# Accelerating Serverless Computing by Harvesting Idle Resources

**Hanfei Yu**[1], Hao Wang[1], Jian Li[2], Xu Yuan[3], Seung-Jong Park[1]
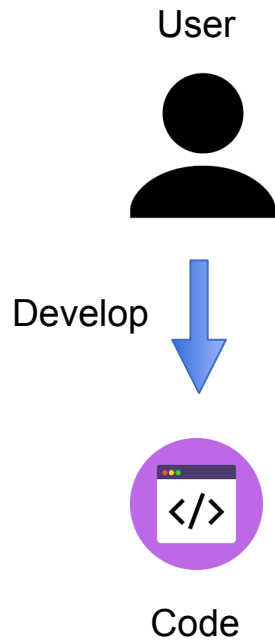
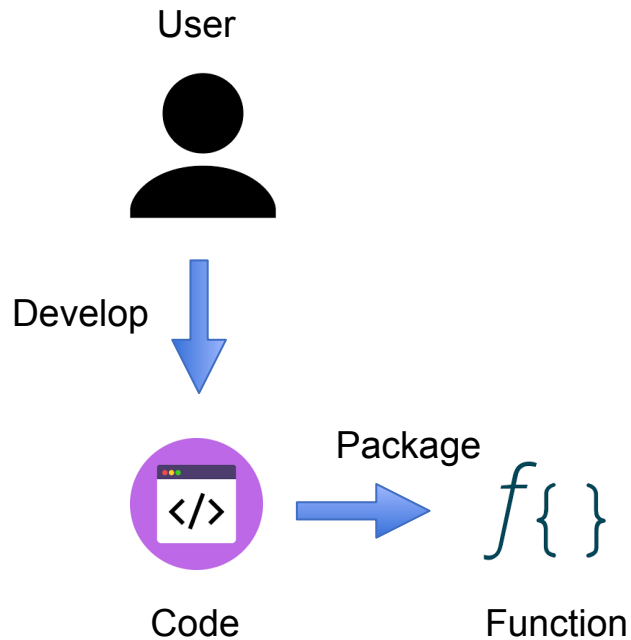[1]Louisiana State University, [2]SUNY-Binghamton University, [3]University of Louisiana at Lafayette

# Serverless Computing



User

Develop

Code

# Serverless Computing

User

Develop

Package

Code

Function

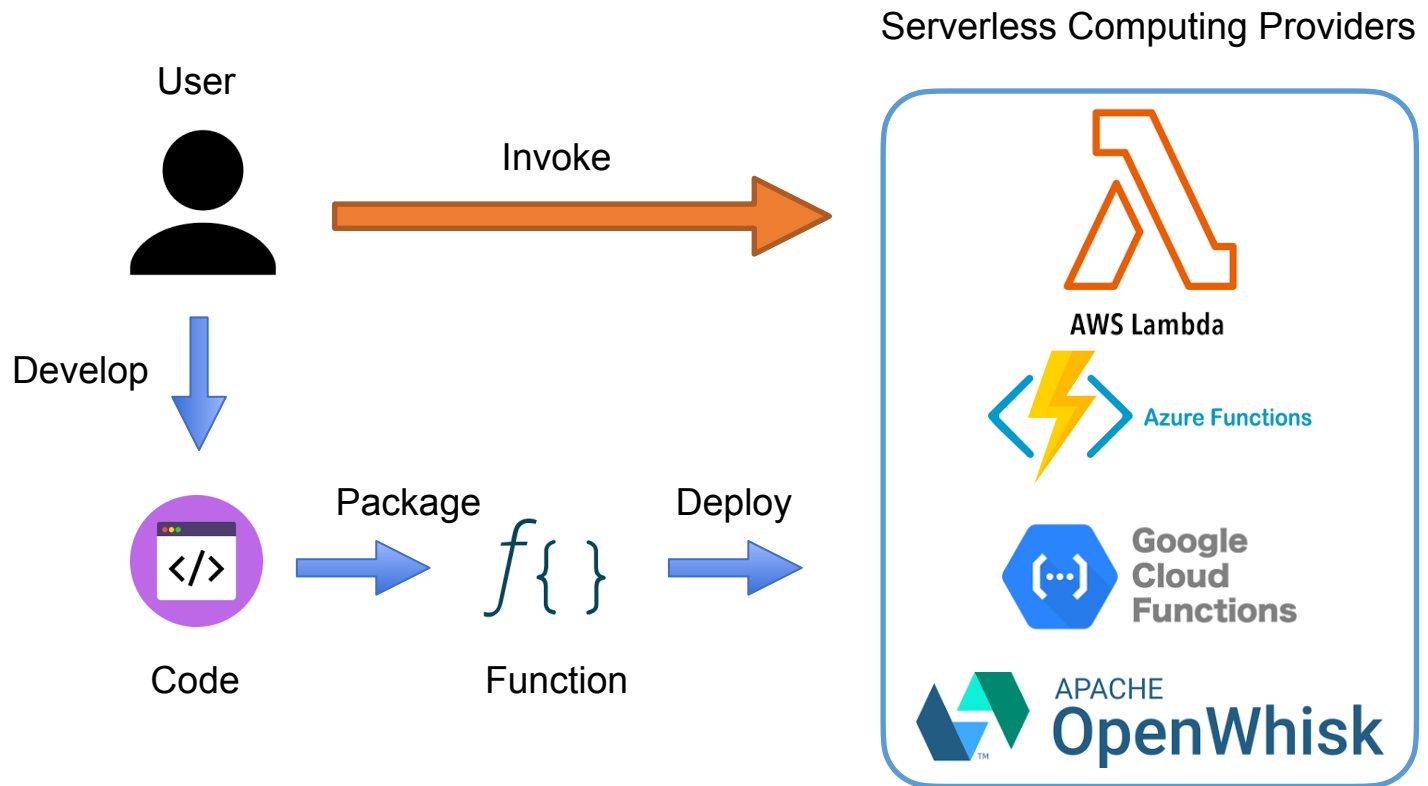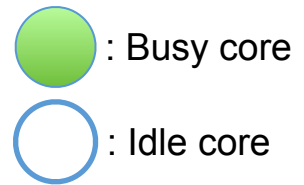# Serverless Computing

