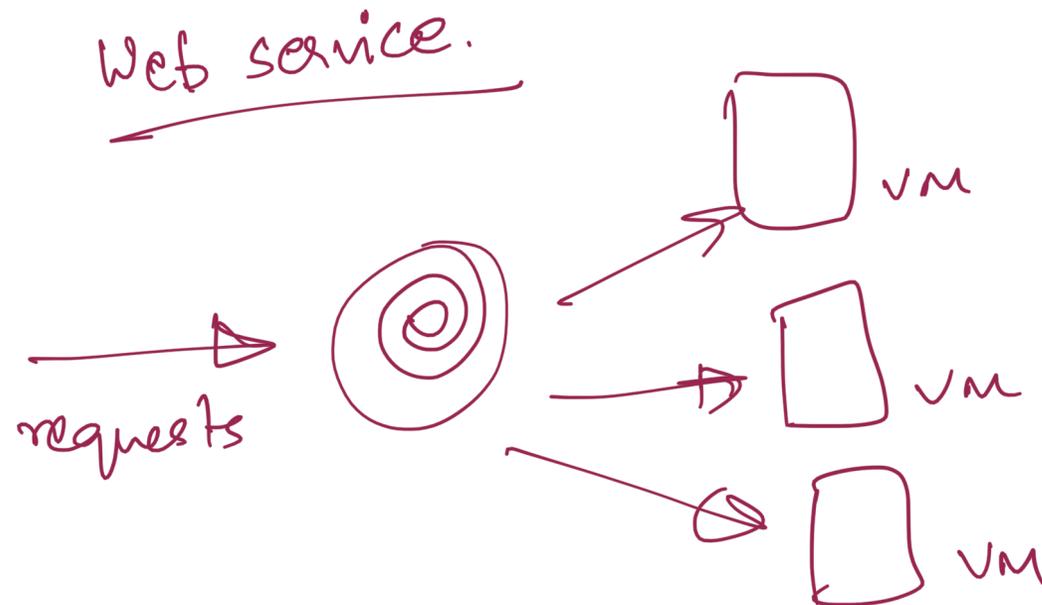




# *FaaS to FaaS* via Scalable Serverless Infrastructure



**Purushottam (Puru) Kulkarni**

Department of Computer Science and Engineering  
Indian Institute of Technology Bombay

# The Cloud Services Story



Infrastructure  
IaaS



Platform  
PaaS



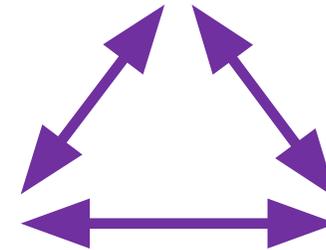
Software  
SaaS



Function  
FaaS

Resource access  
and virtualization

Ease of setup  
Performance  
Cost  
Security



Isolation  
Resource control  
Multiplexing

Software-defined  
techniques

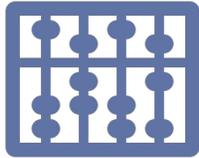
VMs/Containers, HCI,  
SDWAN, SDN, VNF,  
SAN, NAS, Data stores  
User space protocol stacks  
'''

Hardware assistance  
& capabilities

NVM, NVMeOF, RDMA,  
SRIOV, GPUDirect,  
Unified addressing  
SmartNICs, VPP  
...

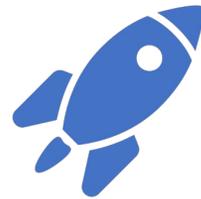
# The FaaS promise

Remove complexity of provisioning and management of cloud resources



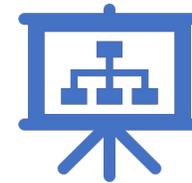
## Function granularity

1-1 mapping of service,  
provisioning and billing



## Serverless!

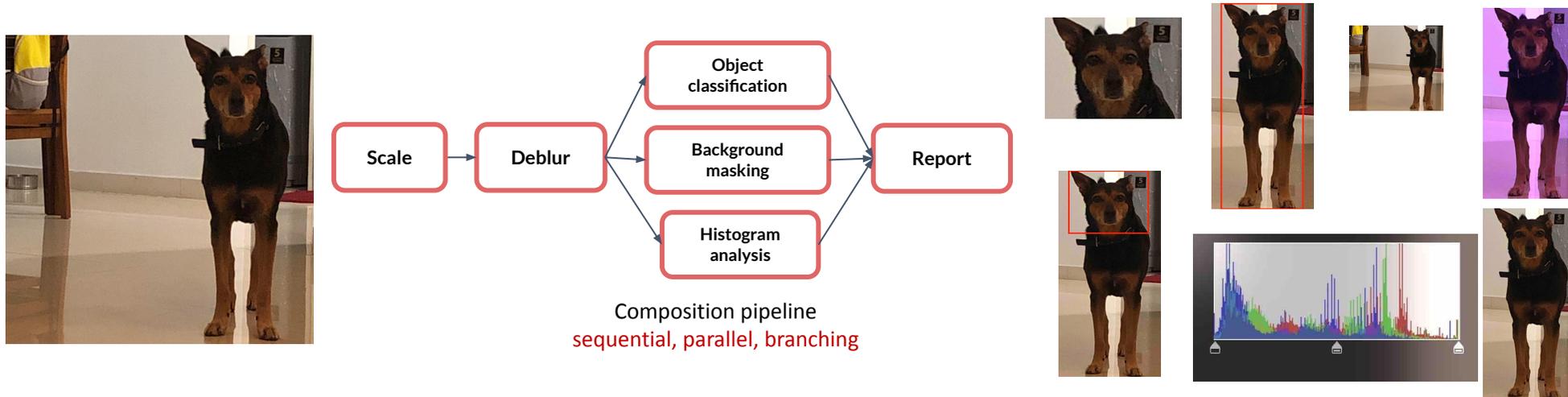
No infrastructure  
management (for clients)



## Agile development

Flexible and dynamic  
composition

# FaaS in action

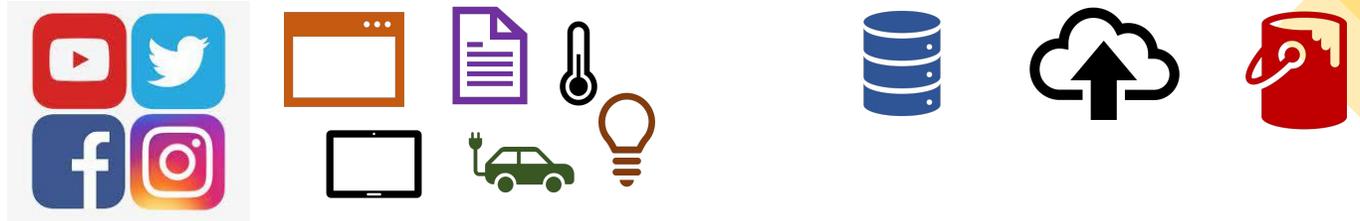


Applications	Functions
Social network	urlshorten followers userinfo postsStorage
Media service	userprofile ads thumbnail getrating getreviews recommender
E-commerce website	discounts catalogue castinfo wishlist payment balance contactinfo
Banking system	imagerecognition classification speed location luminosity
Video analytics	resize histanalysis encrypt decrypt
ML tasks	packetanalysis rulematching
Network functions	

