# AccelFlow: Orchestrating an On-Package Ensemble of Fine-Grained Accelerators for Microservices
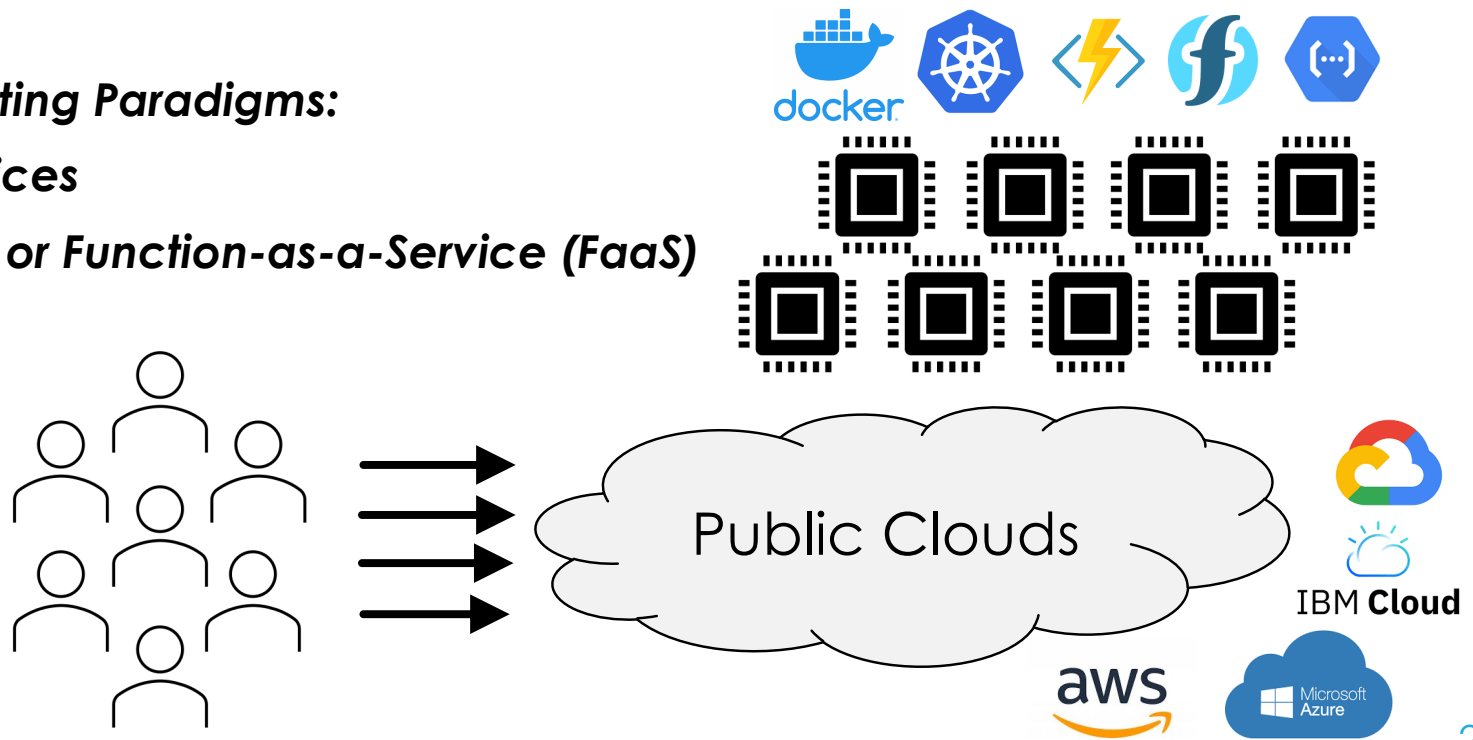
## HPCA 2026, Sydney, Australia

Jovan Stojkovic*, Abraham Farrell, Zhangxiaowen Gong[†], Christopher J. Hughes[†], Josep Torrellas

University of Illinois at Urbana-Champaign, [†]Intel, *Joining UT Austin in Fall 2026
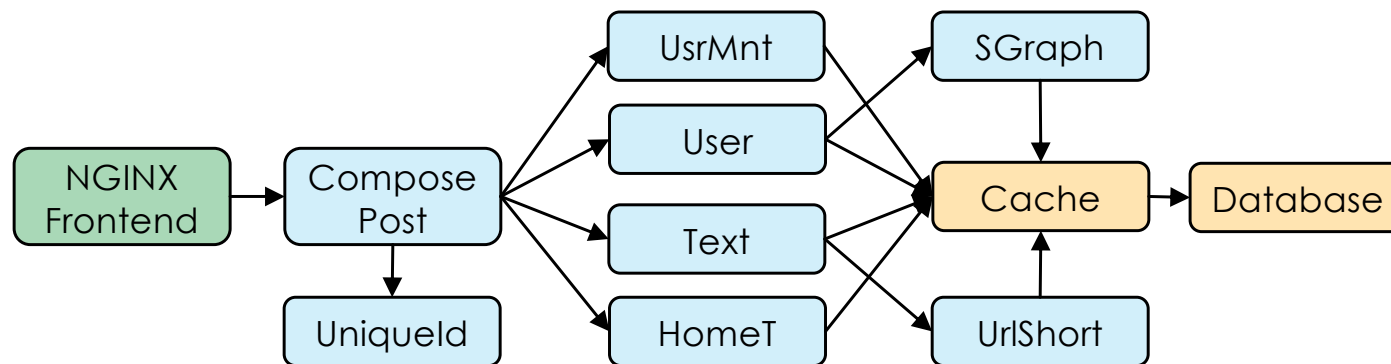
# The Growth of Cloud Computing