# Serverless Computing

# Function Static Configuration

User

Function

Static Config →

$f\{\}$ ◯◯

2 CPU cores

# Varying Input Data Size

# Varying Input Data Size

User

Function

Static Config →

$f\{\ \}$ ○ ○

2 CPU cores

User

Invoke

Large data

Invocation #1

Response Latency

$f\{\ \}$ ● ●

5 sec

Small data

Invocation #2

$f\{\ \}$ ● ○

1 sec

# Harvesting & Acceleration



: Busy core

: Idle core

User

Invoke

Large data

Invocation #1

$f\{\}$  5 sec

Small data

Invocation #2

$f\{\}$  1 sec

Response Latency

# Harvesting & Acceleration

# Harvesting & Acceleration

# Realistic Applications



EG: email generation
KNN: K nearest neighbors

# Realistic Applications

**Performance stops growing when supplying more resources!**



EG: email generation
KNN: K nearest neighbors

# Realistic Harvesting & Acceleration



(a) CPU allocation

(b) Function response latency

EG: email generation
IR: image recognition
ALU: arithmetic logic units
KNN: K nearest neighbors

# Realistic Harvesting & Acceleration

**Latency can be reduced with supplying harvested resources!**



(a) CPU allocation

(b) Function response latency

EG: email generation
IR: image recognition
ALU: arithmetic logic units
KNN: K nearest neighbors

# Realistic Harvesting & Acceleration

**Latency can be reduced with supplying harvested resources!**



(a) CPU allocation

(b) Function response latency

EG: email generation
IR: image recognition
ALU: arithmetic logic units
KNN: K nearest neighbors

**Careful harvesting does not degrade performance**

# General Rebalance Cases

Idle core

Idle memory

# Function index

\* Invocation index

**Case 1**

$f_{\{\#1\}}^{*1}$ → ← $f_{\{\#1\}}^{*2}$

Within the same function
Both donator and receiver

Serverless Cluster

# General Rebalance Cases

Idle core

Idle memory

\# Function index

\* Invocation index

**Case 1**

$f_{\{\#1\}}^{*1}$  →  ←  $f_{\{\#1\}}^{*2}$  □

Within the same function
Both donator and receiver

**Case 2**

$f_{\{\#1\}}^{*1}$  →  →  $f_{\{\#1\}}^{*2}$  □  ○

Within the same function
One donator and one receiver

Serverless Cluster

# General Rebalance Cases

○ Idle core

□ Idle memory

\# Function index

\* Invocation index

## Case 1

○ $f^{*1}_{\{\#1\}}$ ⟹ $f^{*2}_{\{\#1\}}$ □

Within the same function
Both donator and receiver

## Case 2

$f^{*1}_{\{\#1\}}$ ⟹ $f^{*2}_{\{\#1\}}$ □ ○

Within the same function
One donator and one receiver

## Case 3

○ $f^{*1}_{\{\#1\}}$ ⟸ $f^{*1}_{\{\#2\}}$ □

Between two functions
Both donator and receiver

Serverless Cluster

# General Rebalance Cases

○ Idle core

□ Idle memory

# Function index

* Invocation index

**Serverless Cluster**

## Case 1

$f^{*1}_{\{#1\}}$ ⇒⇐ $f^{*2}_{\{#1\}}$ □

Within the same function
Both donator and receiver

## Case 2

$f^{*1}_{\{#1\}}$ ⇒⇒ $f^{*2}_{\{#1\}}$ □ ○

Within the same function
One donator and one receiver

## Case 3

○ $f^{*1}_{\{#1\}}$ ⇐⇒ $f^{*1}_{\{#2\}}$ □

Between two functions
Both donator and receiver

## Case 4

□ ○ $f^{*1}_{\{#1\}}$ ⇐⇐ $f^{*1}_{\{#2\}}$

Between two functions
One donator and one receiver

# Dynamic Decisions

Perspective of a serverless platform:

- Varying **functions**

# Dynamic Decisions

Perspective of a serverless platform:

- Varying **functions**

- Varying **invocations** per functions

# Dynamic Decisions

Perspective of a serverless platform:

- Varying **functions**

- Varying **invocations** per functions

- Varying **input data** per invocation

# Dynamic Decisions

Perspective of a serverless platform:

- Varying **functions**

- Varying **invocations** per functions

- Varying **input data** per invocation

- Every invocation requires an **allocation decision**

# Dynamic Decisions

Perspective of a serverless platform:

- Varying **functions**

- Varying **invocations** per functions

- Varying **input data** per invocation

- Every invocation requires an **allocation decision**

**A series of sequential allocation decisions**

# Dynamic Decisions

Perspective of a serverless platform:

- Varying **functions**

- Varying **invocations** per functions

- Varying **input data** per invocation

- Every invocation requires an **allocation decision**

**A series of sequential allocation decisions**

**Markov Decision Process (MDP)**

# Deep Reinforcement Learning

Perspective of a serverless platform:

- Varying **functions**

- Varying **invocations** per functions

- Varying **input data** per invocation

- Every invocation requires an **allocation decision**

**A series of sequential allocation decisions**

**Markov Decision Process (MDP)**

**Deep Reinforcement Learning**

# Deep Reinforcement Learning



Latency Rewards

Invoke $f\{\}$

Environment

Agent

Allocation Actions

# Freyr

# Freyr Workflow

**Embedding**

**Scoring**

**Selection**

Platform

Inflight invocations

Available CPU

Available Memory

Concatenate

State vector $s^1$

Function

Avg CPU peak

Avg memory peak

…

Baseline exec time

Alloc option 1

Platform

Platform info…

Function

Function info…

Alloc option N

Concatenate

State vector $s^N$

**Score network**

$q^1$

$\vdots$

$q^N$

Actor network

Softmax

$\mathbb{P}(\text{option})$

Safeguard

Best Alloc Option

$b^1$

$\vdots$

$b^N$

Critic network

Mean

$\bar{b}$

# Freyr Workflow



**State information from the platform and the function**

# Freyr Workflow



**Proximal Policy Optimization (PPO)**

# Freyr Workflow



**Safeguard**
- **Filter invalid allocation options**
- **Return resources when detecting a potential full usage**

# Freyr Architecture

Func req 1 🟩 🟦   …   Func req N 🟩 🟦

**Frontend**

**Controller**

**Load Balancer**

*Freyr*

Safeguard — DRL Agent

*allocation*

Database → *state* → KV Storage

*(results, usage)*

Distributed Message Queue
(Pub/Sub)

*(CPU, mem) allocation*

**Invoker**

λ λ λ ← Waiting Queue

λ → Container

🟩🟩🟩🟩⬜⬜⬜⬜ CPU
🟦🟦🟦🟦🟦🟦⬜⬜ Mem

# Freyr Architecture

**Frontend receives function invocations from users**

# Freyr Architecture



**Controller collects and sends states to the DRL agent**

# Freyr Architecture



Func req 1    ...    Func req N

Frontend

**Controller**

Load Balancer

Database

*(results, usage)*

***Freyr***

Safeguard — DRL Agent

*allocation*

*state* → KV Storage

**Agent predicts an allocation and sends it back to controller**

Distributed Message Queue
(Pub/Sub)

*(CPU, mem) allocation*

**Invoker**

λ λ λ ← Waiting Queue

λ → Container

CPU

Mem

# Freyr Architecture



Func req 1 ☐ ☐  …  Func req N ☐ ☐

Frontend

**Controller**
Load Balancer

*Freyr*
Safeguard — DRL Agent
*allocation*

Database → *state* → KV Storage

**Controller then forwards the function invocation with its decision to an Invoker**

*(results, usage)*

Distributed Message Queue (Pub/Sub)

*(CPU, mem) allocation*

**Invoker**
λ λ λ ← Waiting Queue
λ → Container
CPU
Mem

# Freyr Architecture



Func req 1 ... Func req N

Frontend

**Controller**
Load Balancer

*Freyr*
Safeguard — DRL Agent
*allocation*

Database → *state* → KV Storage

Distributed Message Queue
(Pub/Sub)

*(results, usage)*

*(CPU, mem) allocation*

**Invoker**
λ λ λ ← Waiting Queue

λ →
Container

CPU
Mem

**Invoker executes
the function**

# Freyr Architecture



Func req 1 · · · Func req N

Frontend

**Controller**
Load Balancer

*Freyr*
Safeguard — DRL Agent

*allocation*

Database — *state* → KV Storage

*(results, usage)*

**Invoker submits the results and usage to database for further predictions**

Distributed Message Queue (Pub/Sub)

*(CPU, mem) allocation*

**Invoker**

λ λ λ ← Waiting Queue

λ → Container

CPU
Mem

# Experiment

## Setup

- 13 VMs, each with 8 CPUs and 32 GB memory
- One user client, one frontend, one controller
- 10 Worker nodes