# The FaaS game plan

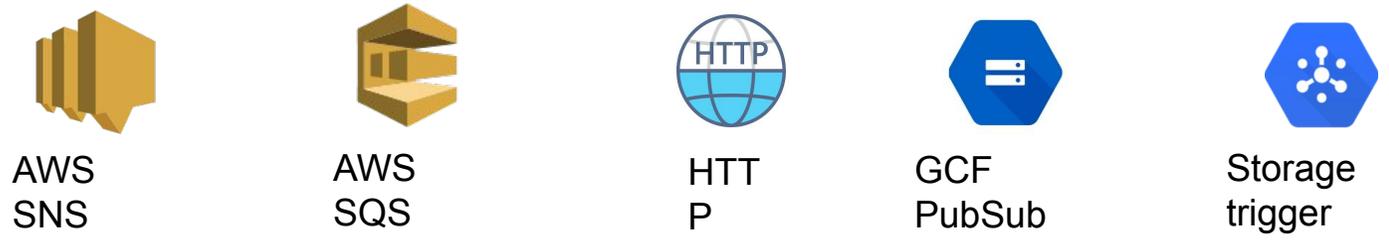
- **Decouple the what, the when and the where**

- Data generation & storage



- Triggers & events

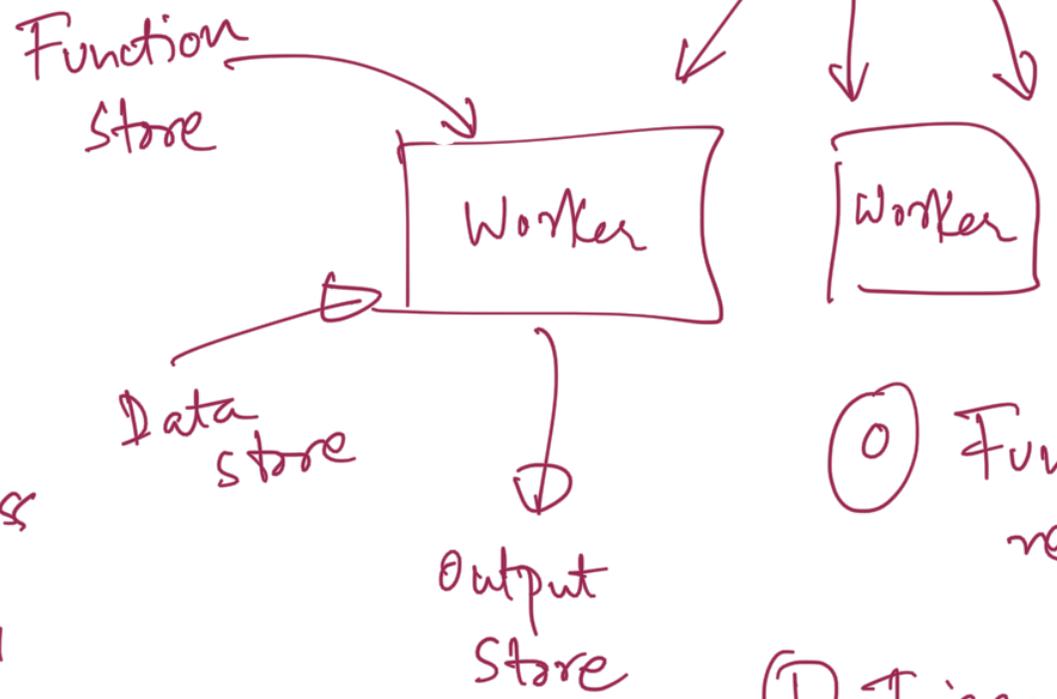
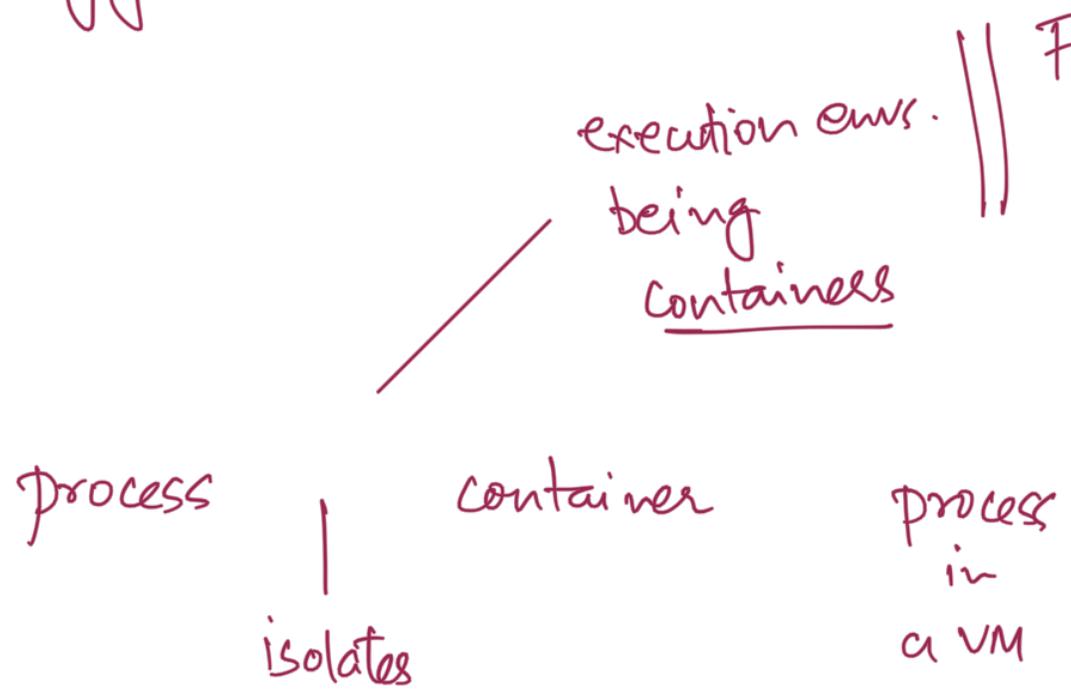
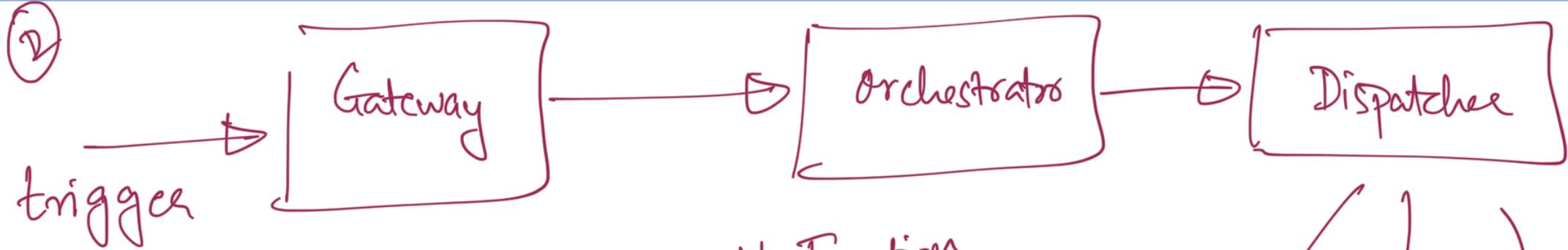
- 



