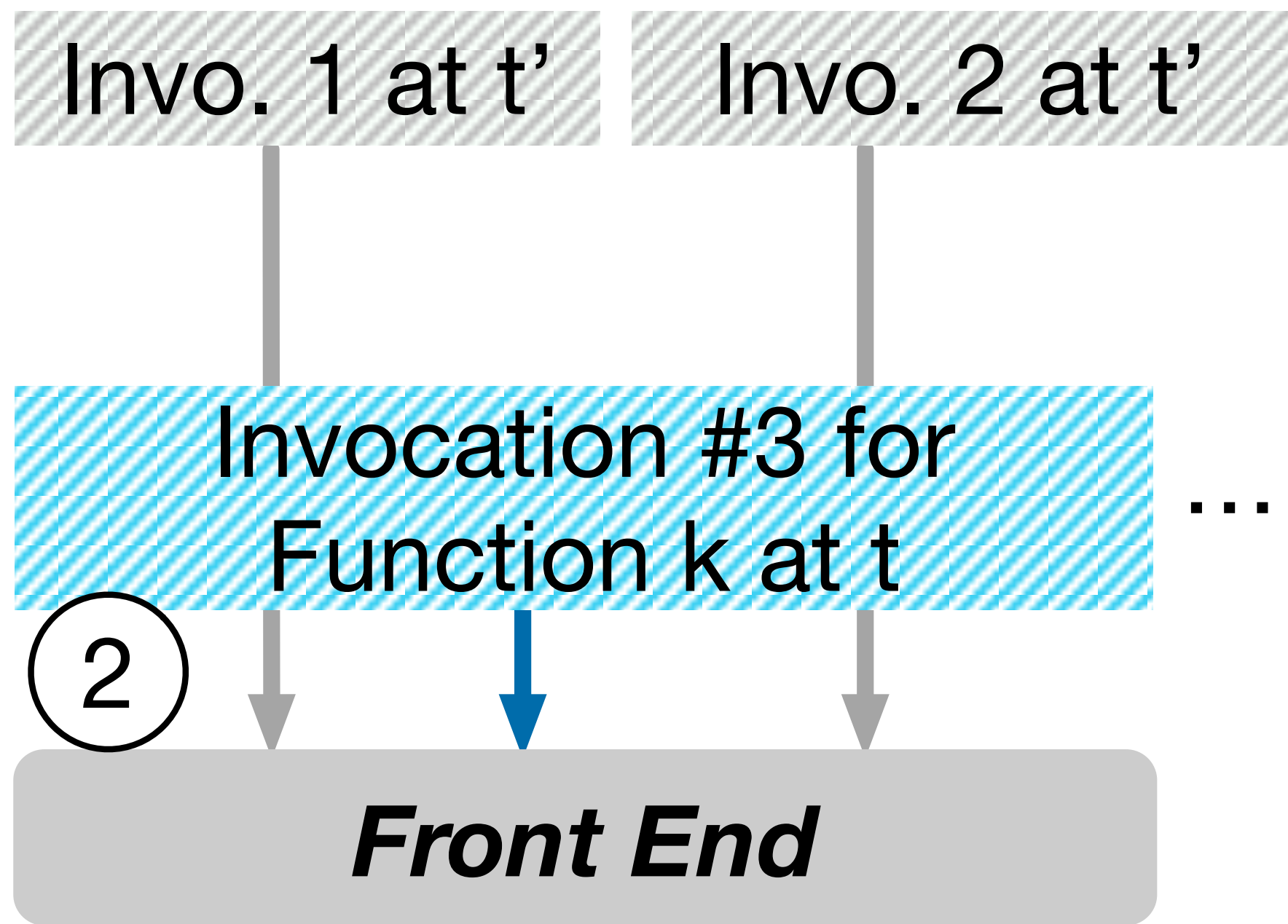Invocation #3 for
Function k at t    ...

②

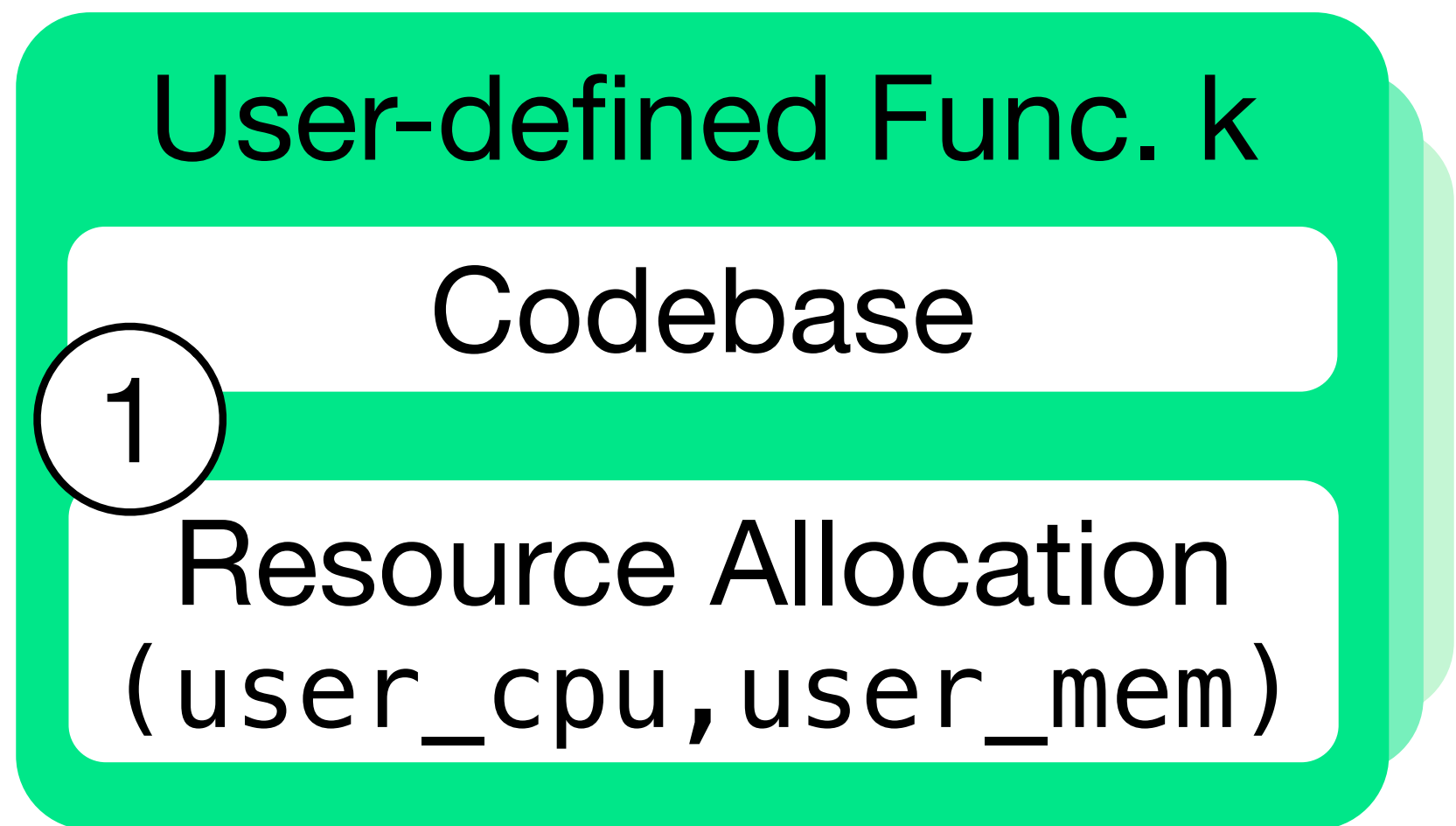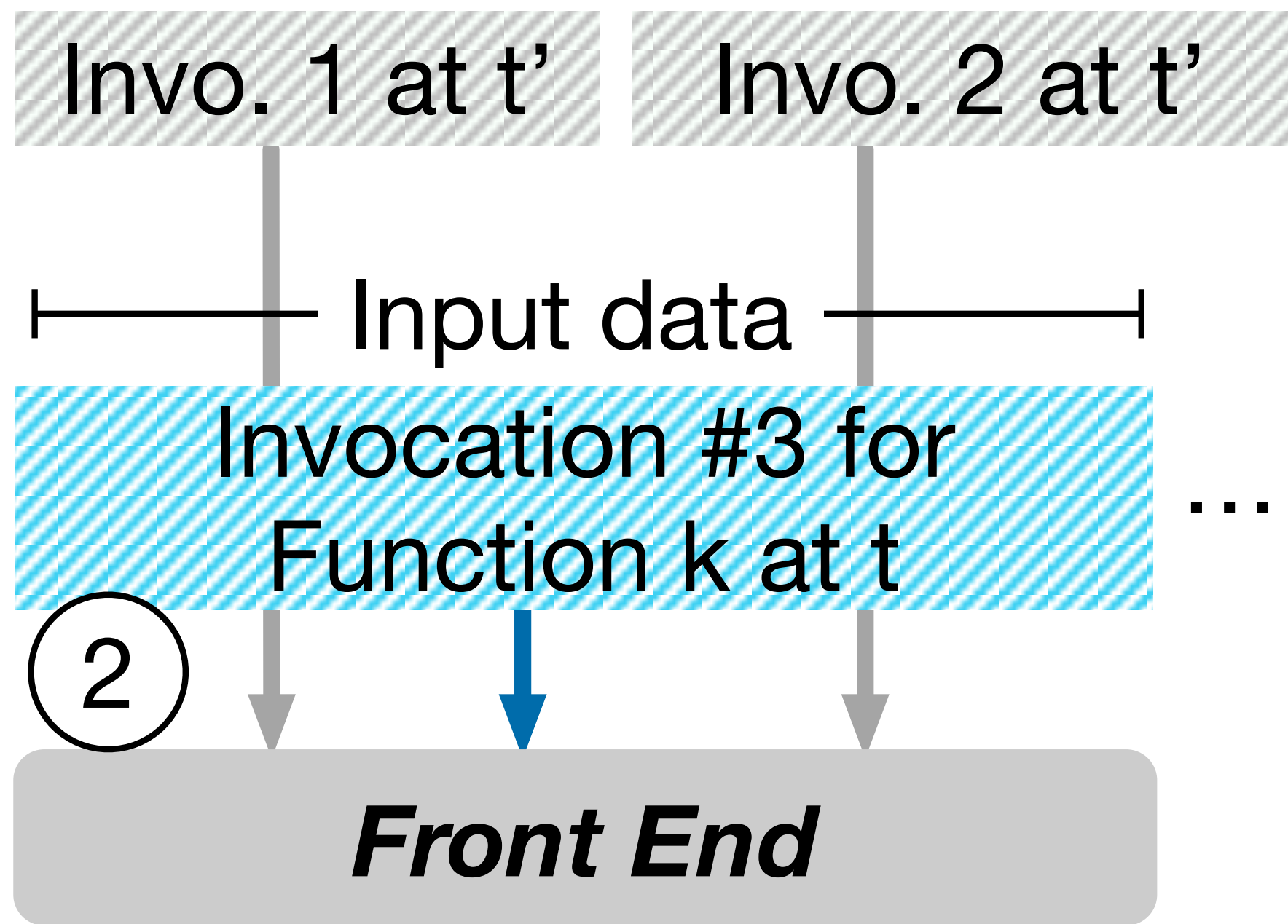# Step 2: Function Invocation