*New Computing Paradigms:*

- *Microservices*

- *Serverless or Function-as-a-Service (FaaS)*

Public Clouds

# Microservices

o Large monolithic applications decomposed into many small interdependent services
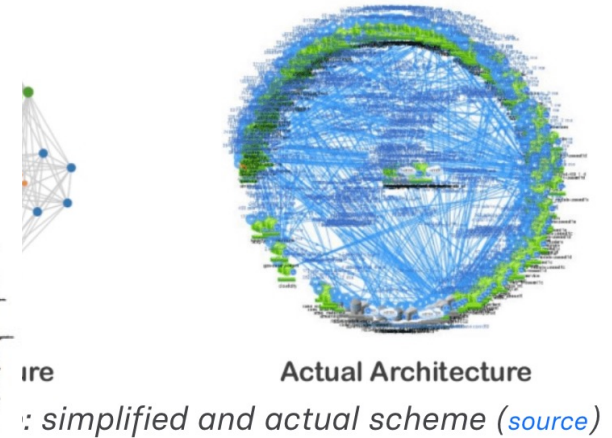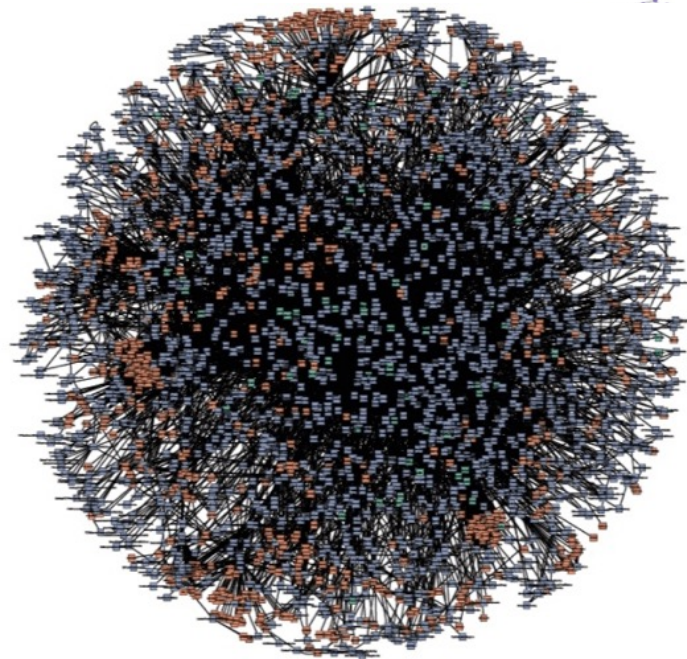
  o Each service implements separate functionality