- Execution platforms



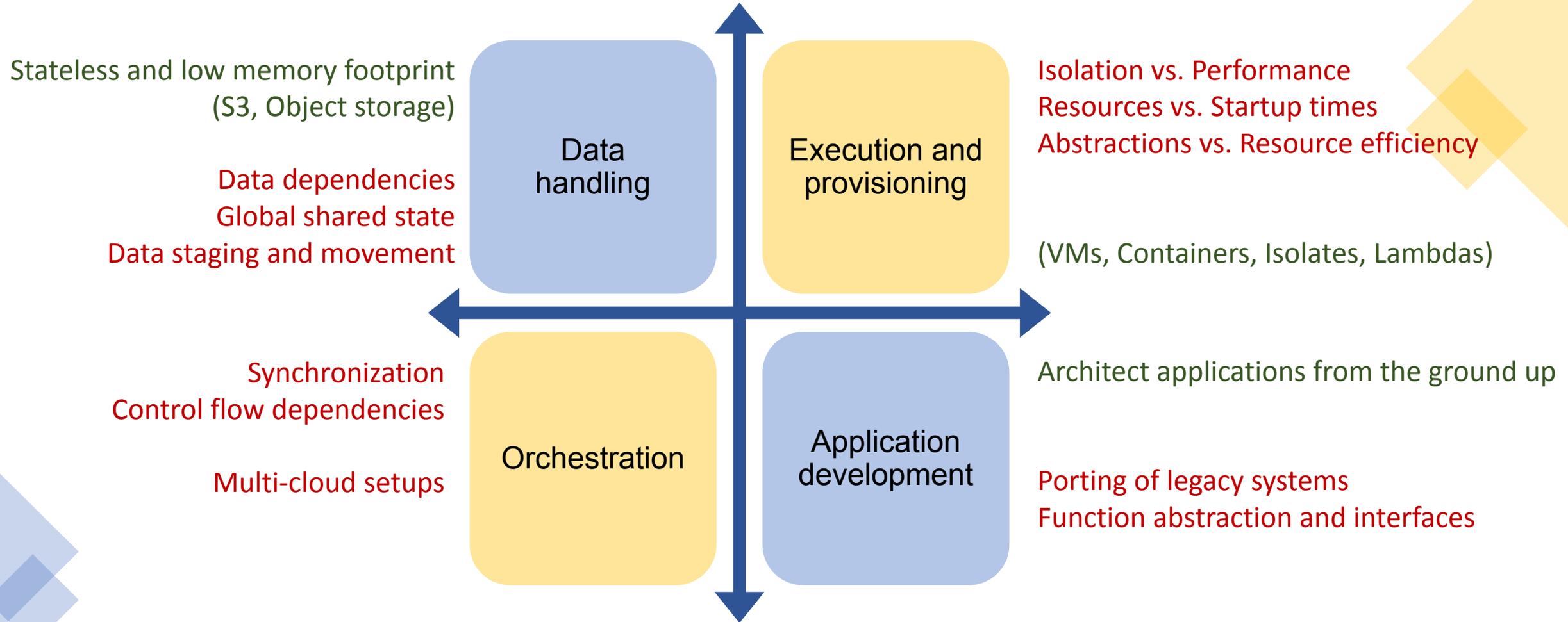
# Ex. serverless platform / arch.



① Function registration.

① Trigger registration + work (flow) specification.<sup>5</sup>

# Challenges

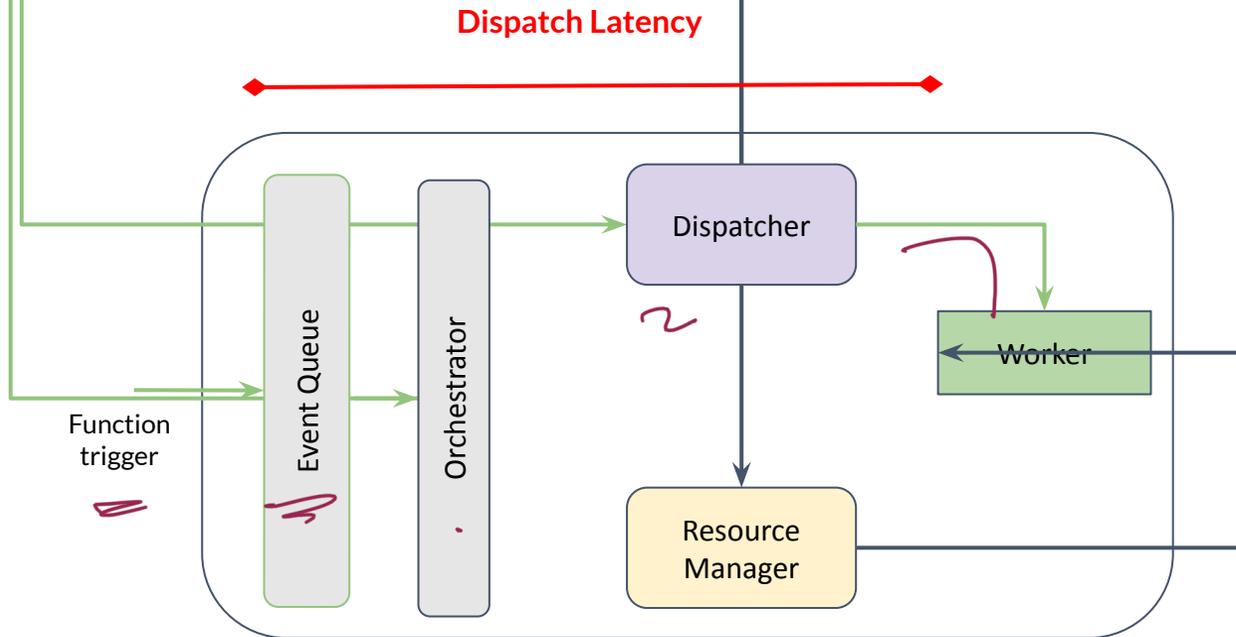


## Two example problems

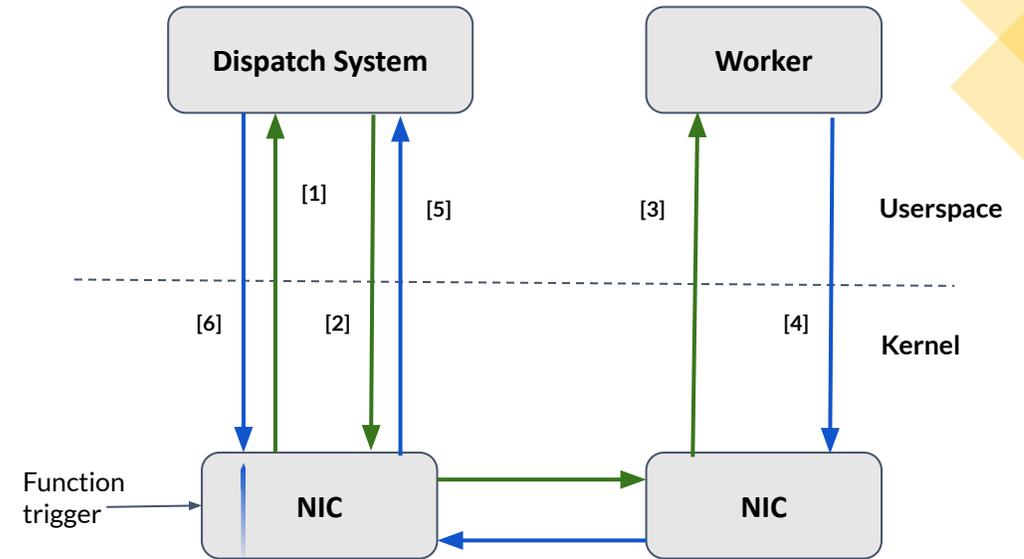
Scalable orchestration of serverless workflows

Acceleration Functions  
-as-a-Service

# Scalable orchestration of serverless workflows



Latency aggregates as FaaS components are executed for each function in workflow



~6 network stack traversals per trigger!