*Front End*

User-defined Func. k

Codebase

①

Resource Allocation
(`user_cpu,user_mem`)

# Step 2: Function Invocation

- Function requests arrive
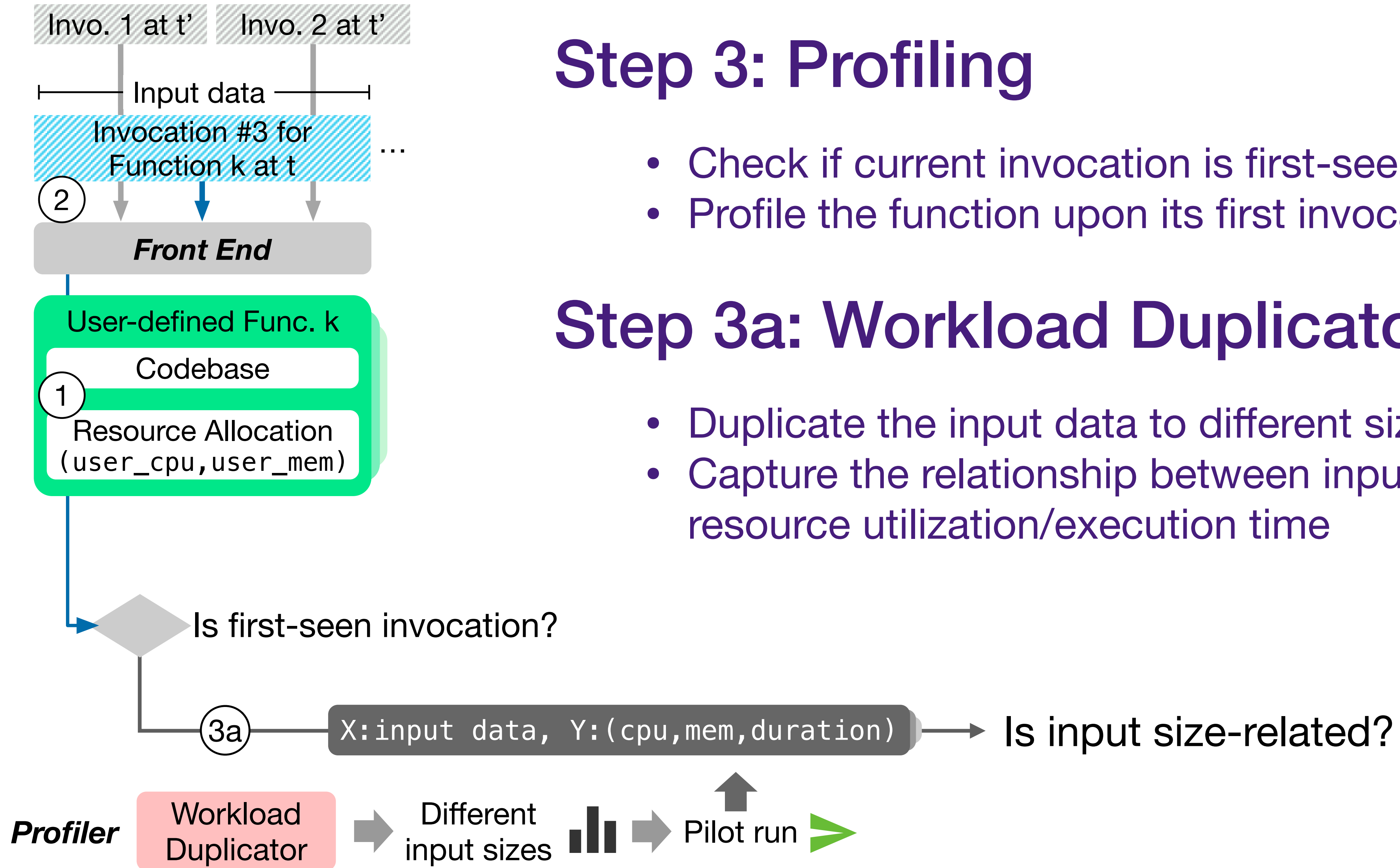
# Step 2: Function Invocation

- Function requests arrive
- Input data of varying sizes

# Step 3: Profiling

- Check if current invocation is first-seen
- Profile the function upon its first invocation

# Step 3a: Workload Duplicator

- Duplicate the input data to different sizes
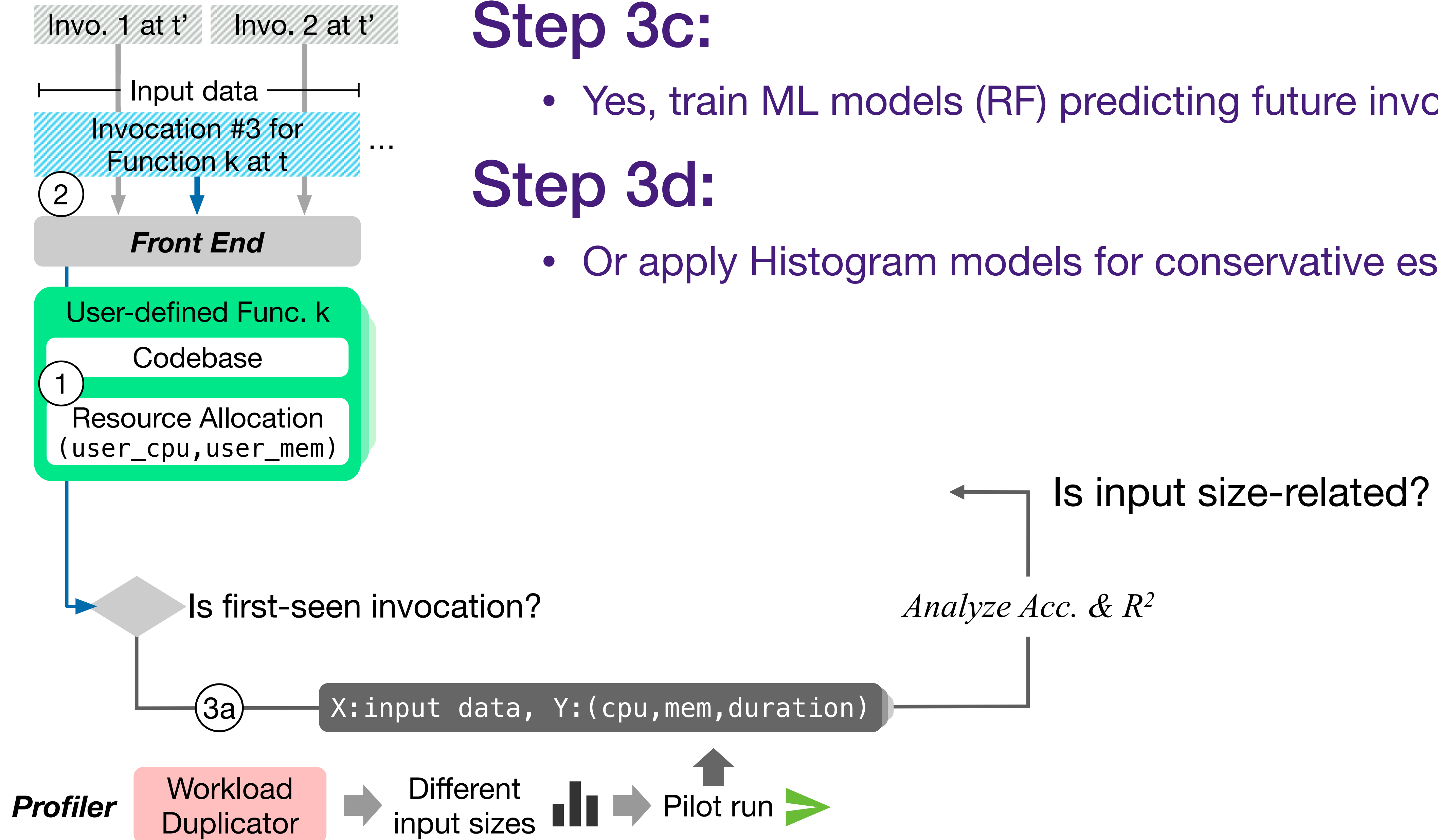- Capture the relationship between input size and resource utilization/execution time
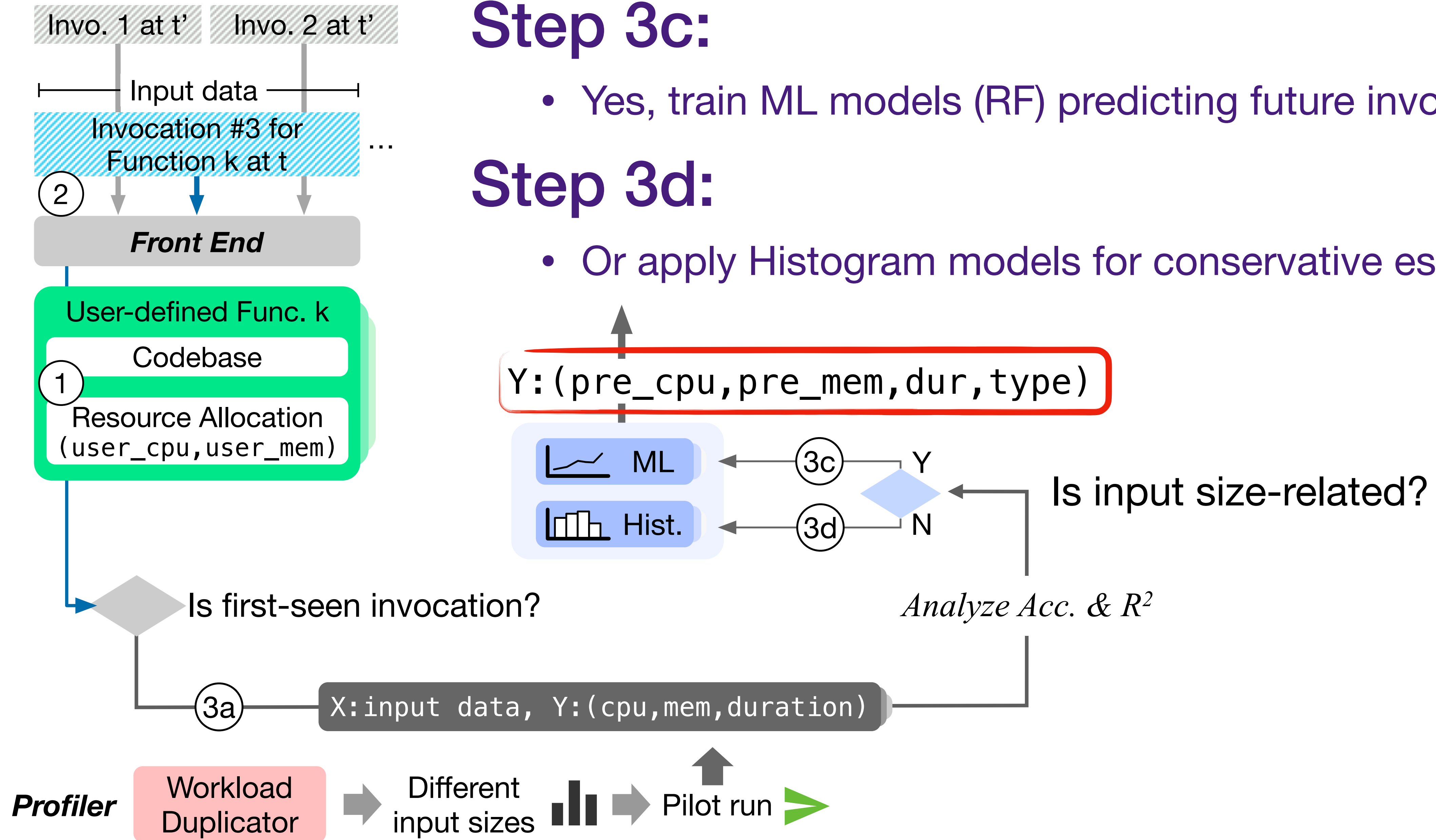
# Step 3c:

- Yes, train ML models (RF) predicting future invocations.

# Step 3d:

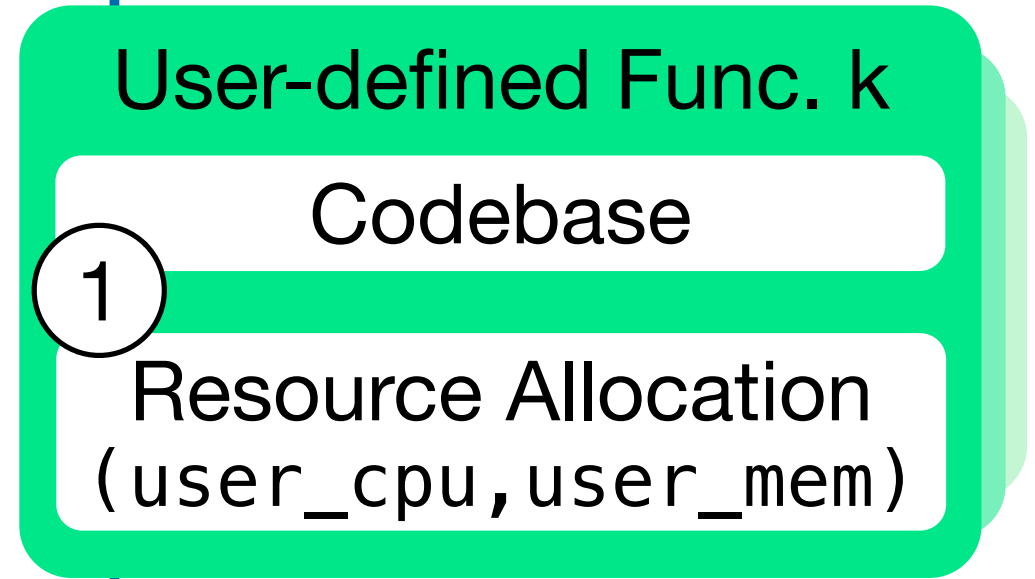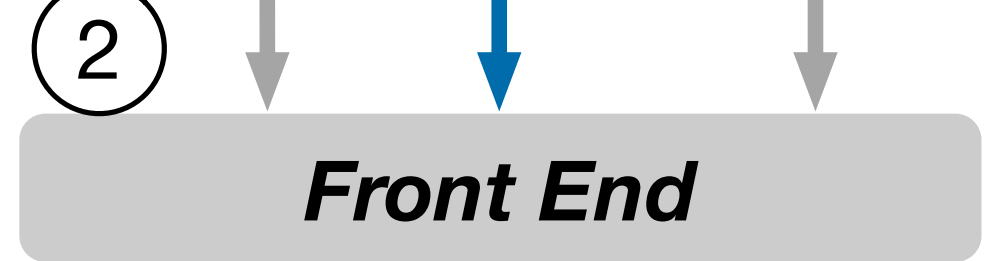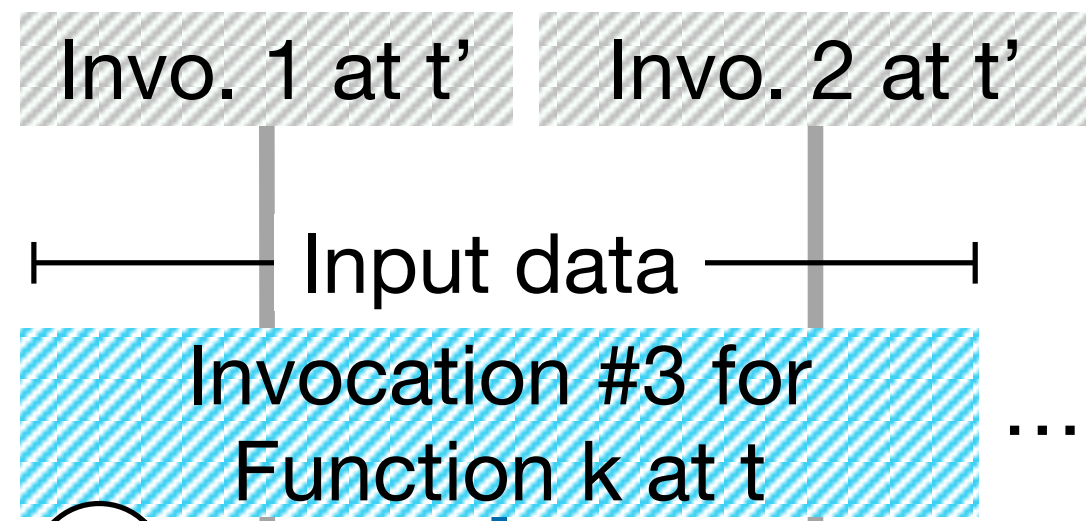- Or apply Histogram models for conservative estimations.

# Step 3c:

- Yes, train ML models (RF) predicting future invocations.

# Step 3d:

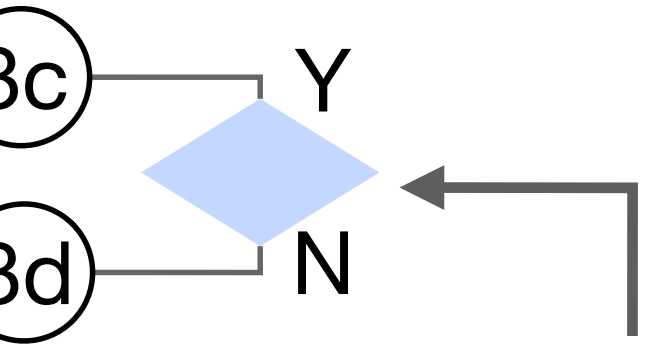- Or apply Histogram models for conservative estimations.

# Step 3b:

- If seen, use built models for estimation.

Invo. 1 at t'  Invo. 2 at t'

Input data

Invocation #3 for Function k at t  ...

② **Front End**

User-defined Func. k

Codebase

① Resource Allocation (user_cpu,user_mem)

*Infer*

3b X:(input data,type)
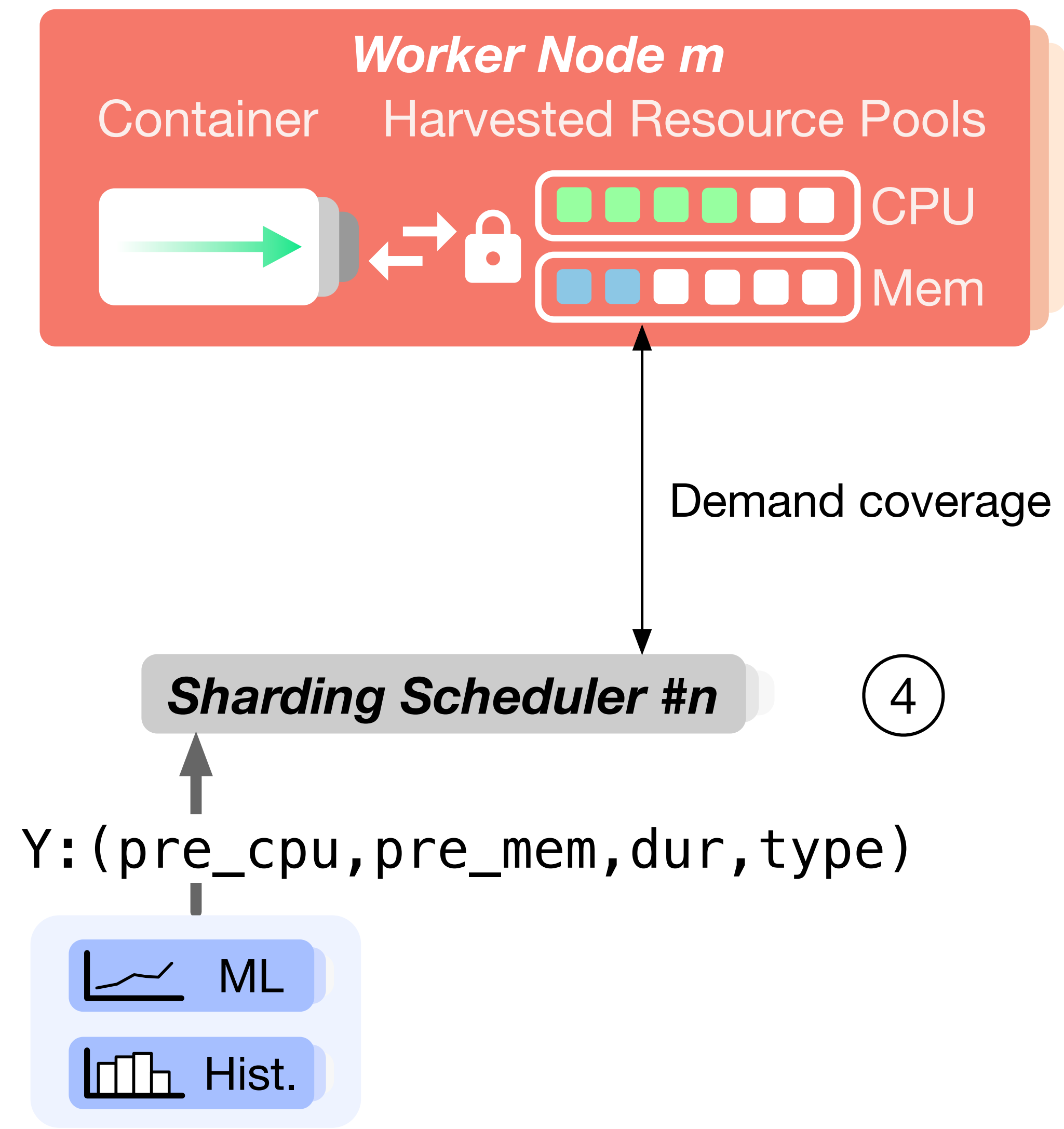
Is first-seen invocation?