# Benefits of Microservices
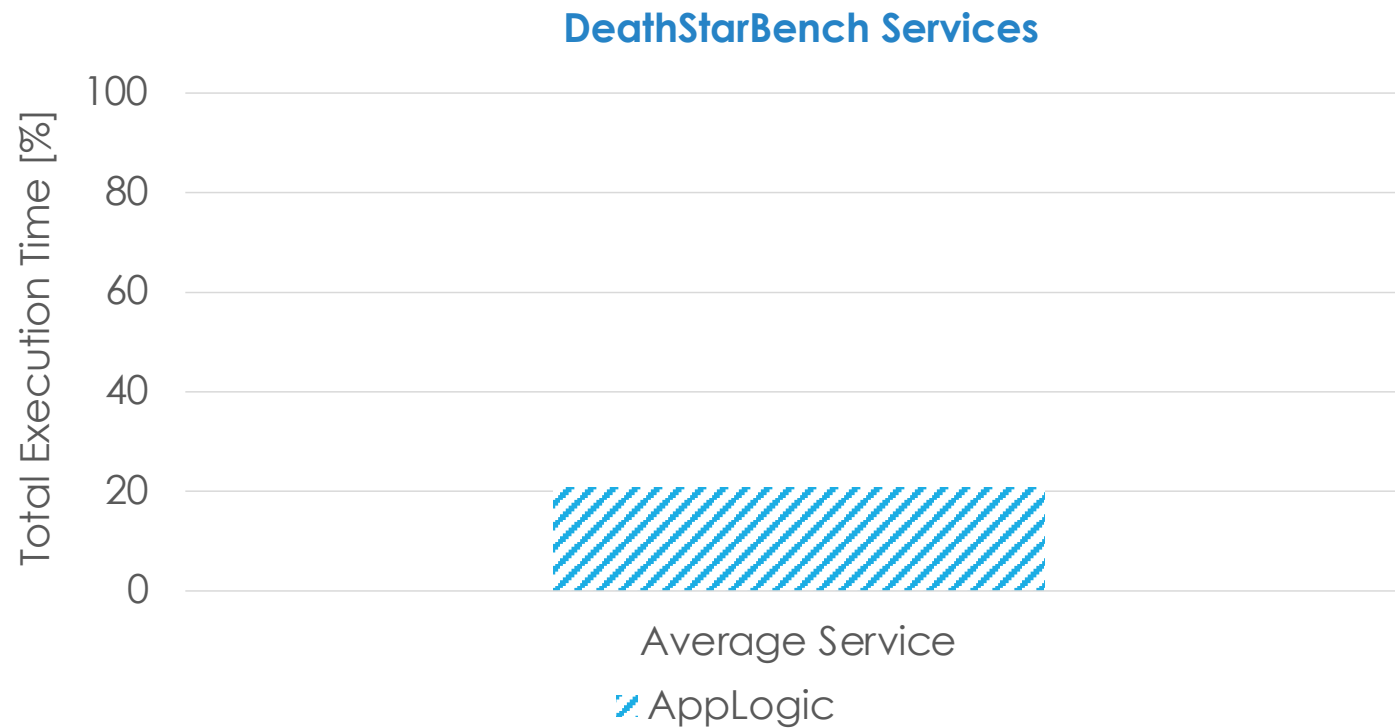
o Scalability ☑

o Design simplicity ☑

o HW management ☑

# Microservices are Widely Used


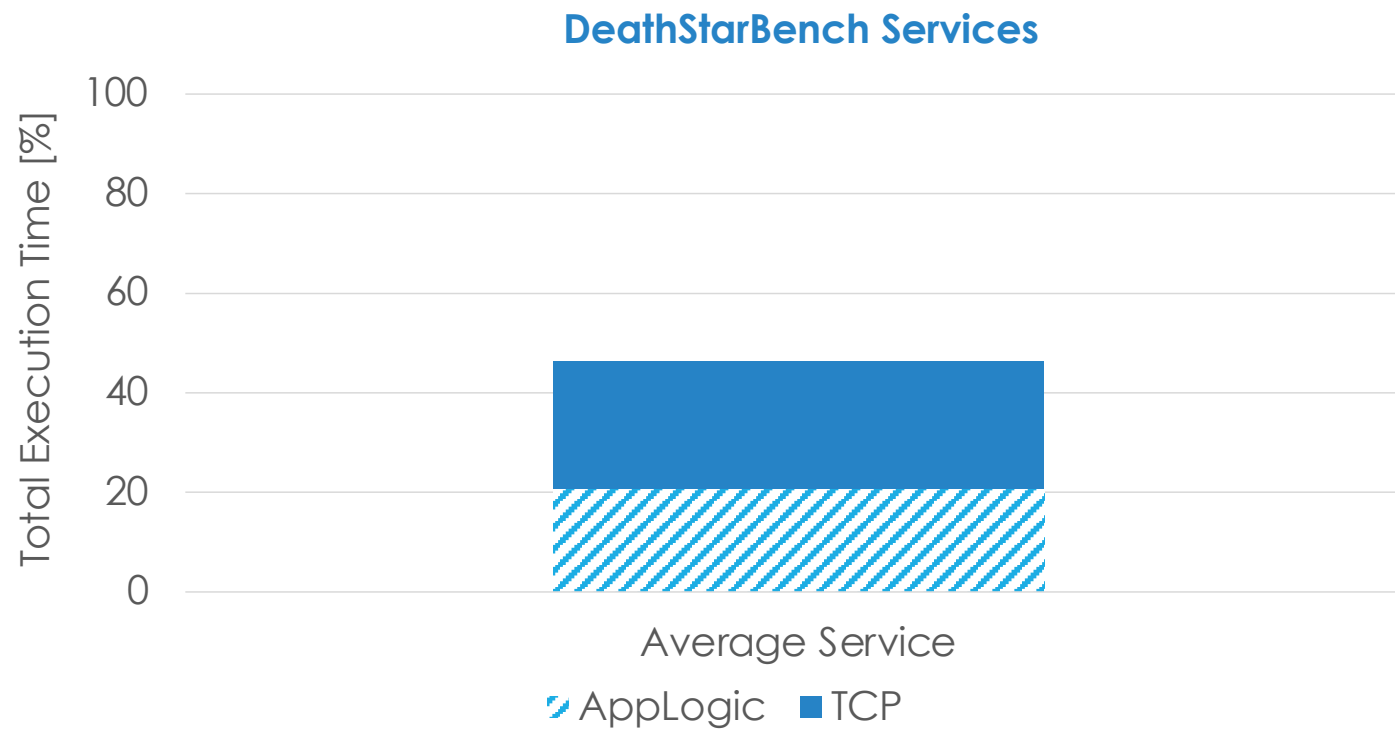
Actual Architecture

*: simplified and actual scheme (source)*

*Structure of microservices at Amazon. Looks almost like a Death Star but is way more powerful.*

5

# Datacenter Tax Dominates Execution
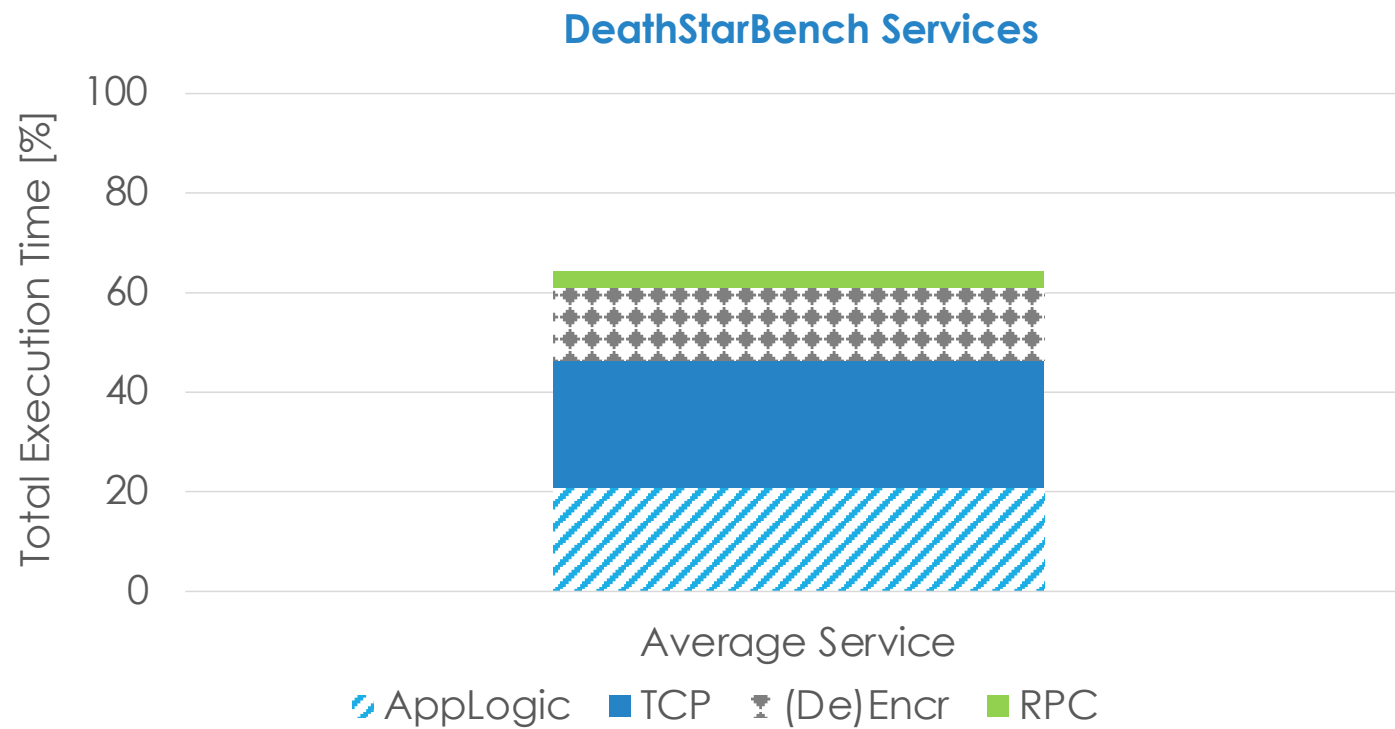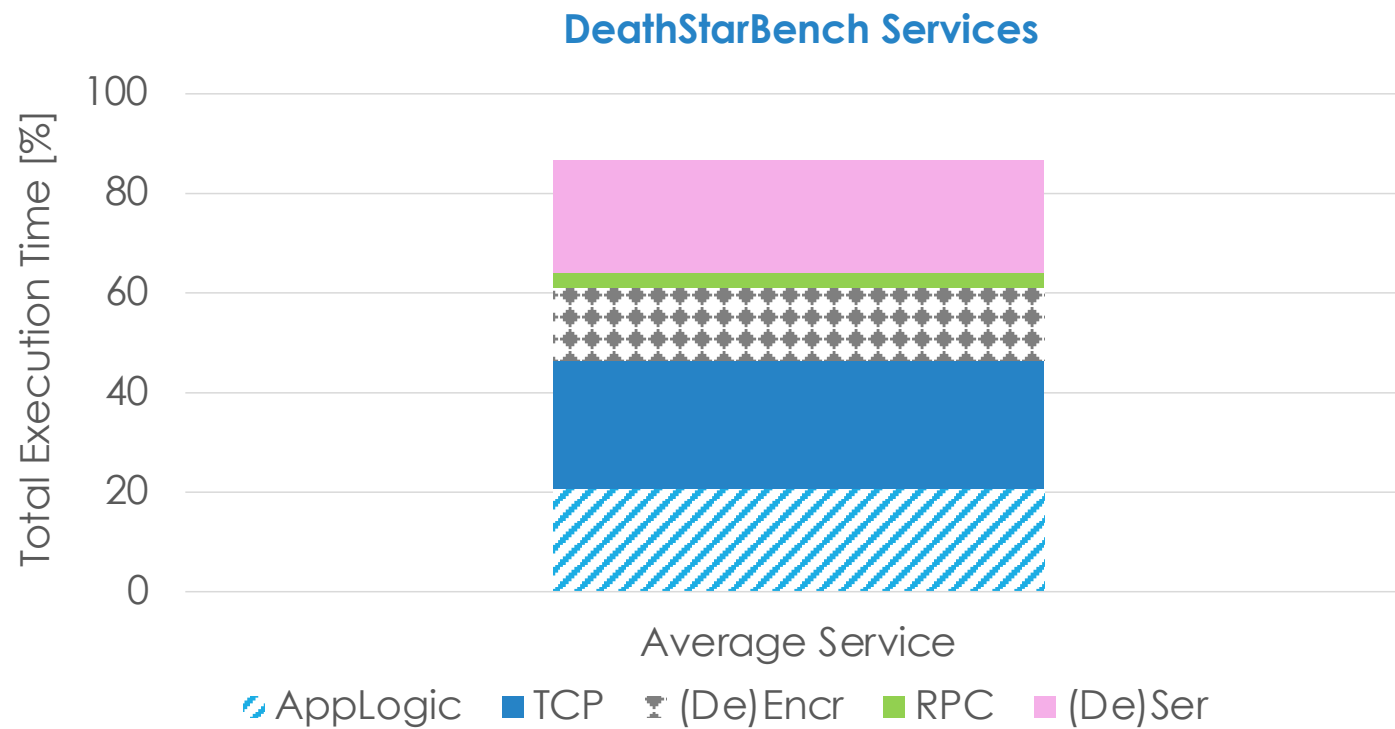
**DeathStarBench Services**