Speedo: Fast dispatch and orchestration of serverless workflows

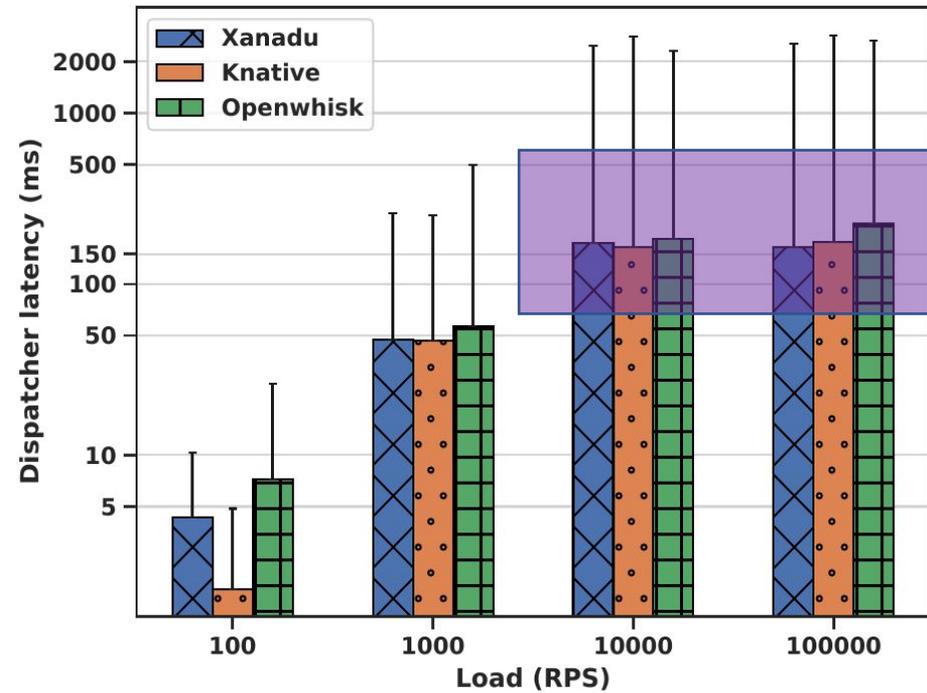
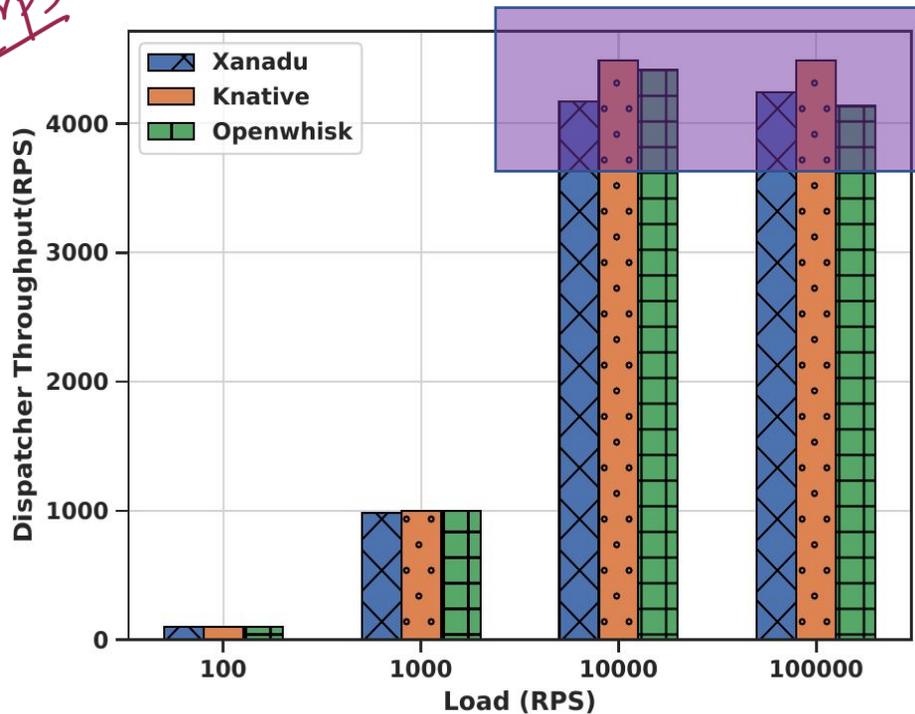
Nilanjan Daw, Umesh Bellur, Purushottam Kulkarni  
ACM Symposium on Cloud Computing (SoCC '21), November 1-4, 2021

# The dispatch bottleneck

- Social network workflow: **~6K posts** per second => **~70K triggers** per second

4000 rps

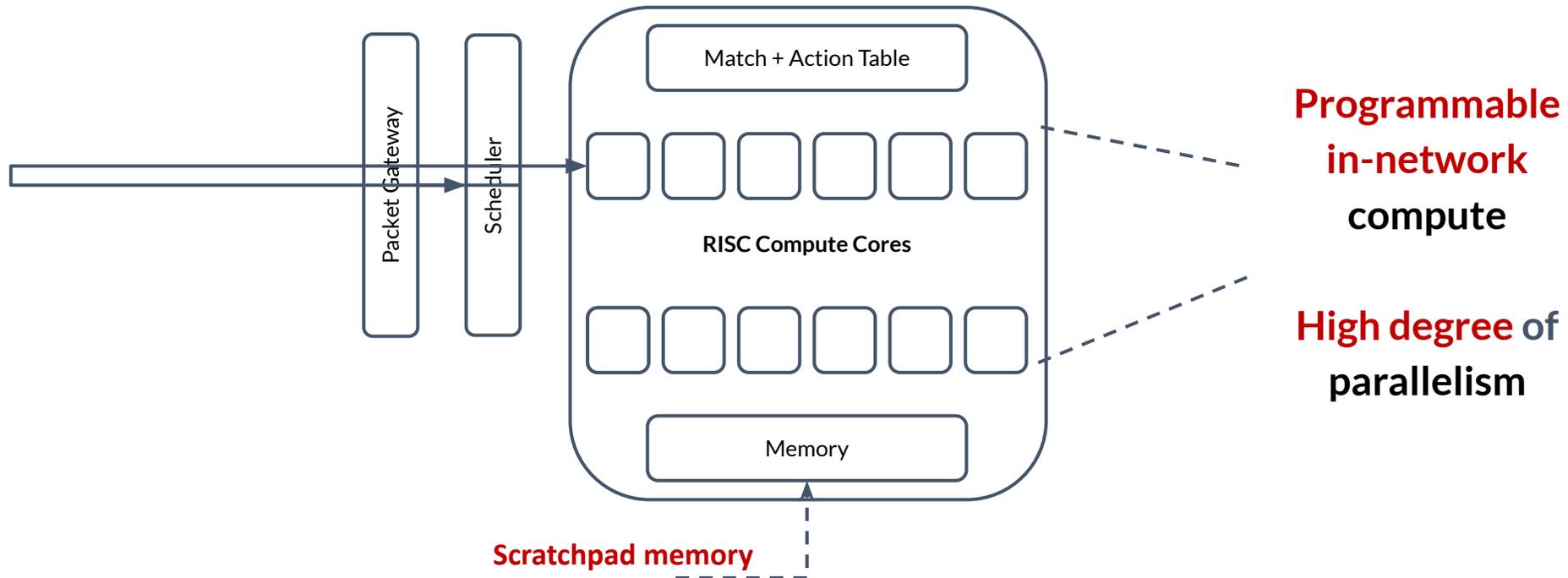
A



Userspace dispatchers saturate at moderate loads (with high dispatch latency)

# Compute at the Edge!

## Modern networking hardware (SmartNICs)



Offload FaaS dispatch system onto SmartNICs for decisions **closer** to the network !