3a X:input data, Y:(cpu,mem,duration)

Y:(pre_cpu,pre_mem,dur,type)

ML ⟵ 3c ⟵ Y

Hist. ⟵ 3d ⟵ N

Is input size-related?

*Analyze Acc. & R²*

*Profiler* Workload Duplicator ⟹ Different input sizes ⟹ Pilot run ▶

**Worker Node m**

Container     Harvested Resource Pools

CPU

Mem

Demand coverage

*Sharding Scheduler #n* ④
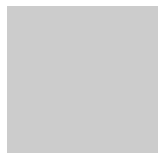
`Y:(pre_cpu,pre_mem,dur,type)`

ML

Hist.

# Step 4: Scheduling

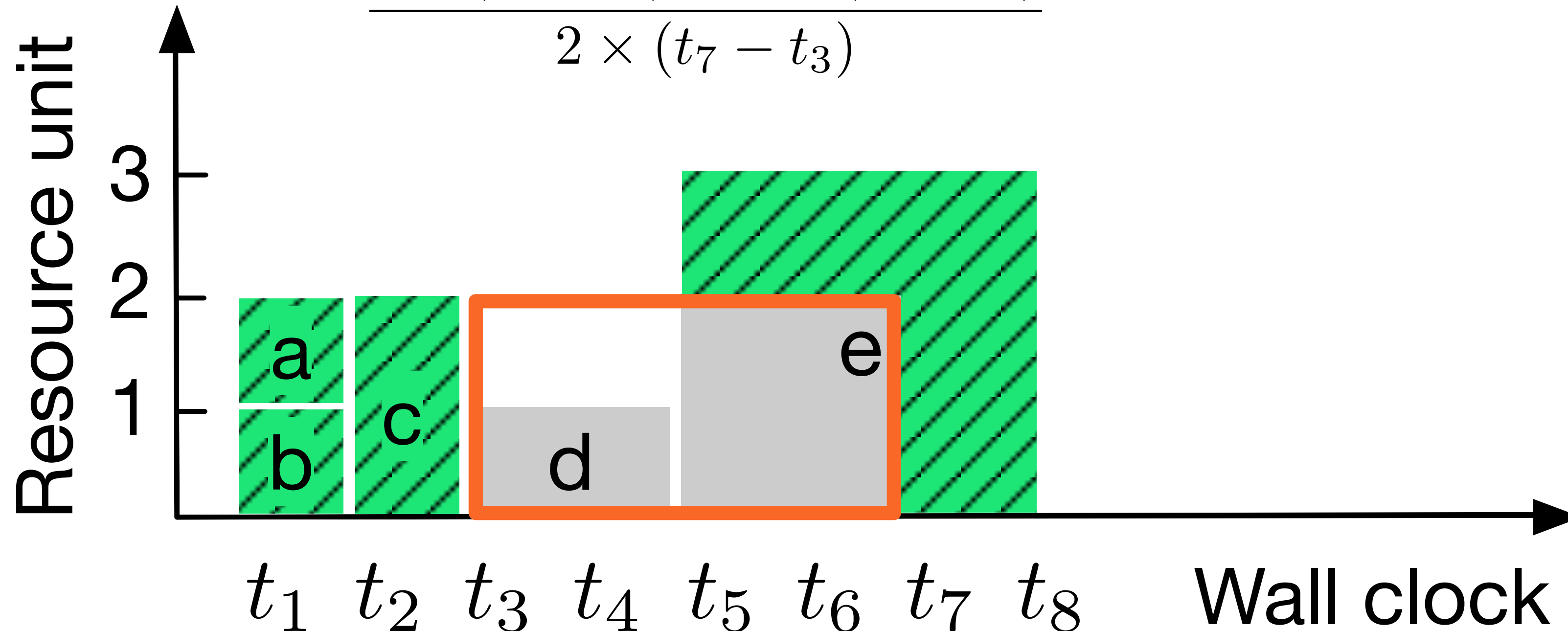- To harvest
- Or to accelerate?
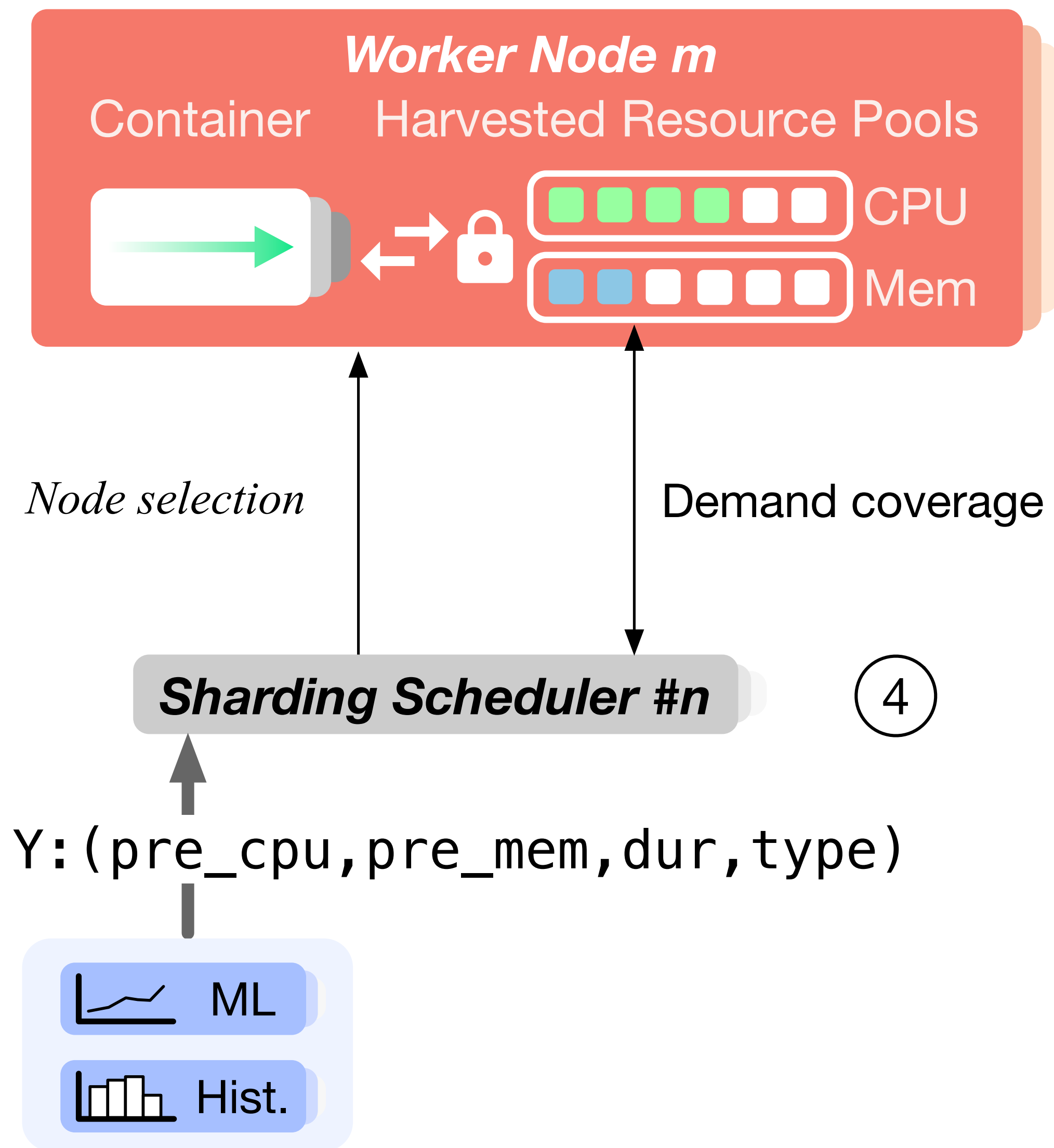
  - Calculate demand coverage

# Demand Coverage

An incoming invocation,
    which demands <u>two extra resource units</u>, $t3-t7$.

Idle resources in pool

Covered resources

Invocation demands

Demand coverage =

$$\frac{1 \times (t_5 - t_3) + 2 \times (t_7 - t_5)}{2 \times (t_7 - t_3)}$$



Resource unit

3

2

1

a

b

c

d

e

$t_1$  $t_2$  $t_3$  $t_4$  $t_5$  $t_6$  $t_7$  $t_8$   Wall clock
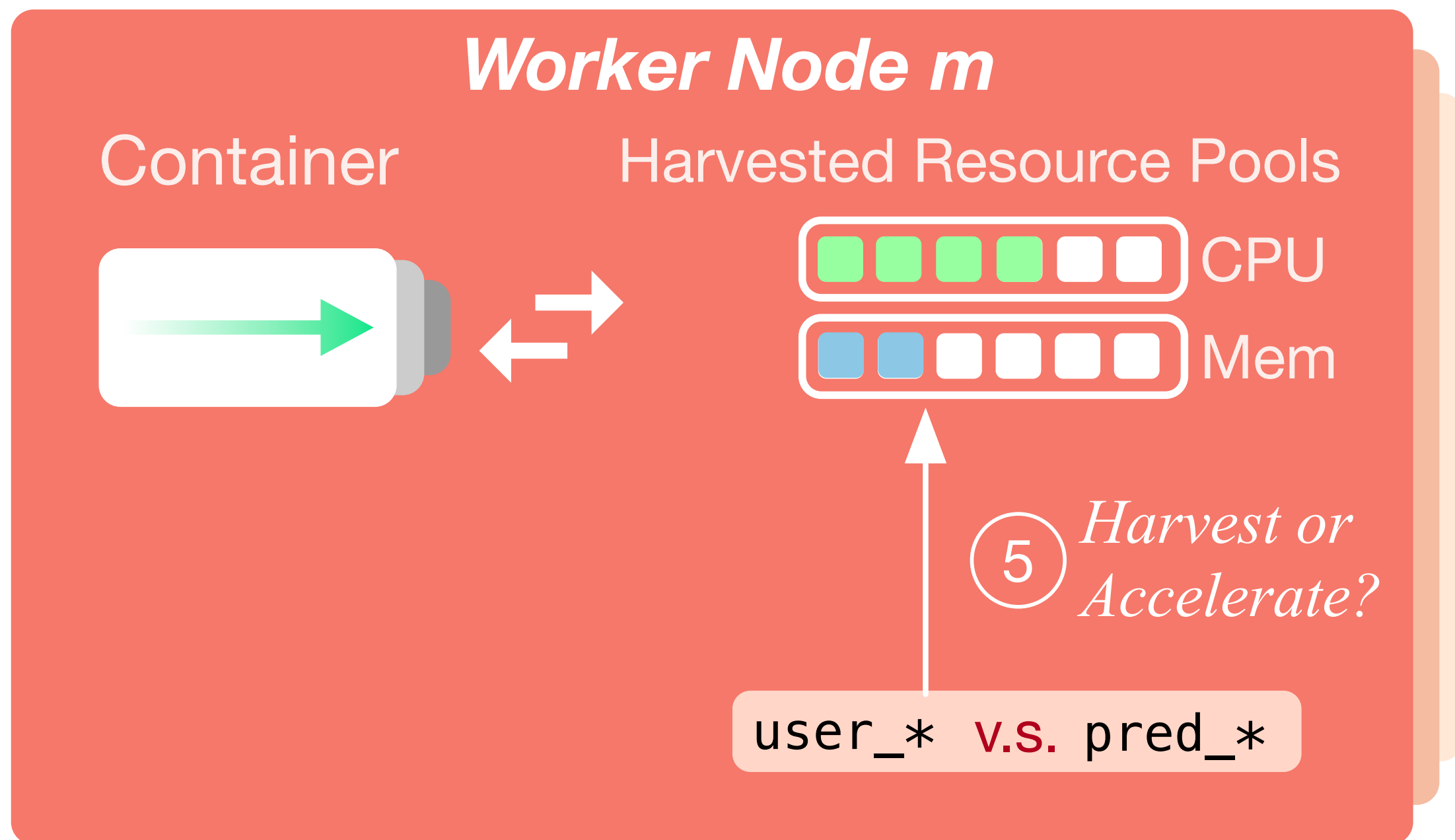
# Step 4: Scheduling