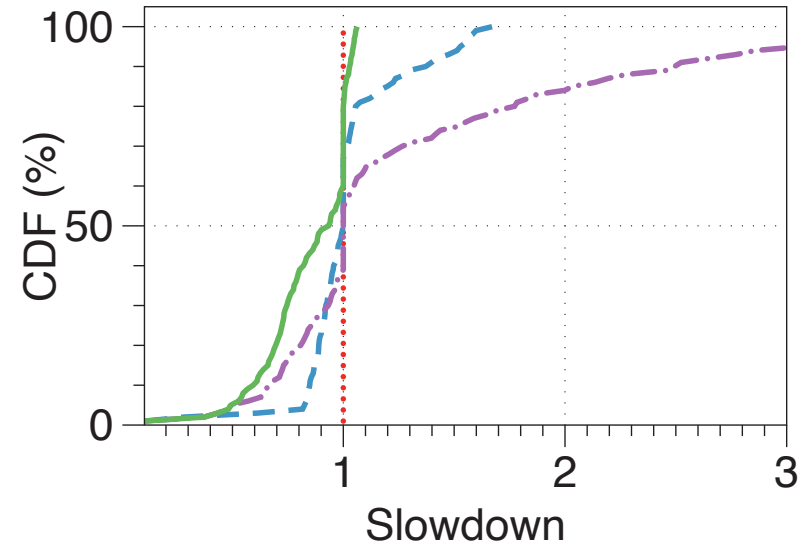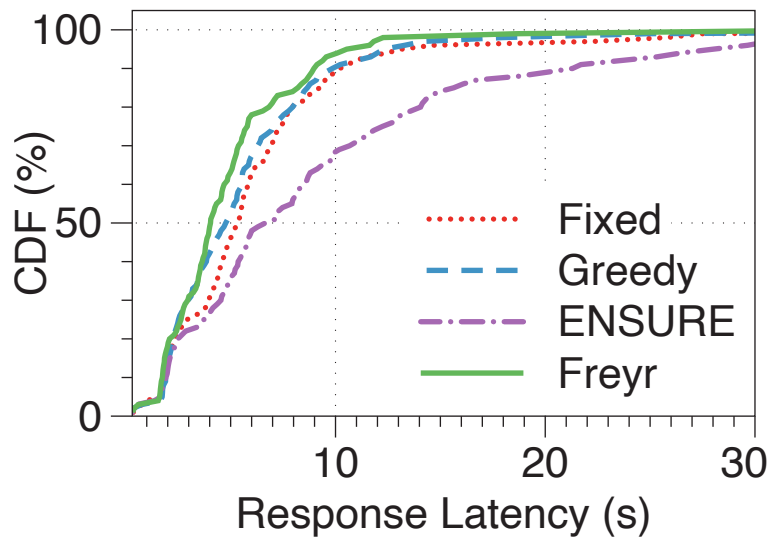## Baselines

- Fixed RM
- Greedy RM
- ENSURE

**Fixed RM**: default OpenWhisk as well as in existing serverless platforms
**Greedy RM**: heuristic
**ENSURE**: Suresh, Amoghavarsha, et al.
"Ensure: Efficient scheduling and autonomous resource management in serverless environments."
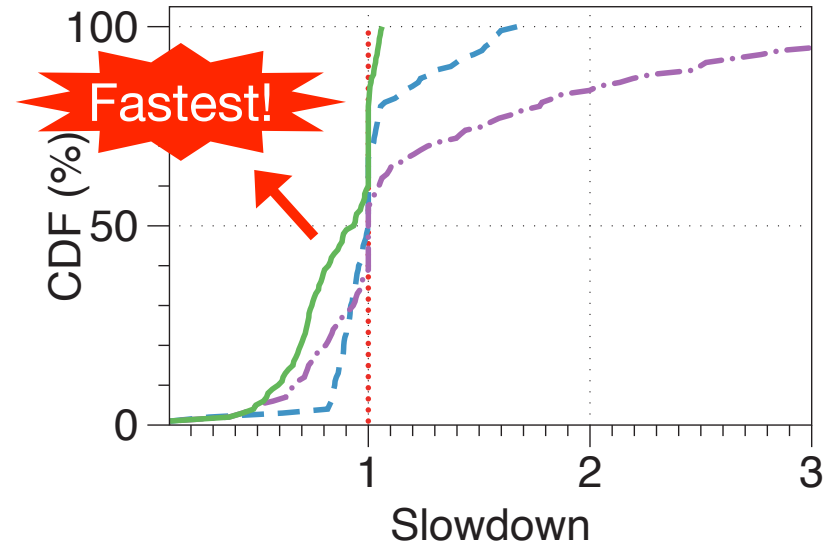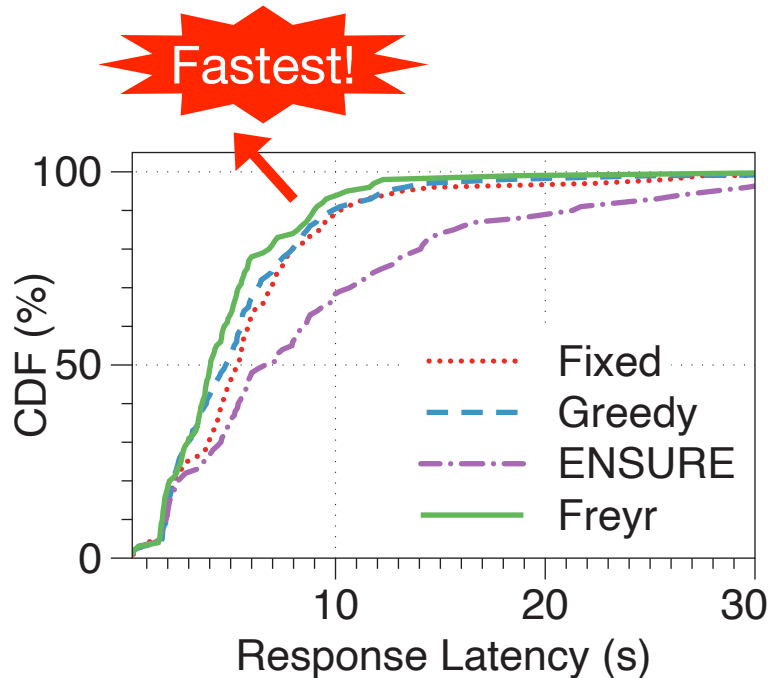(*ACSOS 2020)*