# Speedo key ideas

Exploit **in-network compute** and **programmability** for FaaS dispatch

**#1**

**Network stack traversal  
is costly**

**Observation**

Dispatch decisions closer to the wire improves latency

**Solution**

Offload the dispatcher to the NIC

**#2**

**Userspace parallelism is  
non-deterministic and  
costly**

**Observation**

RISC cores on SmartNICs provide scalability

**Solution**

Exploit hardware parallelism on SmartNICs for scalability

**#3**

**Dispatcher is not  
stateless**

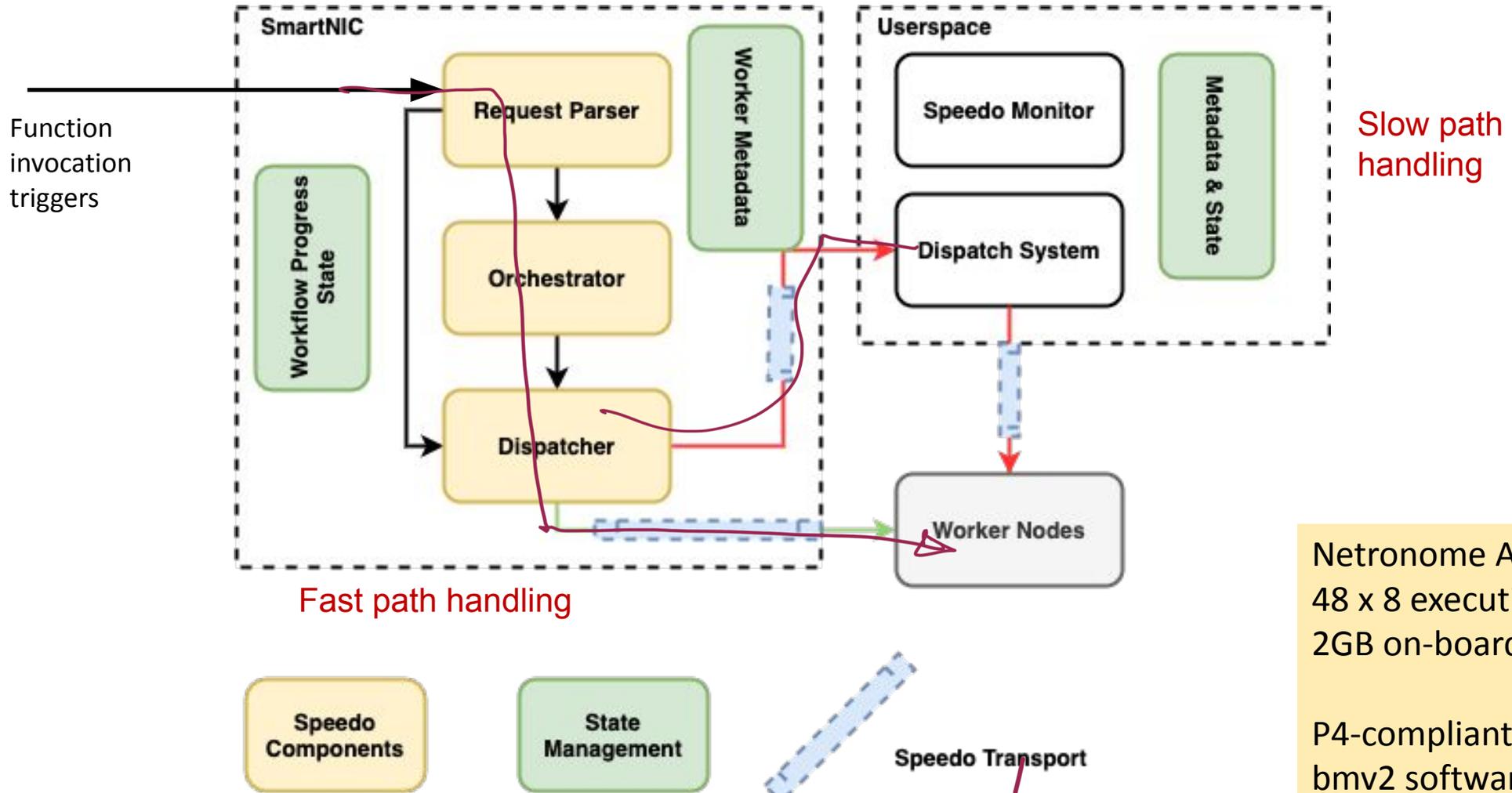
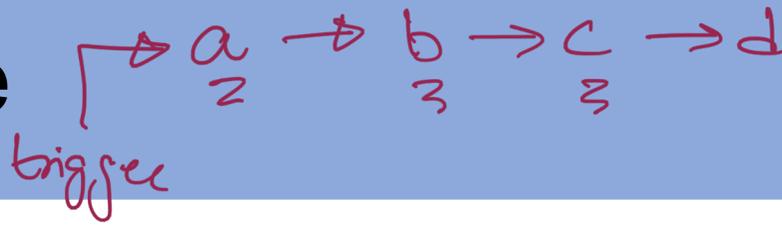
**Observation**