Y-axis: Total Execution Time [%] (0, 20, 40, 60, 80, 100)

X-axis: Average Service

Legend: AppLogic

# Datacenter Tax Dominates Execution

**DeathStarBench Services**



7

# Datacenter Tax Dominates Execution

## DeathStarBench Services



Chart: Total Execution Time [%] vs Average Service. Stacked bar shows AppLogic (0–20%), TCP (20–47%), and (De)Encr (47–60%).

Legend: AppLogic, TCP, (De)Encr

8

# Datacenter Tax Dominates Execution

**DeathStarBench Services**



Chart: Total Execution Time [%] (y-axis, 0 to 100) for Average Service (x-axis), stacked bar showing AppLogic, TCP, (De)Encr, and RPC components.

Legend: AppLogic | TCP | (De)Encr | RPC

# Datacenter Tax Dominates Execution

**DeathStarBench Services**



*Chart: Total Execution Time [%] (y-axis, 0 to 100) vs Average Service (x-axis). Legend: AppLogic, TCP, (De)Encr, RPC, (De)Ser*

10

# Datacenter Tax Dominates Execution



**DeathStarBench Services**

Total Execution Time [%]

100
80
60
40
20
0

Average Service

AppLogic    TCP    (De)Encr    RPC    (De)Ser    (De)Cmp

11

# Datacenter Tax Dominates Execution



**DeathStarBench Services**

Total Execution Time [%]

100 80 60 40 20 0

Average Service

AppLogic   TCP   (De)Encr   RPC   (De)Ser   (De)Cmp   LdB

12

# Datacenter Tax Dominates Execution

**DeathStarBench Services**



Lots of CPU cycles spent on datacenter tax operations

Total Execution Time [%]

100
80
60
40
20
0

Average Service

AppLogic ■ TCP (De)Encr ■ RPC ■ (De)Ser (De)Cmp ■ LdB

13

# Datacenter Tax Reported by Major Hyperscalers

## Google 2015



**Figure 4: 22-27% of WSC cycles are spent in different components of "datacenter tax".**

## Facebook 2020



**Figure 1.** Breakdown of cycles spent in core application logic vs. orchestration work: orchestration overheads can significantly dominate.

## Google 2023



**Figure 3: High-Level Application-Level Cycle Breakdown**

14

# Many Proposals for Individual Accelerators



5

# Many Proposals for Individual Accelerators



6

# How to Orchestrate Many Accelerators?

o Many individual accelerators proposed – how to manage them?

# Orchestrate Many Accelerators: CPU-Centric

TCP | Decr | RPC | Dser | Dcmp | LdB

Initiate

# Orchestrate Many Accelerators: CPU-Centric

TCP  Decr  RPC  Dser  Dcmp  LdB

Interrupt

# Orchestrate Many Accelerators: CPU-Centric

| TCP | Decr | RPC | Dser | Dcmp | LdB |

Initiate

# Orchestrate Many Accelerators: CPU-Centric

| TCP | Decr | RPC | Dser | Dcmp | LdB |
|-----|------|-----|------|------|-----|

Interrupt

# Repeated Interrupts → High Overhead

Fraction of time spent orchestrating accelerators



22

# Orchestrate Many Accelerators: Direct Chain

| TCP | Decr | RPC | Dser | Dcmp | LdB |
|-----|------|-----|------|------|-----|

Multiple sources of the tax execute in sequence!

# Orchestrate Many Accelerators: Direct Chain

TCP · Decr · RPC · Dser · Dcmp · LdB

Initiate

# Orchestrate Many Accelerators: Direct Chain

TCP → Decr → RPC → Dser → Dcmp → LdB

# Orchestrate Many Accelerators: Direct Chain

| TCP | Decr | RPC | Dser | Dcmp | LdB |

Interrupt

# Direct Chaining Significantly Reduces Overheads

Fraction of time spent orchestrating accelerators



27

# Challenges of Direct Chaining

o Control-flow divergences

# Challenges of Direct Chaining

o Control-flow divergences

Receive function request

TCP → Decr → RPC → Dser → Compressed? —Y→ Dcmp

N → LdB

LdB → CPU

# Challenges of Direct Chaining

o Control-flow divergences

o Data format transformations

# AccelFlow: Accelerator Orchestration Framework

Processor package

# AccelFlow: Accelerator Orchestration Framework

o Ensemble of accelerators

Processor package

# AccelFlow: Accelerator Orchestration Framework

o Ensemble of accelerators

o Direct inter-accelerator chaining

Processor package

# AccelFlow: Accelerator Orchestration Framework

o Ensemble of accelerators

o Direct inter-accelerator chaining

o Sequence of accelerators stored in **Software "Traces"**

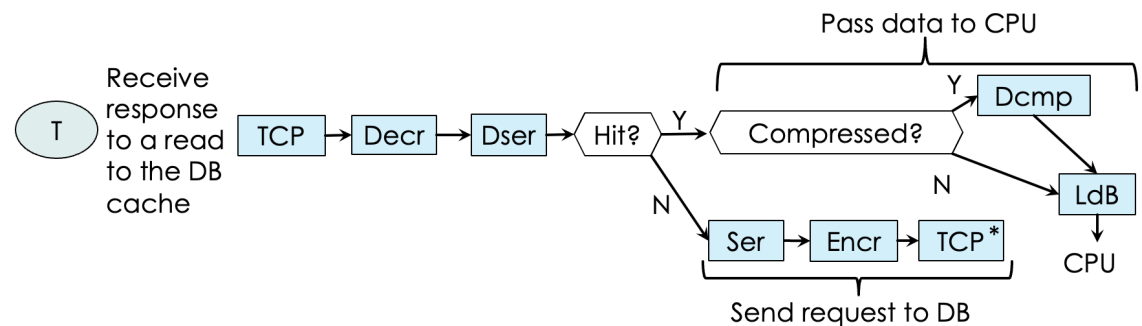# AccelFlow: Accelerator Orchestration Framework