# Function Execution Speedup



**Response latency**: function invocation end-to-end latency
**Slowdown**: relative performance compared to user-defined resources.
Larger than1.0 means degradation, less than 1.0 means speedup

# Function Execution Speedup



**Response latency**: function invocation end-to-end latency
**Slowdown**: relative performance compared to user-defined resources.
Larger than1.0 means degradation, less than 1.0 means speedup
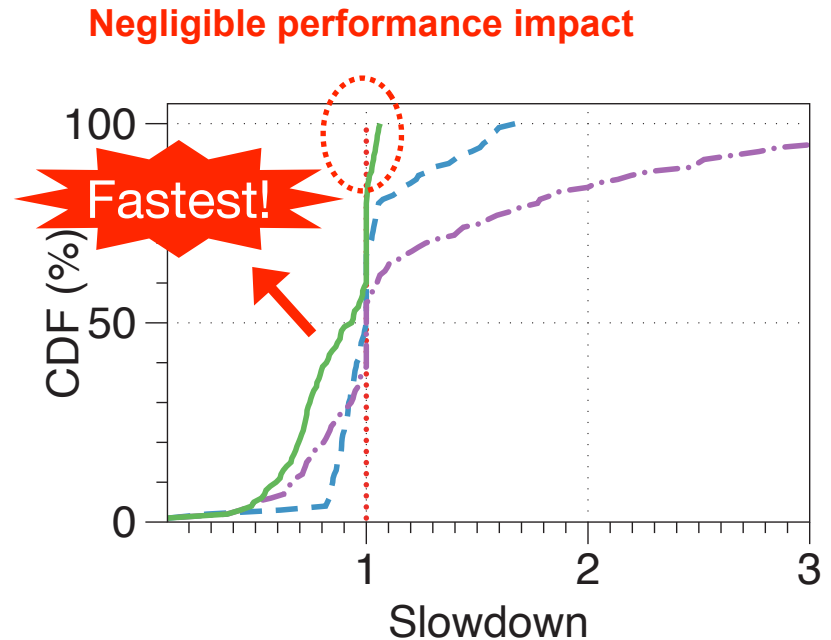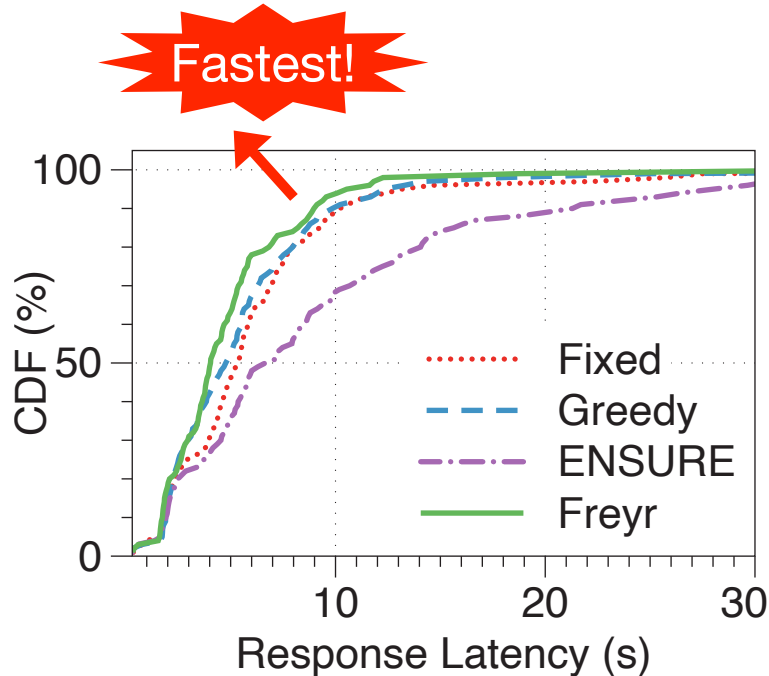
# Function Execution Speedup