- Calculate demand coverage in realtime
- Select the node with the highest score
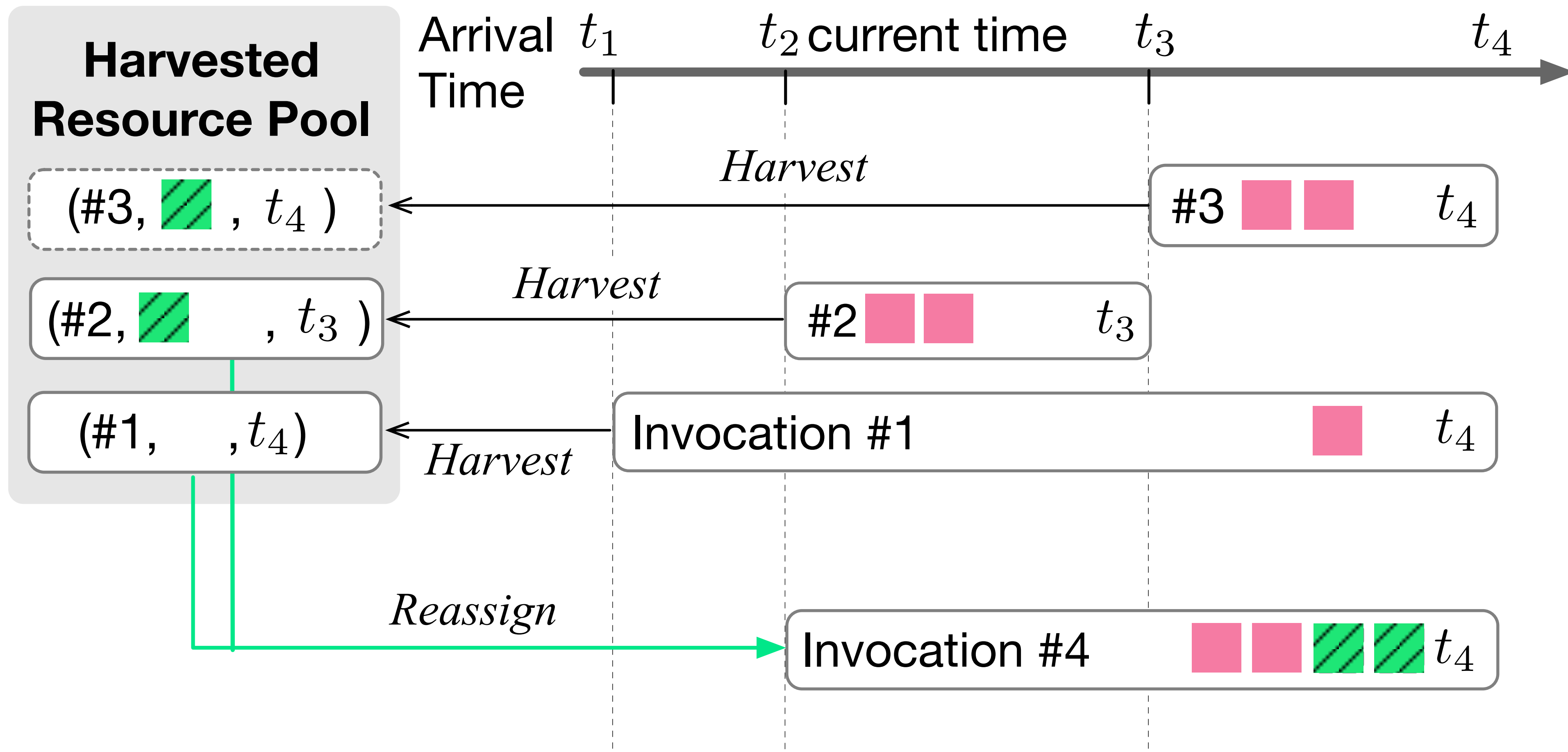- Same score, then consider locality

Coverage score =

$$\alpha \times D_{cpu} + (1-\alpha) \times D_{memory}$$
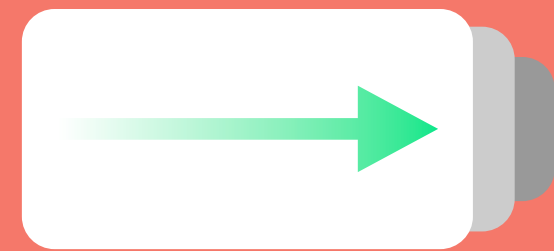
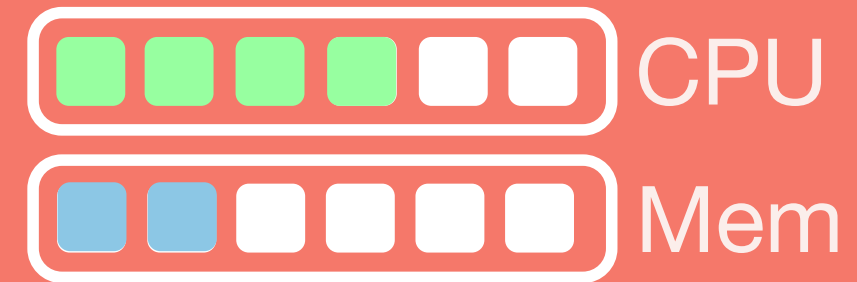Demand coverage of CPU

Demand coverage of Memory

**Worker Node m**

Container          Harvested Resource Pools

CPU

Mem

⑤ *Harvest or Accelerate?*

`user_*` `v.s.` `pred_*`

# Step 5: Harvesting & Acceleration

**Harvested Resource Pool**

(#3, 🟩 , $t_4$ )

(#2, 🟩 , $t_3$ )

(#1, , $t_4$)

Arrival Time

$t_1$    $t_2$ current time    $t_3$    $t_4$

*Harvest*

#3 🟥 🟥    $t_4$

*Harvest*

#2 🟥 🟥    $t_3$

Invocation #1    🟥    $t_4$

*Harvest*

*Reassign*

Invocation #4    🟥 🟥 🟩 🟩 $t_4$

🟩 Idle [ ] More resources needed 🟥 Busy $t$ Estimated completion time

**SafeGuard**

# Safeguard in Realtime

**Updating observations**

# Implementation

Libra is prototyped on top of Apache OpenWhisk using Scala

| **Profiler** | **Scheduler** | **Harvest Pool** | **Safeguard** |
|:---:|:---:|:---:|:---:|
| Python | OpenWhisk's | OpenWhisk's | Docker container |
| Scikit-learn | Load Balancer | Invoker | Linux `cgroups` |

# Evaluation

**Baselines
for scheduling**

OpenWhisk default

Min-Worker-Set [5]

Join-the-Shortest-Queue

Round-robin

**Baselines
for harvesting**

OpenWhisk default

Freyr [4]

**Metrics**

Function response latency

Resource utilization

**Benchmarks**

SeBS [6]

ServerlessBench [7]

ENSURE [8]

**Traces**

Azure Functions traces [9]

1K+ invocation traces

[4] Yu, Hanfei, et al. "Accelerating serverless computing by harvesting idle resources." WWW 22

[5] Zhang, Yanqi, et al. "Faster and cheaper serverless computing on harvested resources." SOSP 21

[6] Copik, Marcin, et al. "Sebs: A serverless benchmark suite for function-as-a-service computing." MIDDLEWARE 21

[7] Yu, Tianyi, et al. "Characterizing serverless platforms with serverlessbench." SoCC 20

[8] Suresh, Amoghavarsha, et al. "Ensure: Efficient scheduling…" ACSOS 20

[9] Shahrad, Mohammad, et al. "Serverless in the Wild…" ATC 20

# Testbed Clusters

**Evaluating harvesting**

3 nodes
72 Intel Xeon
E5-2670 CPU cores
72 GB memory

**Evaluating scheduling**

6 nodes
160 Intel Xeon
E5-2420 CPU cores
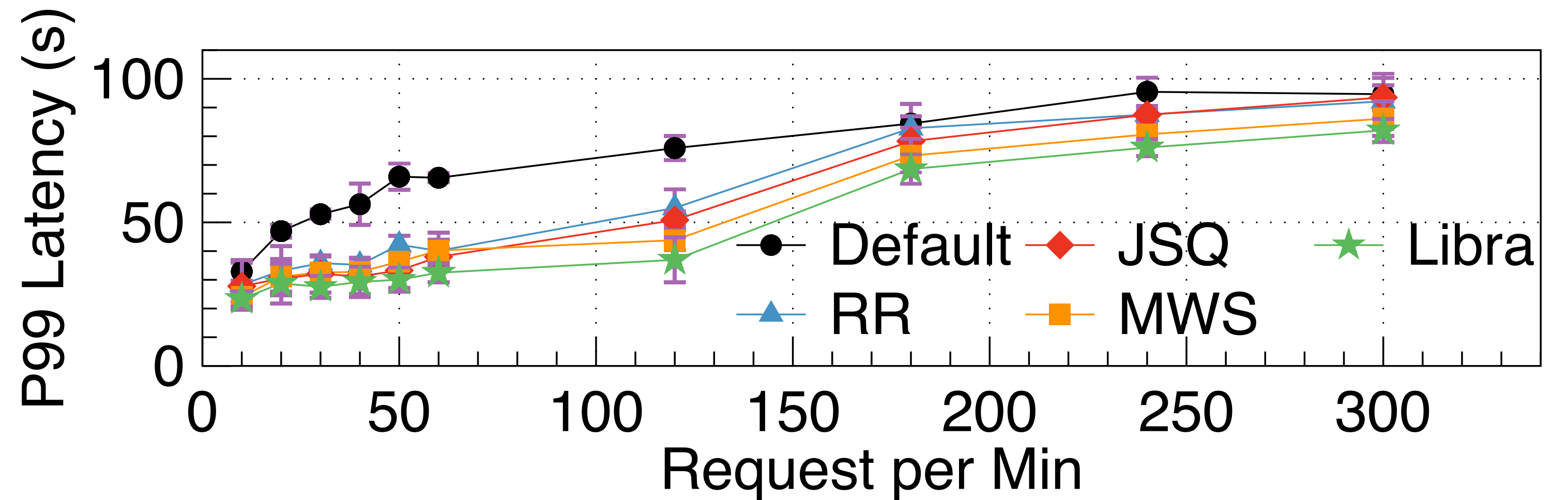160 GB memory