- Ensemble of accelerators
- Direct inter-accelerator chaining
- Sequence of accelerators stored in software "traces"
- Standard interface
  - Input and output queues and dispatchers

Processor package

# Input Dispatcher

o Schedules the requests from Input Queue to PEs

o Fetches large input payloads from memory

o Simple Finite State Machine

# Output Dispatcher

o Forward the request + data to next accelerator or to the CPU

o Compute branch conditions

o Perform data transformations

# Programming AccelFlow

o AccelFlow API allows programmers to construct new traces:

  o Define a linear chain of accelerators

```python
from AFlow import Trace, seq, branch, transform
trace = Trace() # Define trace
pipeline = seq( # Compose trace
  "TCP", "Decr", "RPC", "Dser",
  branch(condition_op="out['compressed'] == 1",
         on_true=seq(trans("JSON", "str"), "Dcmp"),
         on_false=None),
  "LdB")
trace.build(pipeline) # Attach pipeline to trace
trace.register(name="func_req") # Register trace
```
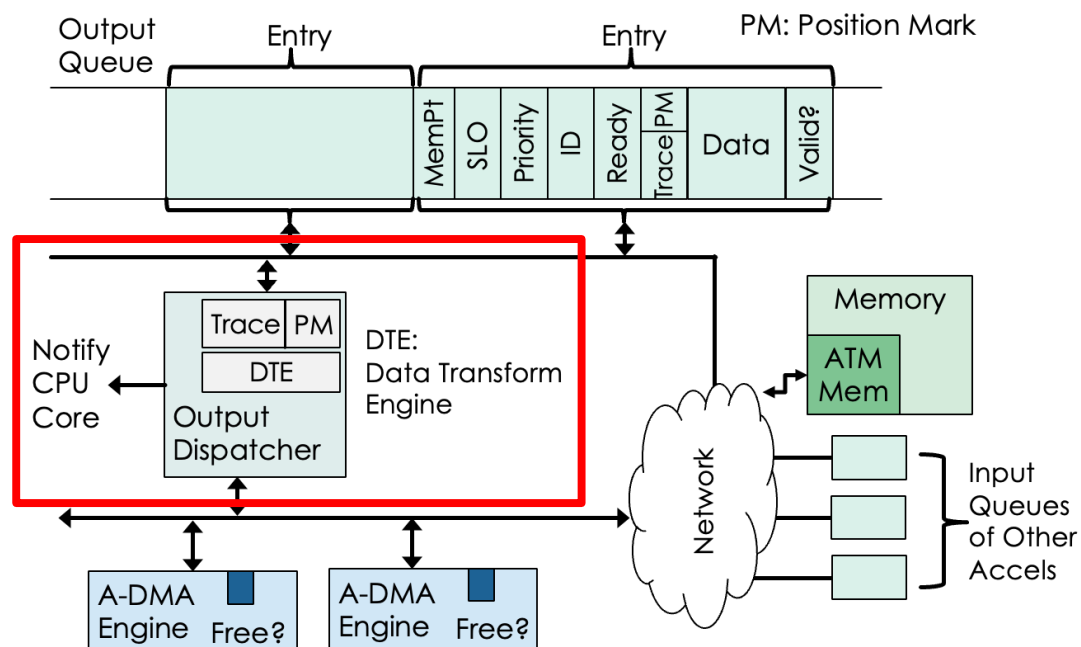
# Programming AccelFlow

o AccelFlow API allows programmers to construct new traces:

- o Define a linear chain of accelerators
- o Add a conditional control flow

```python
from AFlow import Trace, seq, branch, transform
trace = Trace() # Define trace
pipeline = seq( # Compose trace
  "TCP", "Decr", "RPC", "Dser",
  branch(condition_op="out['compressed'] == 1",
         on_true=seq(trans("JSON", "str"), "Dcmp"),
         on_false=None),
  "LdB")
trace.build(pipeline) # Attach pipeline to trace
trace.register(name="func_req") # Register trace
```

# Programming AccelFlow

o AccelFlow API allows programmers to construct new traces:

  o Define a linear chain of accelerators

  o Add a conditional control flow

  o Transform the format of the data from one representation to another

```python
from AFlow import Trace, seq, branch, transform
trace = Trace() # Define trace
pipeline = seq( # Compose trace
  "TCP", "Decr", "RPC", "Dser",
  branch(condition_op="out['compressed'] == 1",
         on_true=seq(trans("JSON", "str"), "Dcmp"),
         on_false=None),
  "LdB")
trace.build(pipeline) # Attach pipeline to trace
trace.register(name="func_req") # Register trace
```
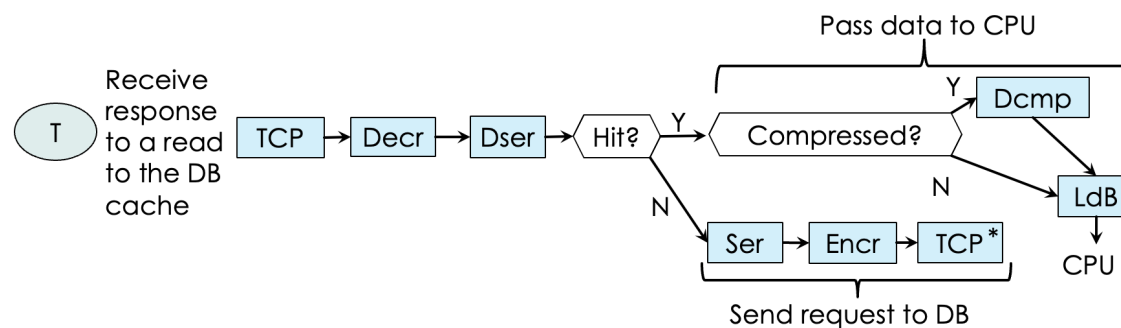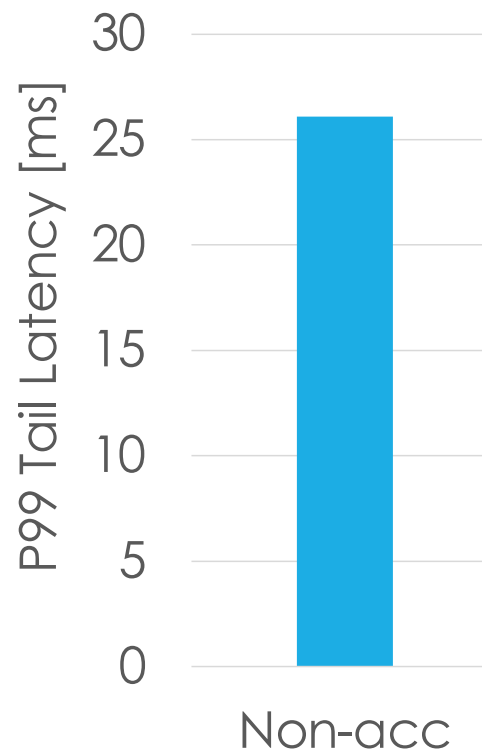
# AccelFlow Summary