Lifetime of workflow state is ephemeral and short lived

**Solution**

Use NIC onboard memory to store workflow state

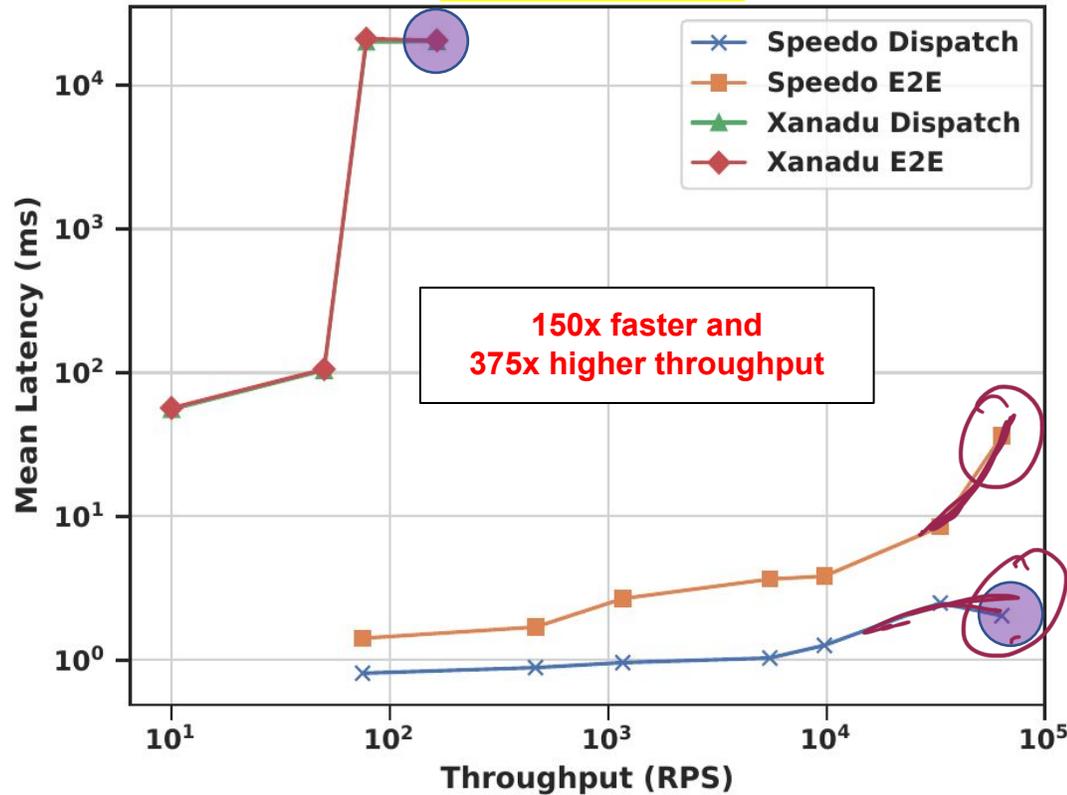
# Speedo Architecture



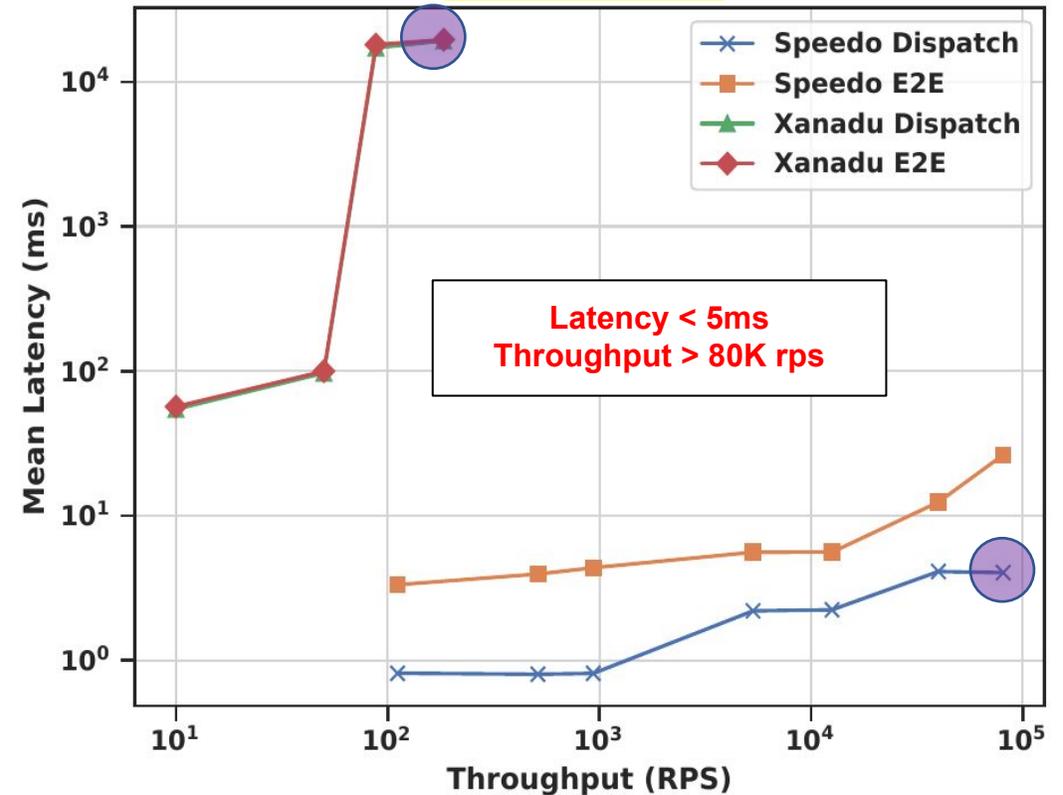
Netronome Agilio NIC  
 48 x 8 execution instances  
 2GB on-board memory  
 P4-compliant  
 bmv2 software switch

# Speedo Performance

## Social Network



## Media service

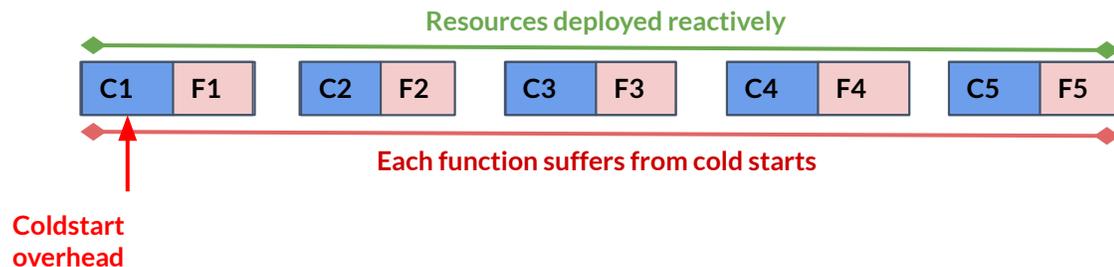


- Speedo **saturates server** worker capacity without saturating dispatch system

# Speedo Summary

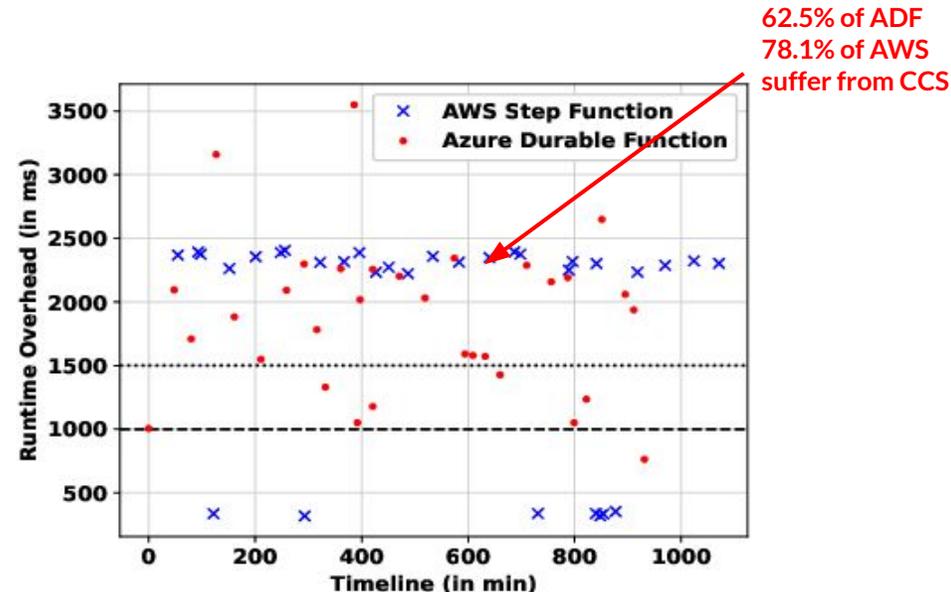
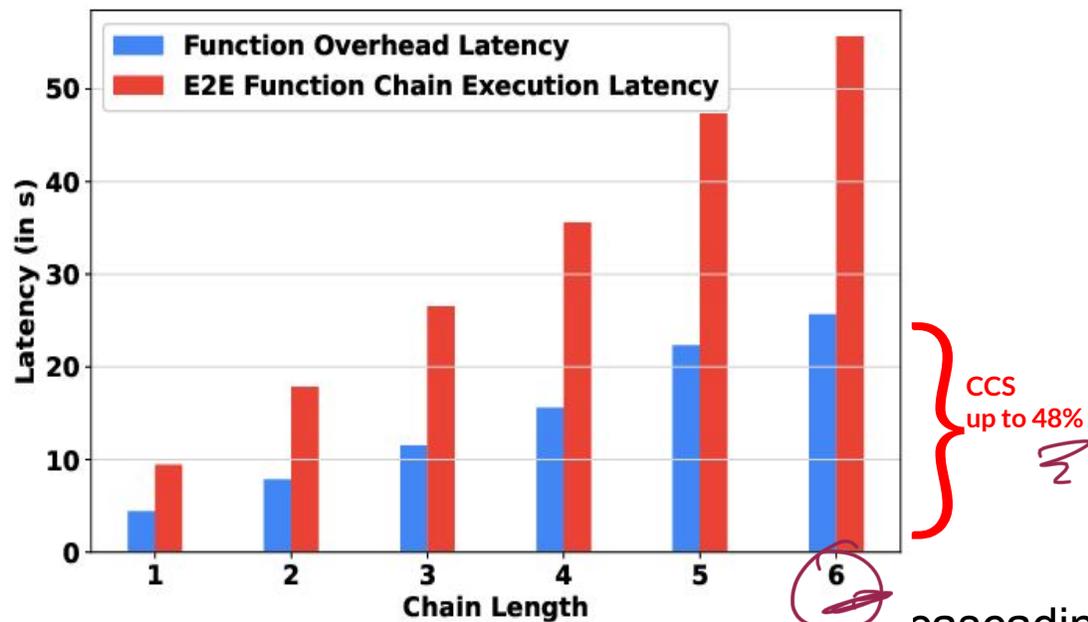
- Speedo exploited compute and programmability of SmartNICs to improve FaaS management components --- dispatcher and orchestrator
  - 177x** improvement in dispatch throughput (compared to userspace dispatch)
  - 50x** average reduction in dispatch latency
- **Future work**
  - Dynamic offloading of compute between host CPU and SmartNIC
  - SmartNIC resource partitioning across applications