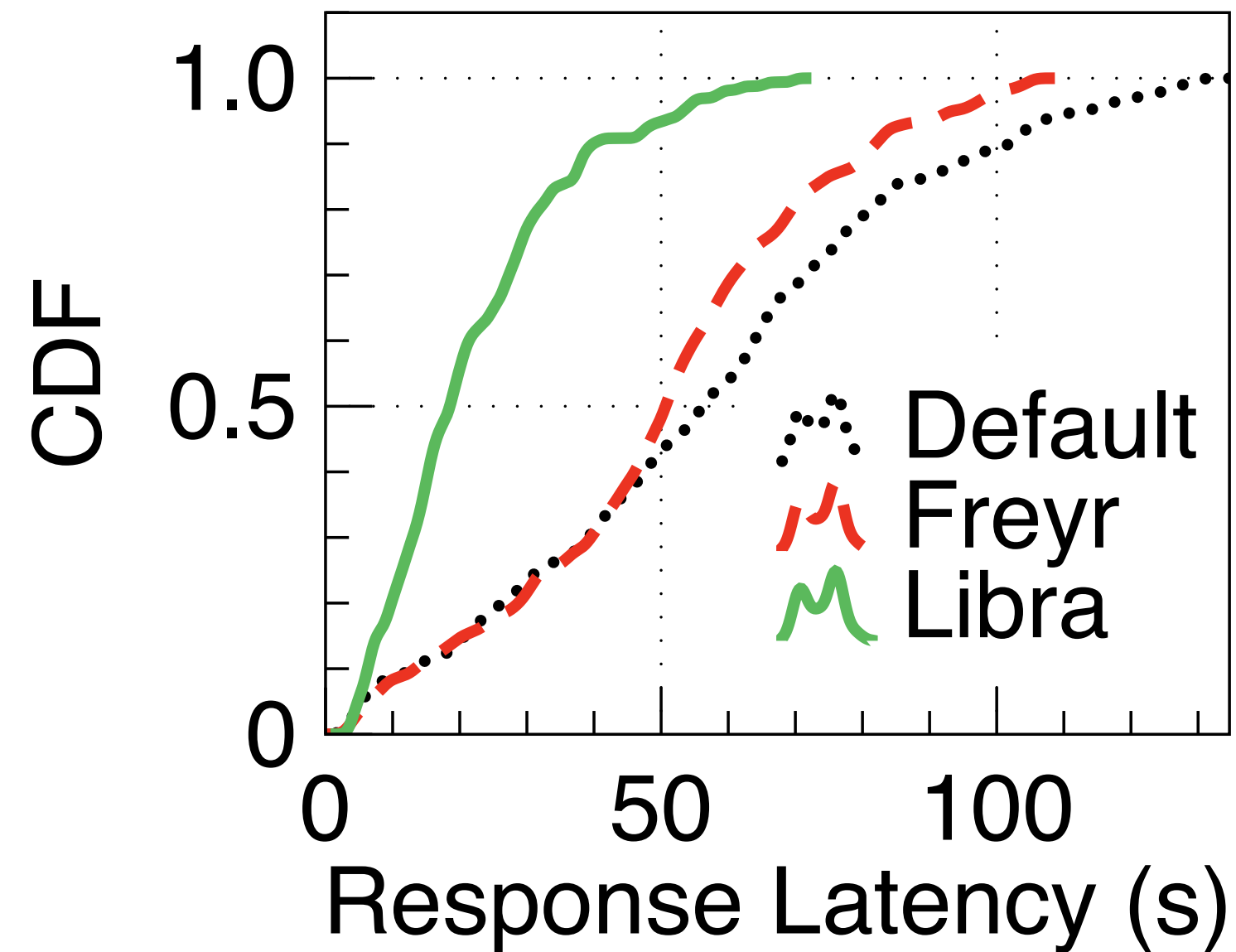**Evaluating scalability**

50 Jetstream nodes
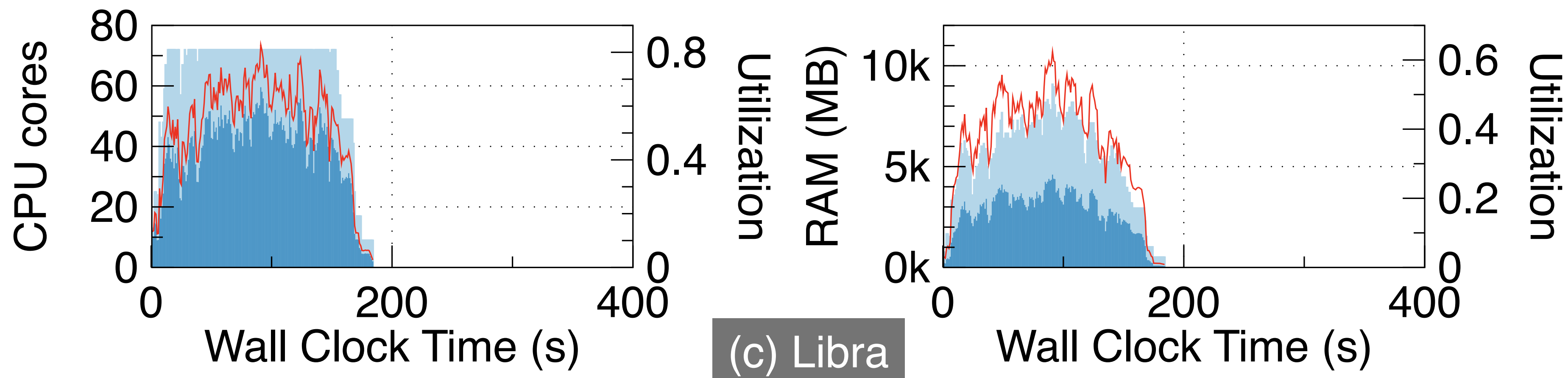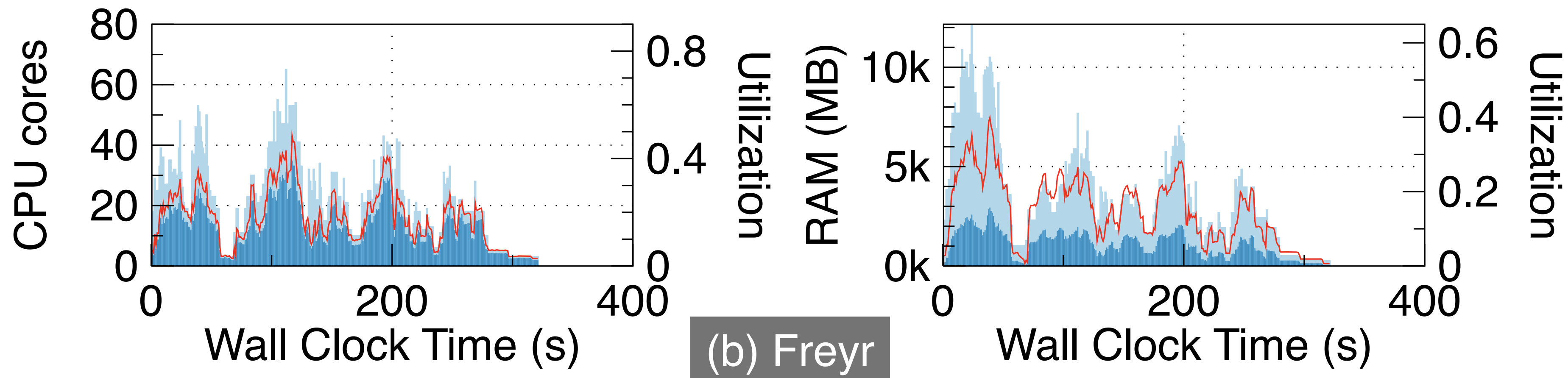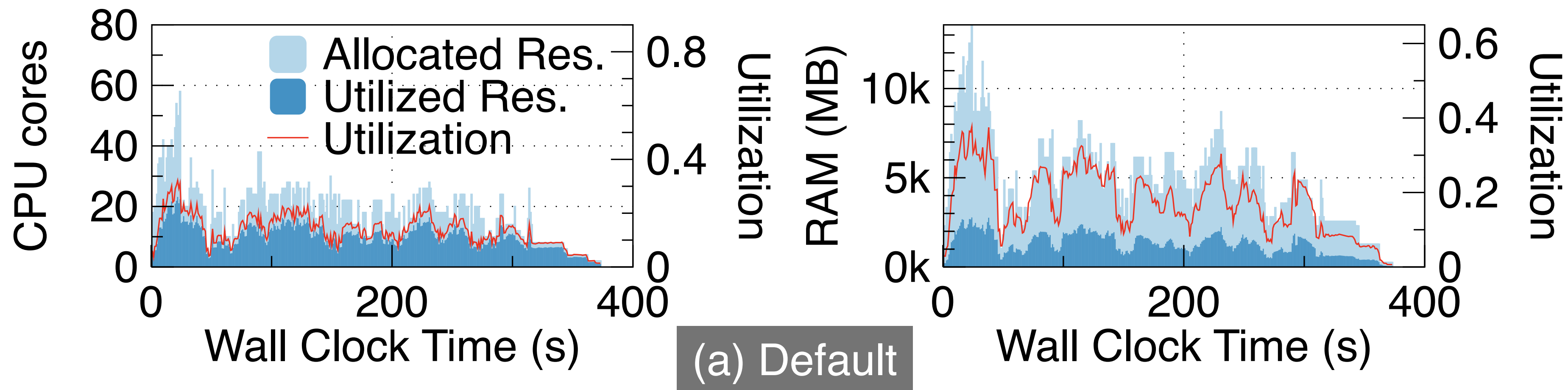1,200 Intel Xeon
E5-2680 CPU cores
1,200 GB memory