o Many on-chip accelerators to reduce datacenter tax

o Accelerators communicate directly with each other

o Small hardware engines

    o Schedule requests onto accelerator PEs

    o Compute branch conditions

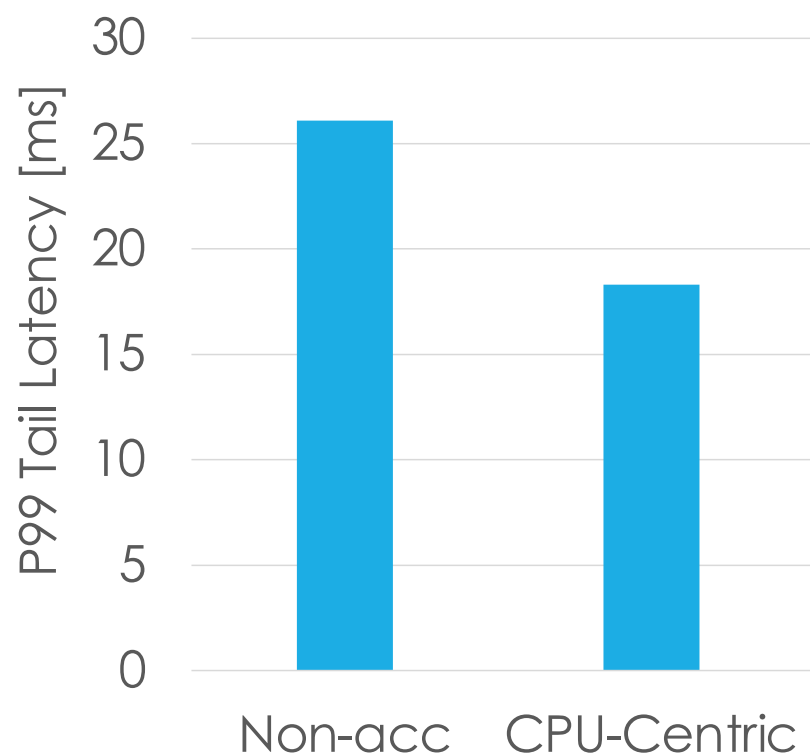    o Perform simple data transformations

# Evaluation Methodology

- Cycle-accurate full-system simulations: SST + QEMU
- DeathStarBench services with Alibaba's production invocation traces
- Systems evaluated
  - **CPU-centric:** accelerators orchestrated by CPU cores
  - **RELIEF (HPCA'24):** accelerators orchestrated by a dedicated and centralized hardware manager
  - **Cohort (ASPLOS'23):** links pairs of accelerators that frequently go together, but otherwise relies on the cores to orchestrate the accelerators
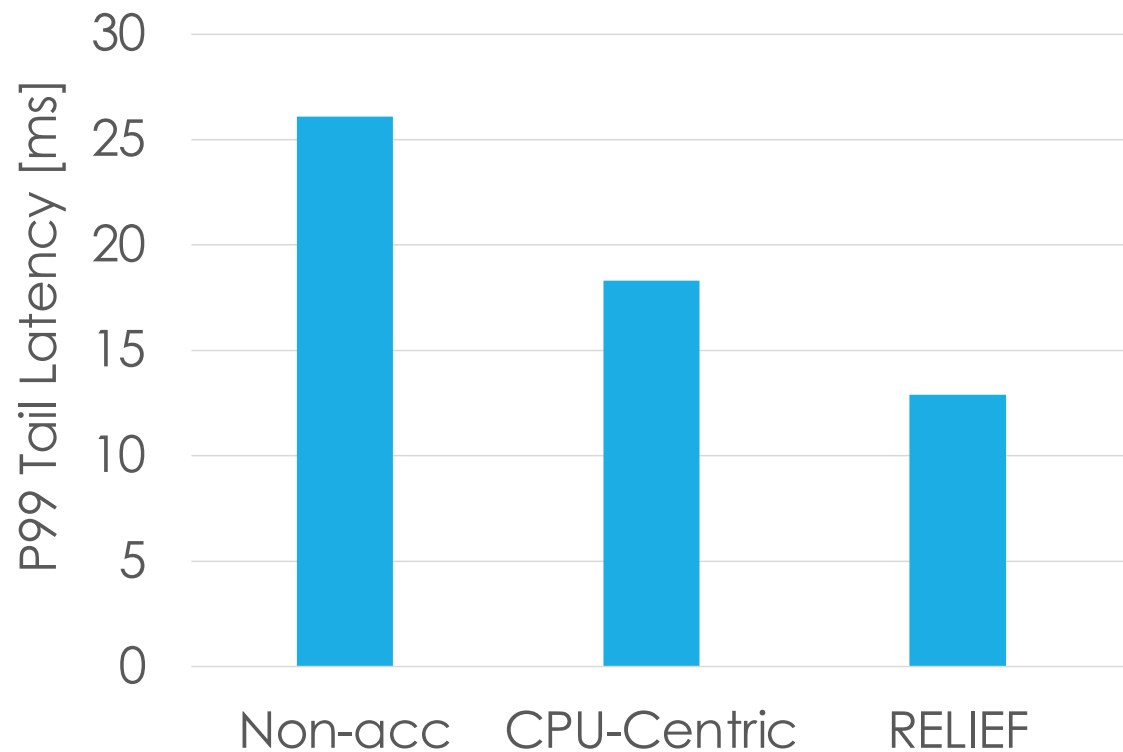  - **AccelFlow:** our proposal
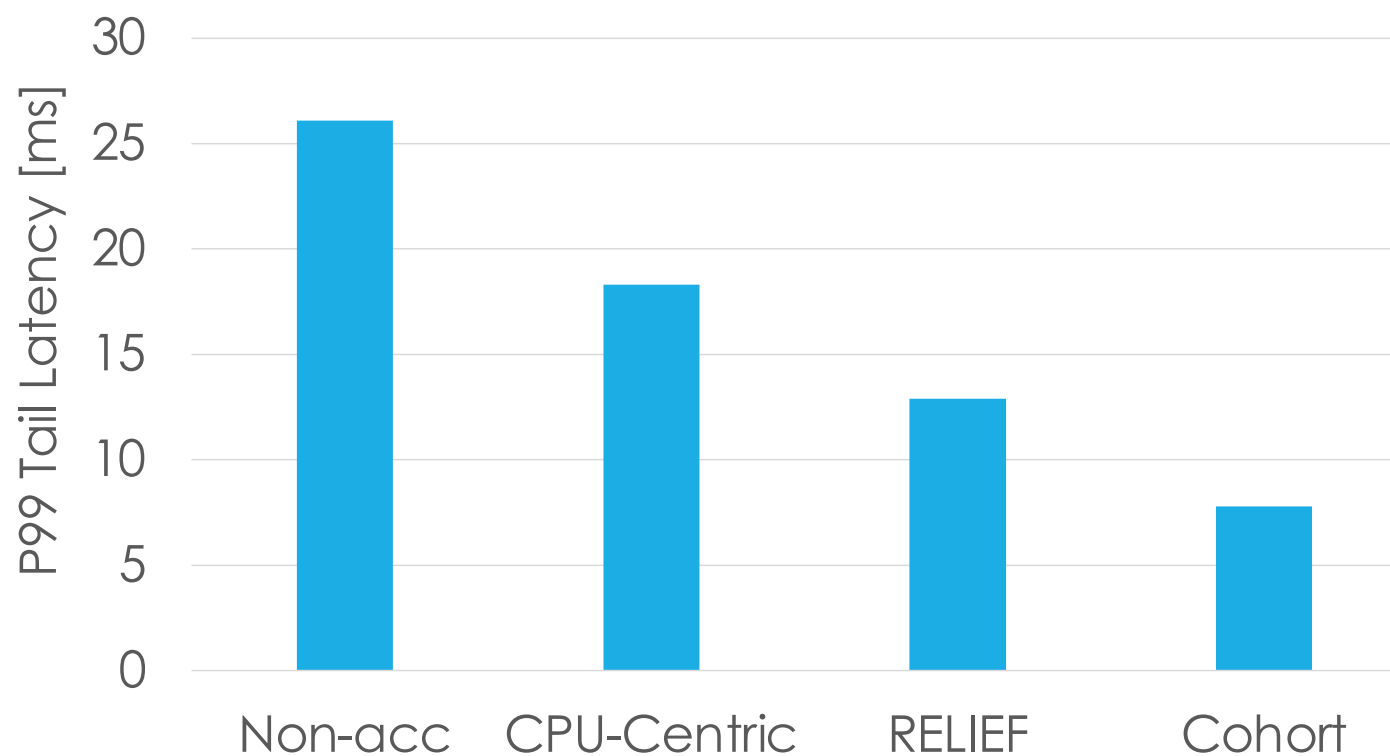
# AccelFlow Significantly Reduces Tail Latency