# Addressing the Cascading Cold start (CCS) problem



## E2E latency breakup

Cold start overhead C1+C2+C3+C4+C5	Function execution times
---------------------------------------	--------------------------



Addressing cascading cold starts in serverless function chain deployments

Nilanjan Daw, Umesh Bellur, Purushottam Kulkarni

21<sup>st</sup> ACM/IFIP International Middleware Conference, December 2020

# Xanadu key ideas

**#1**

**Cascading cold starts cause performance degradation**

**Observation**

Pre-deploying resources can prevent cascading overheads

**Solution**

Speculatively pre-deploy resources

**#2**

**Cost of speculative deployment still significant**

**Observation**

Application runtime tends to remain stable

**Solution**

Deploy resources Just-in-time to reduce resource costs

**#3**

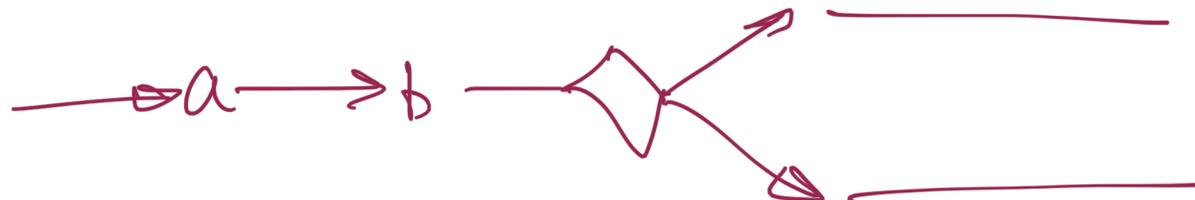
**Naive pre-deployment is resource expensive**

**Observation**

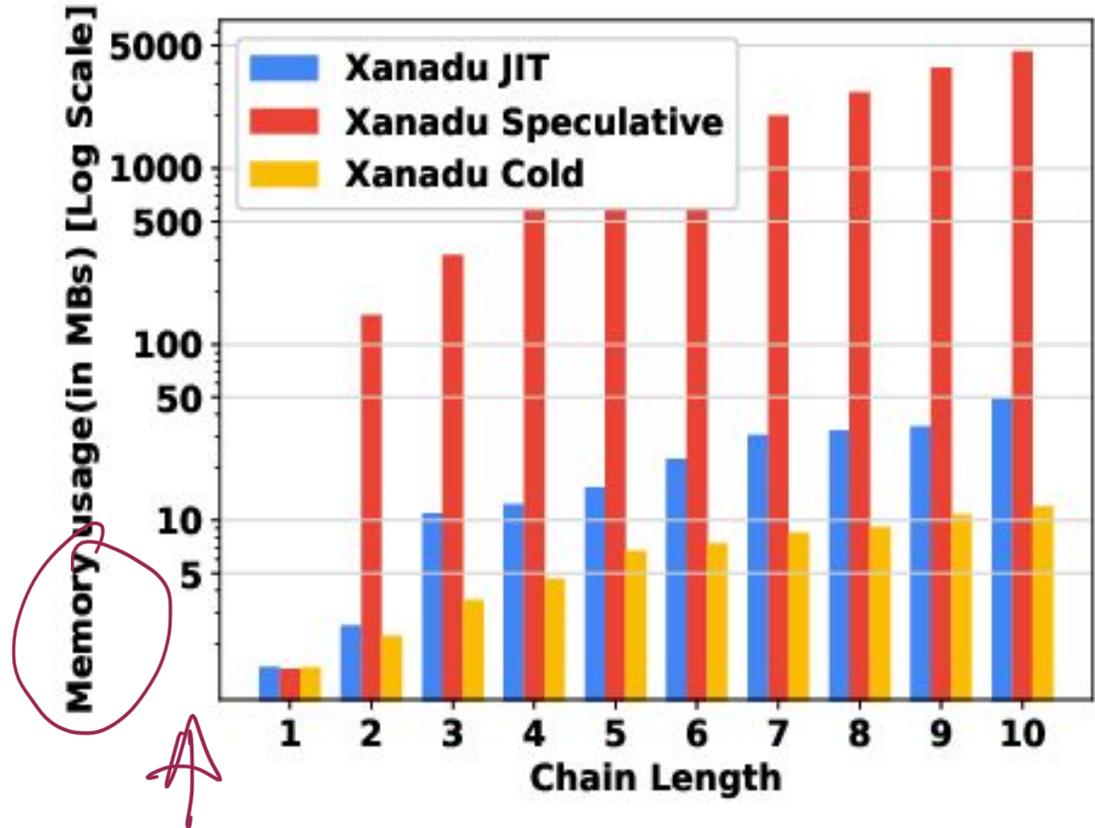
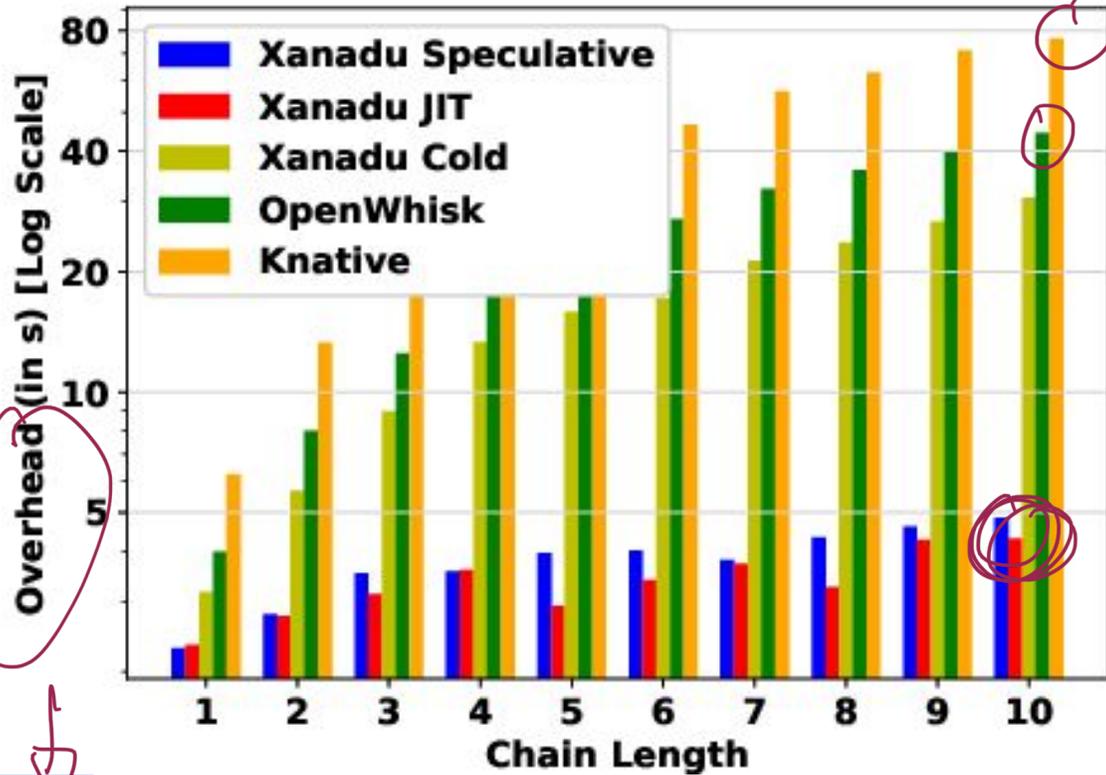
Application behavior tends to follow a stable path

**Solution**

Detect the Most Likely Path (MLP) for speculative deployment



# JIT + MLE advantages



# Xanadu Summary

- Cold start latency mitigation intersects several dimensions
  - Execution platforms (VMs, containers, isolates, lamdas)
  - Resources reservation vs. Costs
  - Proactive vs. Reactive management
  - **Data provisioning (on the network)**
    - Standardized data pipelines need of the hour!