# Function Response Latency

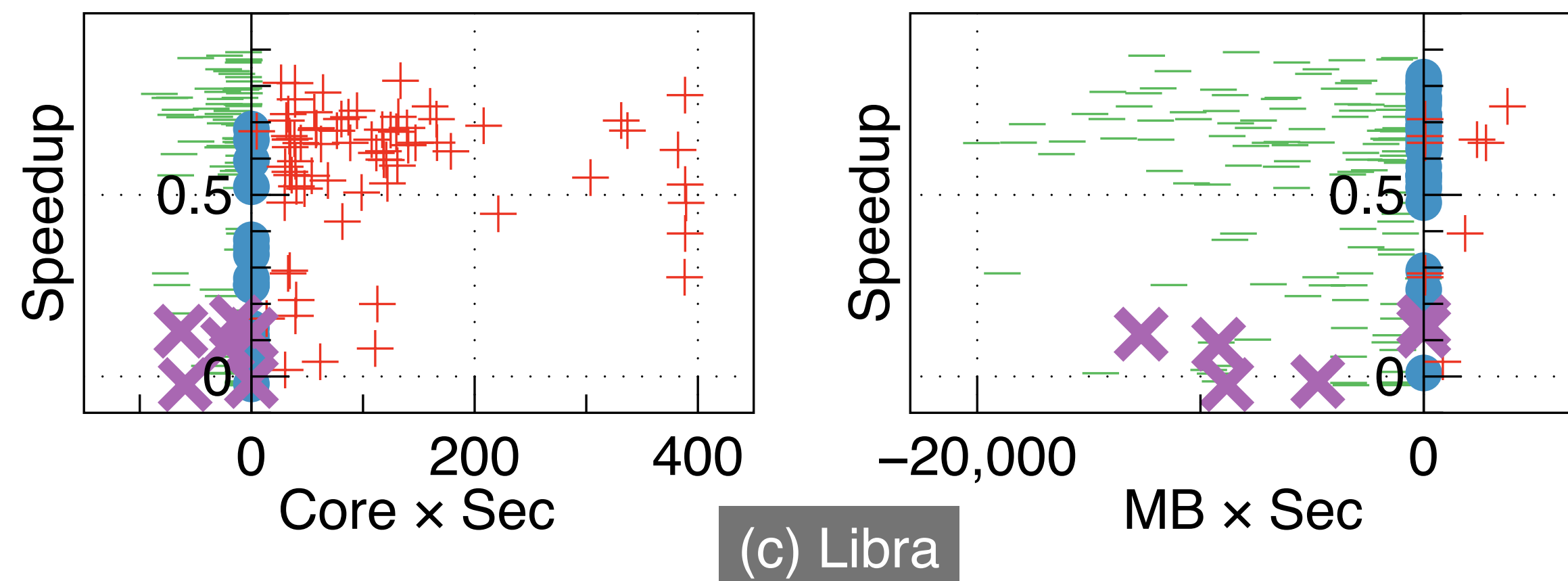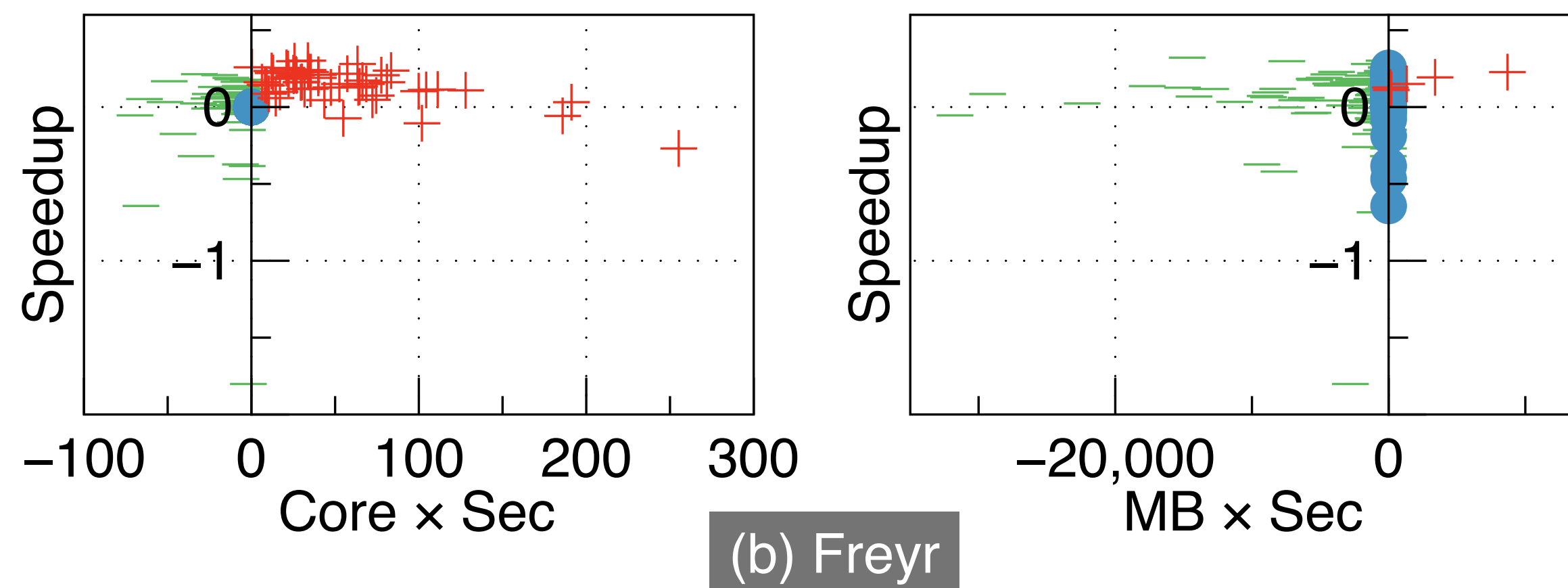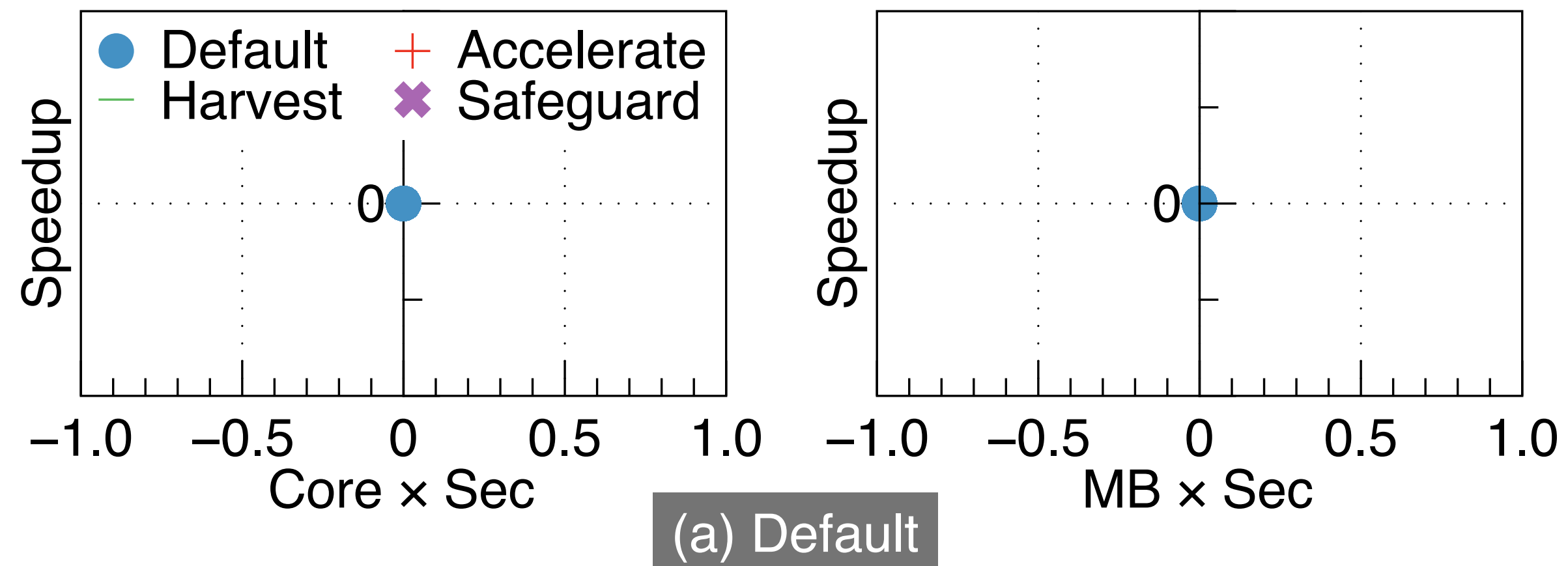**Libra provides lowest function response latency**