43

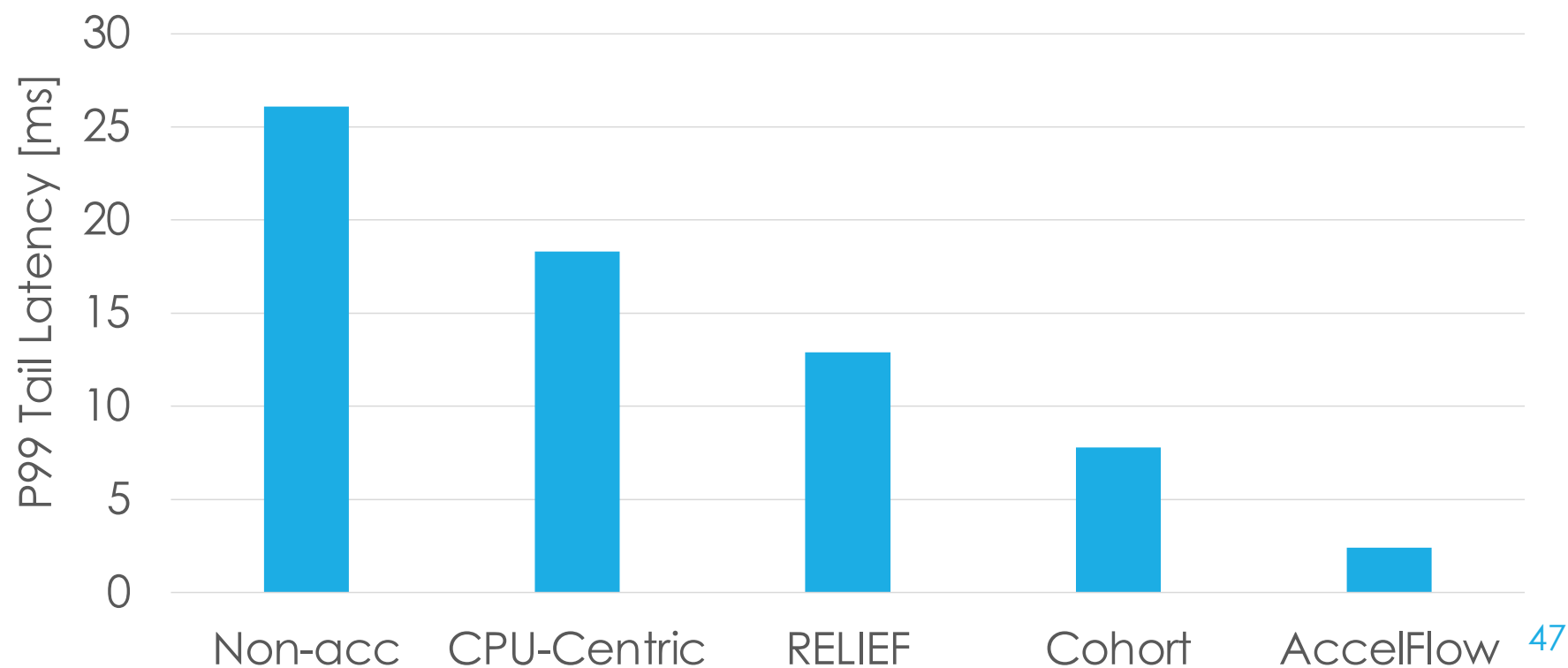# AccelFlow Significantly Reduces Tail Latency