**Response latency**: function invocation end-to-end latency
**Slowdown**: relative performance compared to user-defined resources. Larger than 1.0 means degradation, less than 1.0 means speedup.
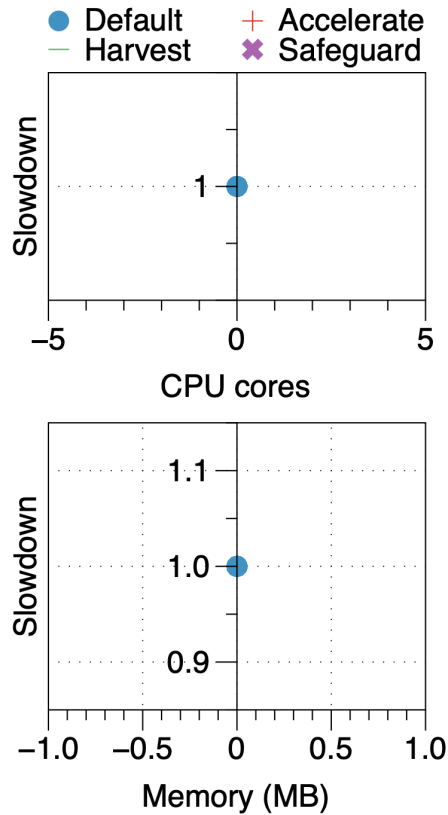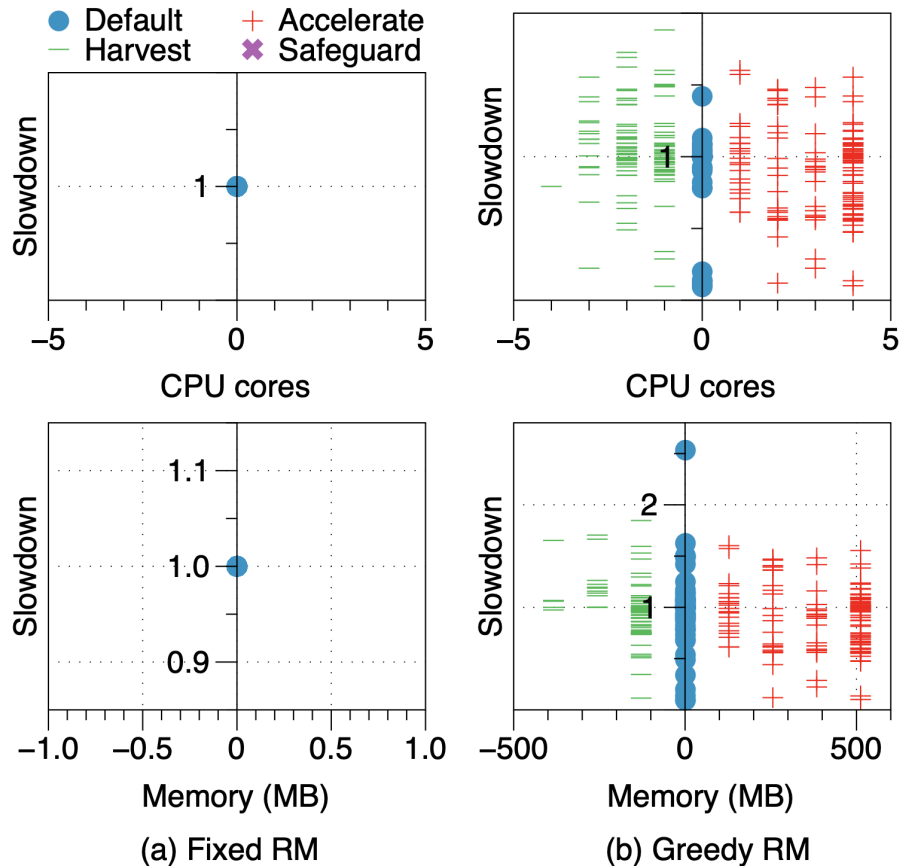
# Resource Allocation