**Default**: OpenWhisk default     **JSQ**: Join-the-Shortest-Queue

**MWS**: Min-Worker-Set     **RR**: Round-robin

(a) Default

(b) Freyr

(c) Libra

**Resource Utilization**

32

(a) Default

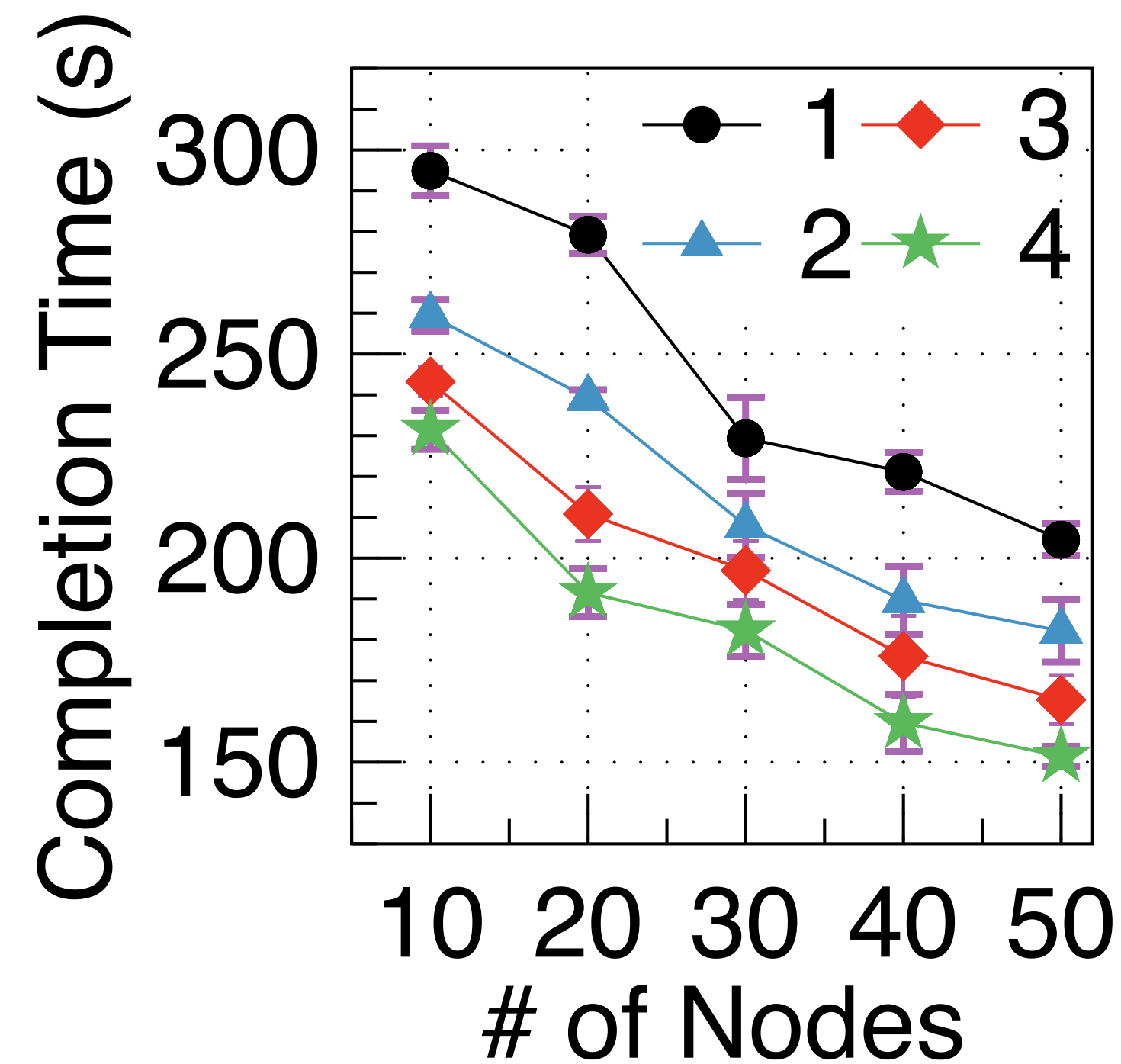(b) Freyr
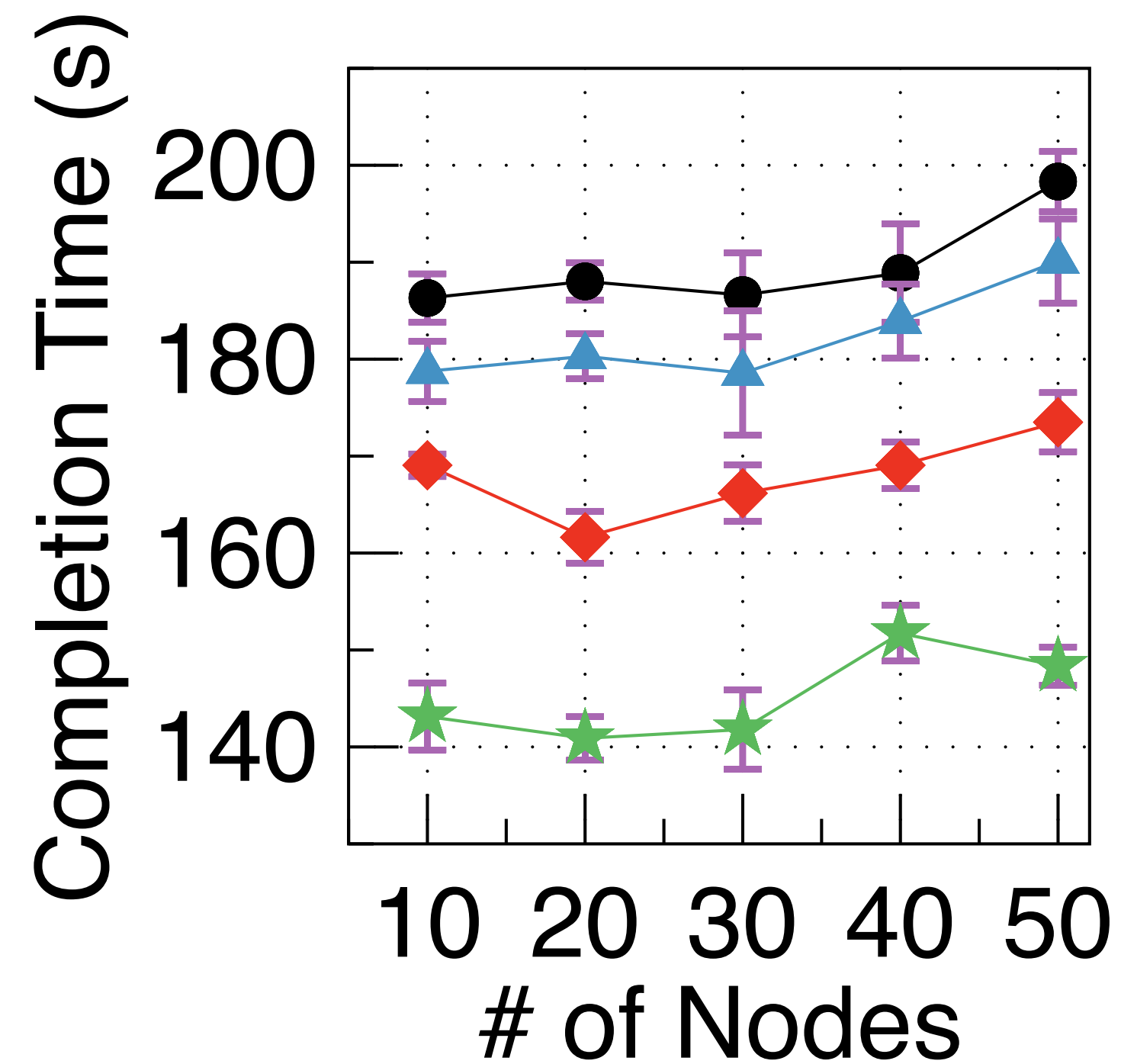
(c) Libra

# Harvesting and Safeguarding

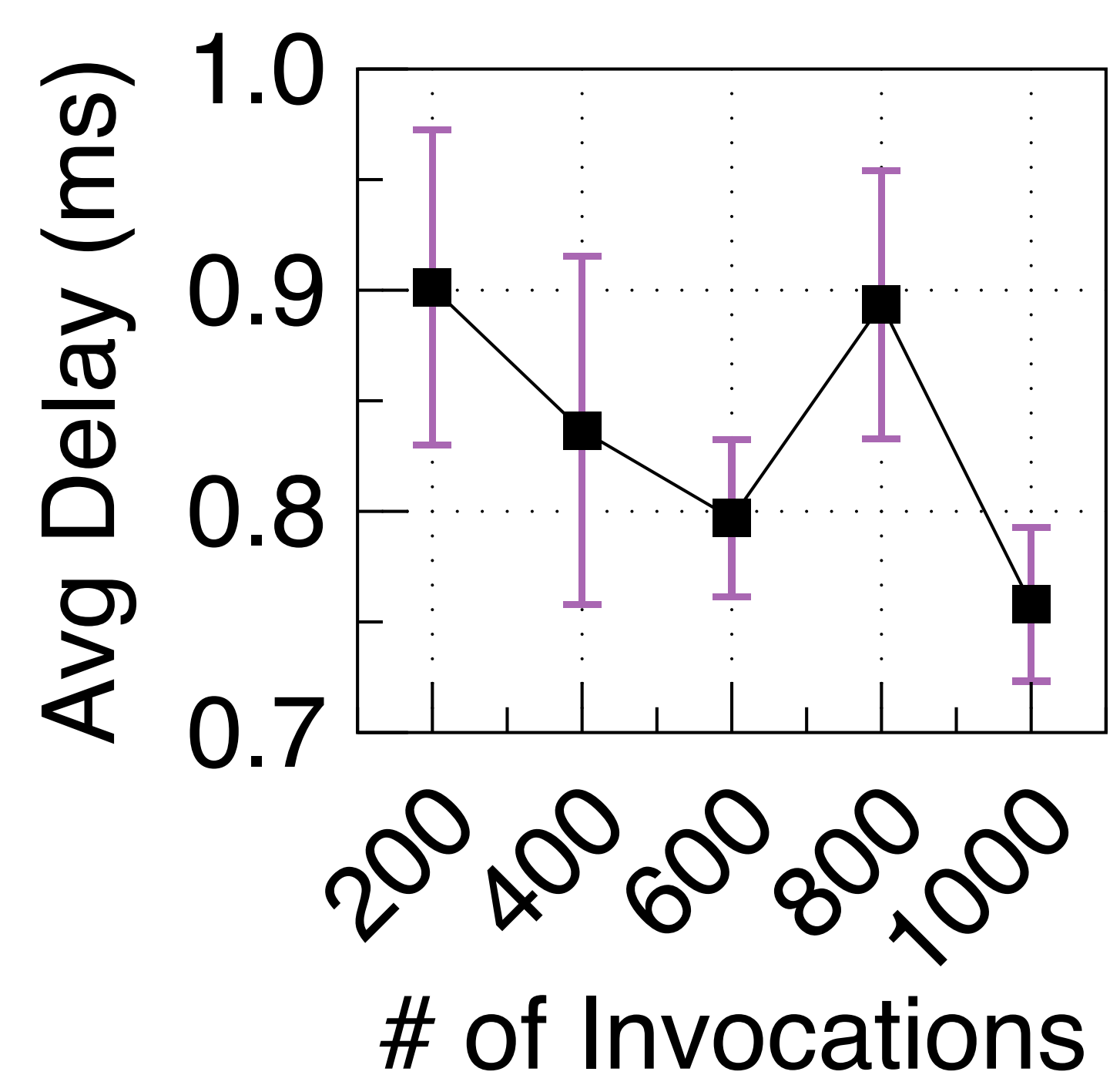$$Speedup := \frac{t^{user} - t^*}{t^{user}}$$

33

# Scalability & Overhead



(a) Strong Scaling

(b) Weak Scaling

(c) Sched. Delay

**Timeliness-aware resource harvesting & scheduling**

**Input data size-awareness**

**Timely safeguard**

# Libra

## 39%
lower function response latency

## 3x
higher resource utilization

**LSU**

**Libra Code Repo:**
https://github.com/IntelliSys-Lab/Libra-HPDC23

**Corresponding Author:**
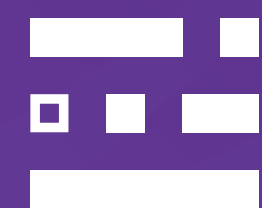Hanfei Yu <hyu25@lsu.edu>
Hao Wang <haowang@lsu.edu>

NSF | Jetstream2 | IntelliSys Lab