44

# AccelFlow Significantly Reduces Tail Latency

# AccelFlow Significantly Reduces Tail Latency



46

# AccelFlow Significantly Reduces Tail Latency
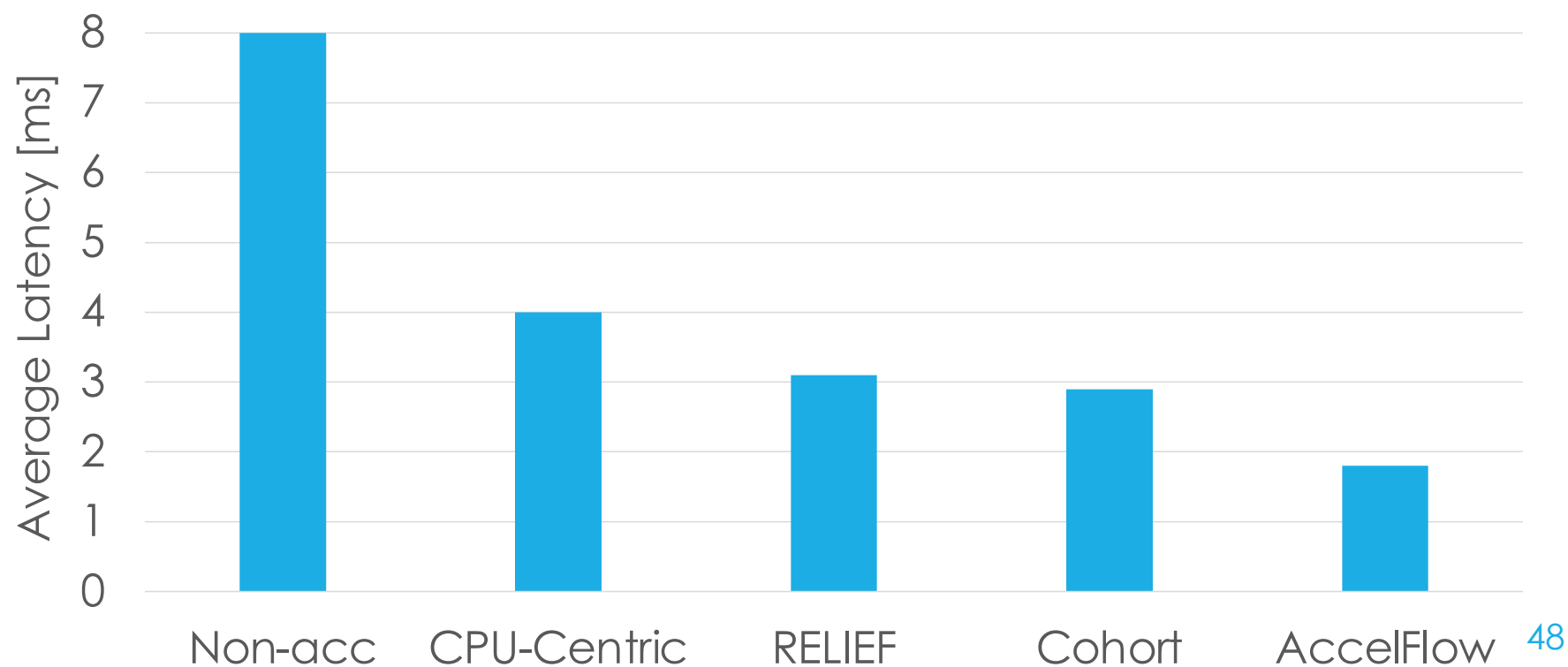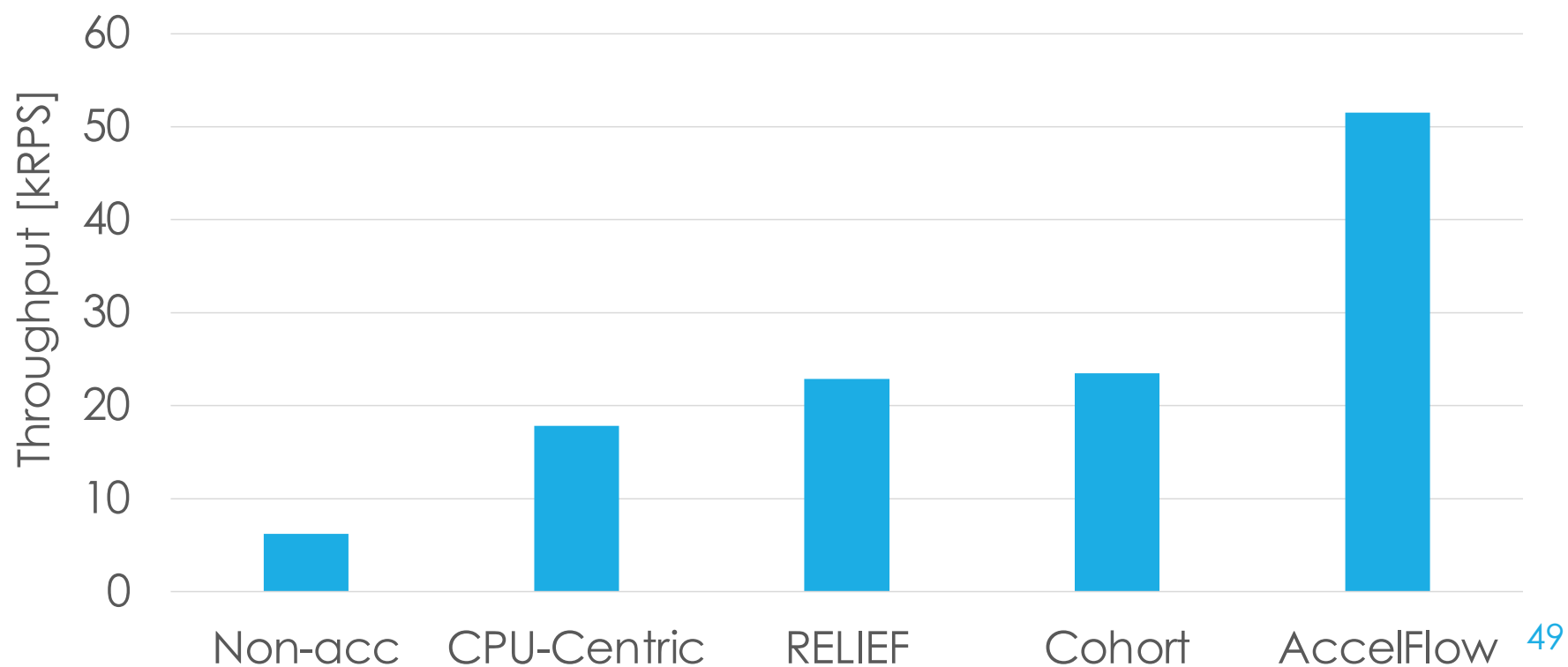


47

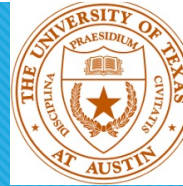# AccelFlow Reduces Average Latency



48

# AccelFlow Improves Throughput



49

# Conclusion

- An ensemble of domain-specific accelerators for "datacenter tax" has the potential to improve the efficiency of microservices

- Realizing these benefits requires an orchestration framework that can keep up with the fine-grained and dynamic microservices

- **AccelFlow:** the first accelerator-orchestration framework for on-chip accelerators targeting microservices

  - 70% lower tail latency

  - 38% lower average latency

  - 2.2x higher throughput

# AccelFlow: Orchestrating an On-Package Ensemble of Fine-Grained Accelerators for Microservices

## HPCA 2026, Sydney, Australia

Jovan Stojkovic*, Abraham Farrell, Zhangxiaowen Gong[†], Christopher J. Hughes[†], Josep Torrellas

University of Illinois at Urbana-Champaign, [†]Intel, *Joining UT Austin in Fall 2026