Default ●    Accelerate +
Harvest —    Safeguard ✖

**Slowdown** vs **CPU cores**

**Slowdown** vs **Memory (MB)**

(a) Fixed RM

**Slowdown**: relative performance compared to user-defined resources. Larger than 1.0 means degradation, less than 1.0 means speedup.
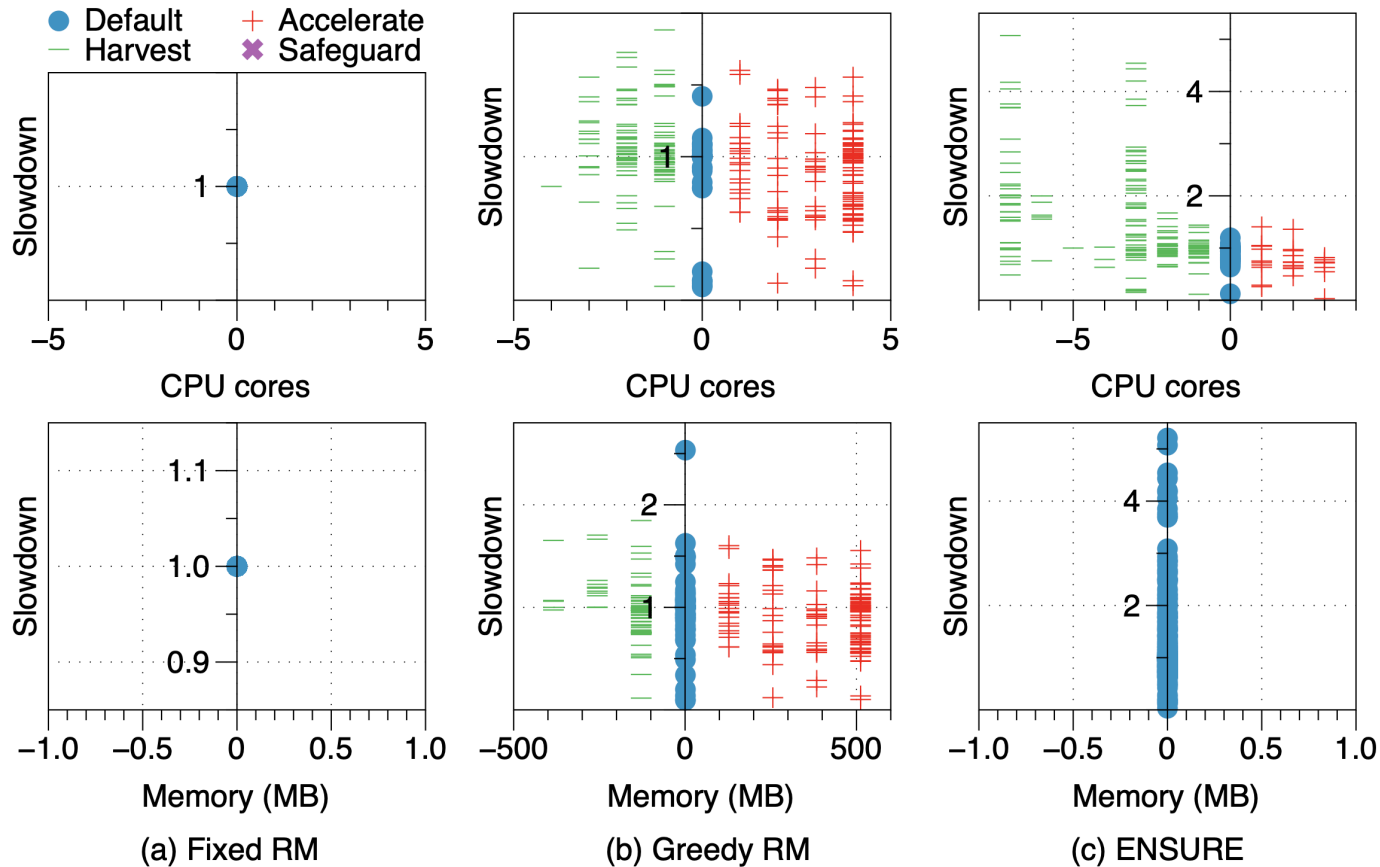
# Resource Allocation



(a) Fixed RM

(b) Greedy RM

**Slowdown**: relative performance compared to user-defined resources. Larger than 1.0 means degradation, less than 1.0 means speedup.

# Resource Allocation



**Slowdown**: relative performance compared to user-defined resources. Larger than 1.0 means degradation, less than 1.0 means speedup.
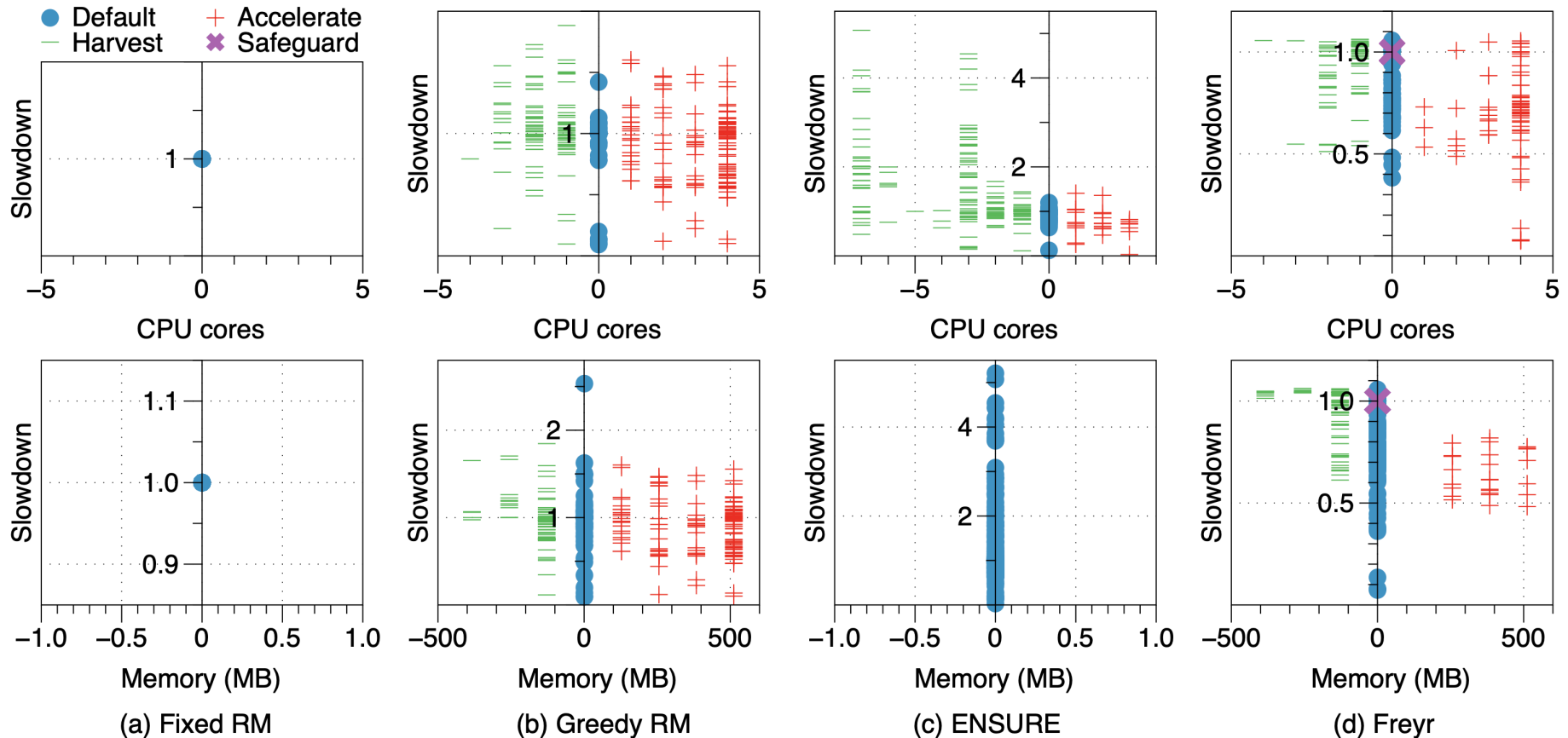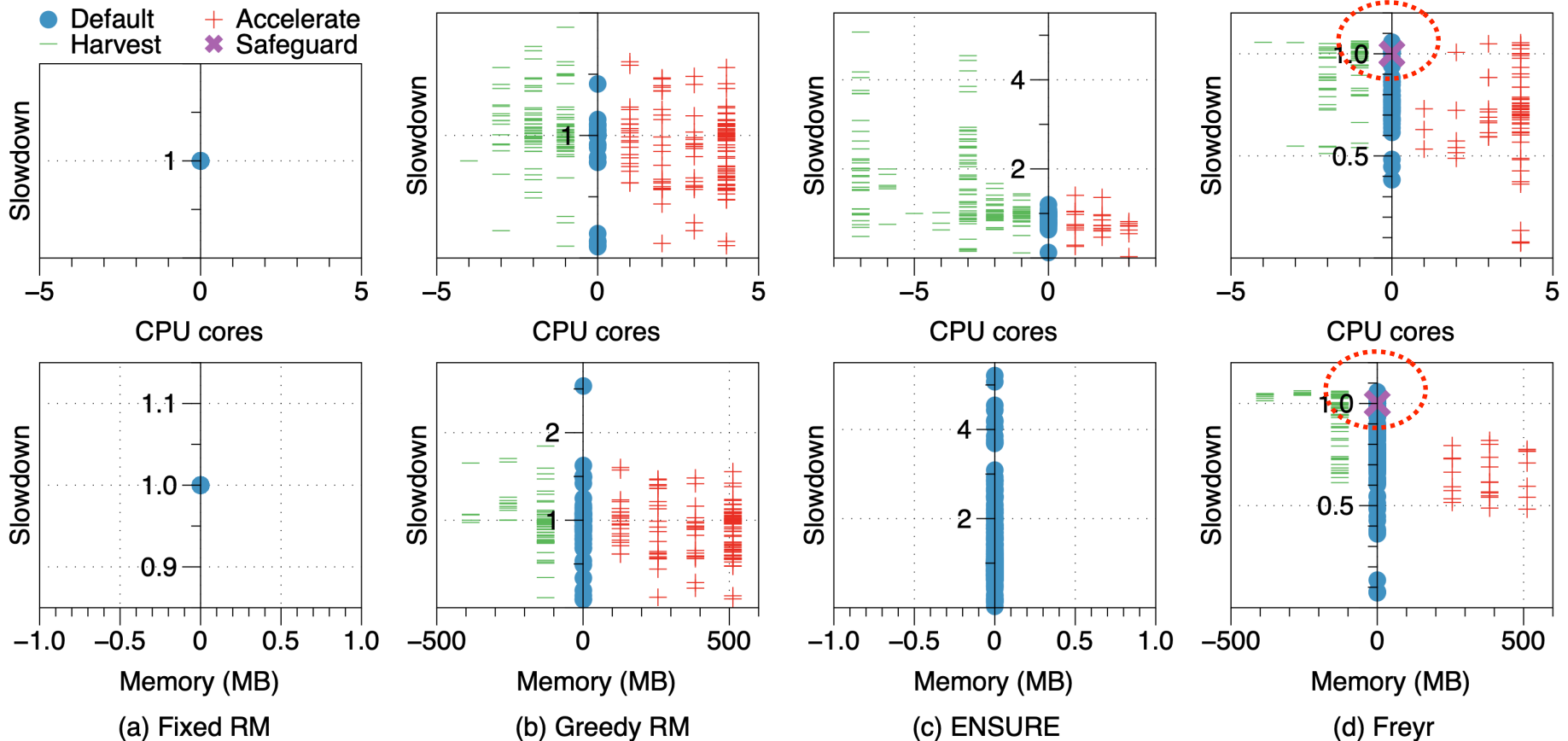
# Resource Allocation

**Slowdown**: relative performance compared to user-defined resources. Larger than 1.0 means degradation, less than 1.0 means speedup.

# Resource Allocation

(a) Fixed RM

(b) Greedy RM

(c) ENSURE

(d) Freyr

**Slowdown**: relative performance compared to user-defined resources. Larger than1.0 means degradation, less than 1.0 means speedup.

# Thank You

**Hanfei Yu**: [hyu25@lsu.edu](mailto:hyu25@lsu.edu)

**IntelliSys Lab**: https://intellisys.haow.ca