## Two example problems

Scalable orchestration of serverless workflows

Acceleration Functions  
-as-a-Service

# The acceleration spectrum

Applications

Graphics, Signal processing, Networking, Crypto, AI, ML, Simulations

API &  
Libraries

CUDA | OPENCL | MPI | OpenACC | TensorFlow | PyTorch | NumPy | P4

Access +  
Virtualization

Direct | Passthrough | GPUDirect  
Para/ Full / Hardware Virtualization | API remoting

Hardware  
accelerators

GPU

TPU

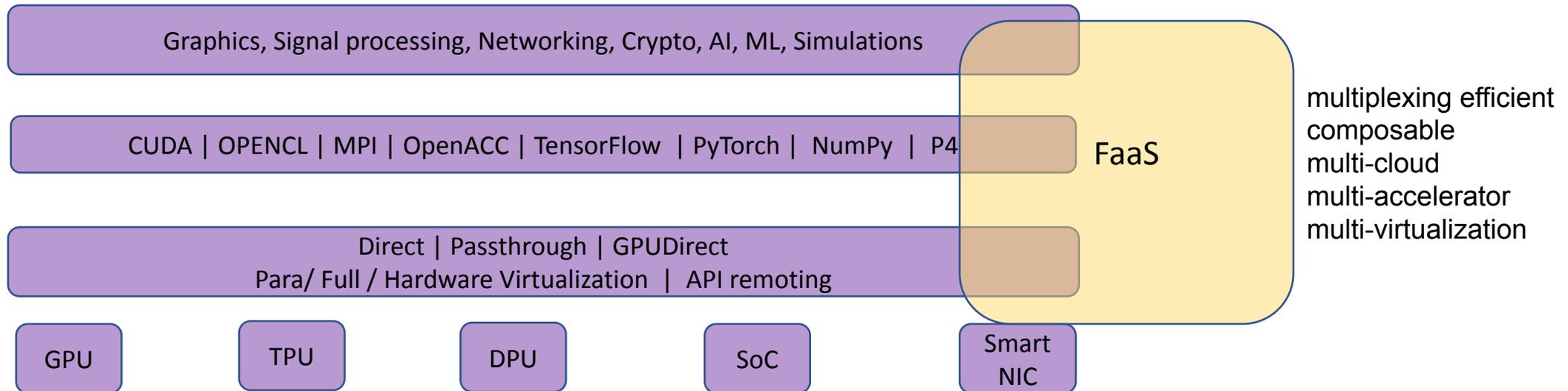
DPU

SoC

SmartNIC

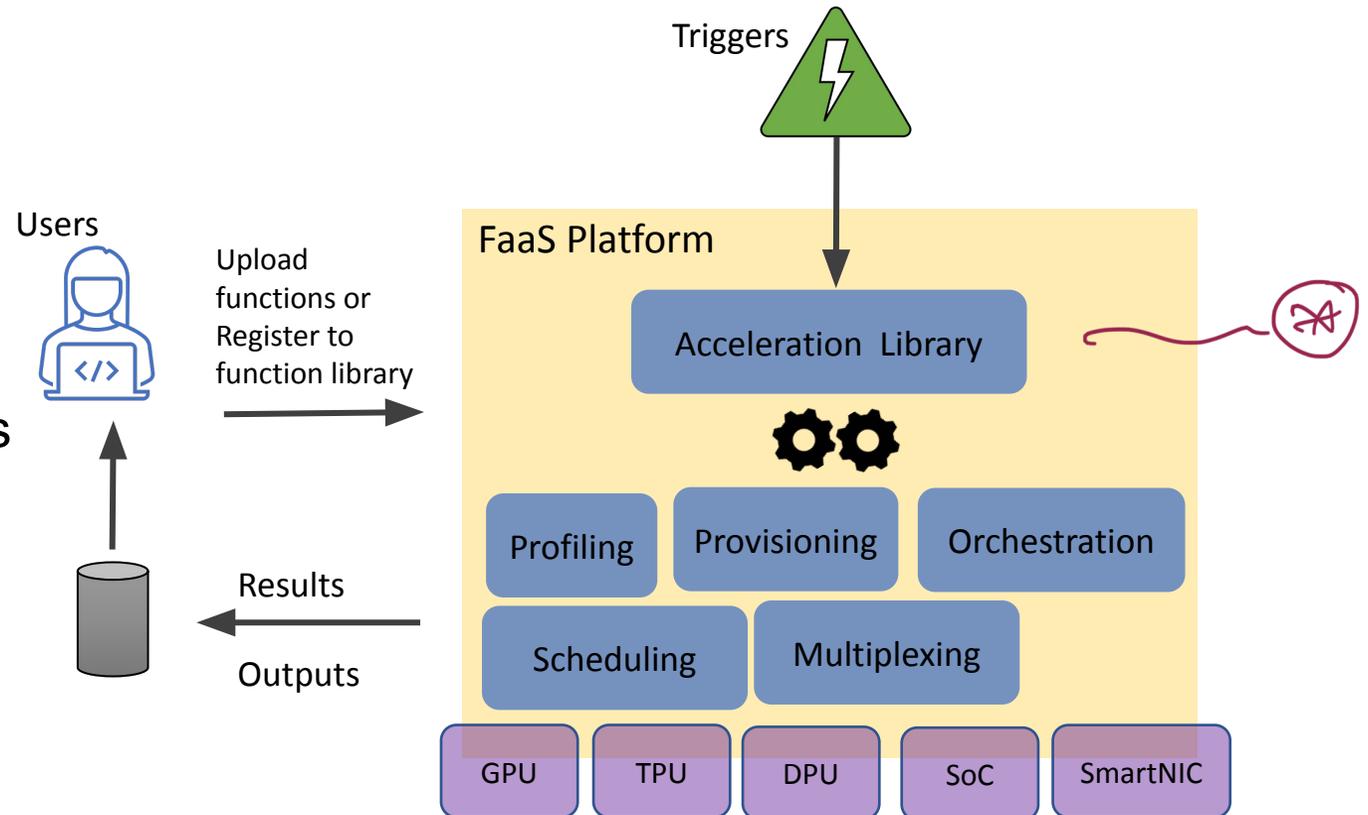
# Challenges with Accelerators

- Exposure as a resource, needs programmatic control and usage
- Reserve and use models are not efficient
- Preemption and multiplexing not first-class primitives
- Accelerator interfaces not generic and lead to accelerator stack lock-in



# Accelerated Functions-as-a-Service

- The **function** abstraction
  - Serverless, composition, ...
- Multi-accelerator backends
- Transparent elastic resource pools
- Suite of optimizations for improved resource usage efficiency



## FaaSter: Accelerated Functions-as-a-Service on Heterogeneous GPUs

# AFaaS with Heterogenous GPUs

## • Issues

- GPU functions run to completion
- Very limited/no control on scheduling
- Multi-GPU parallelism per task missing

## • Key enabler

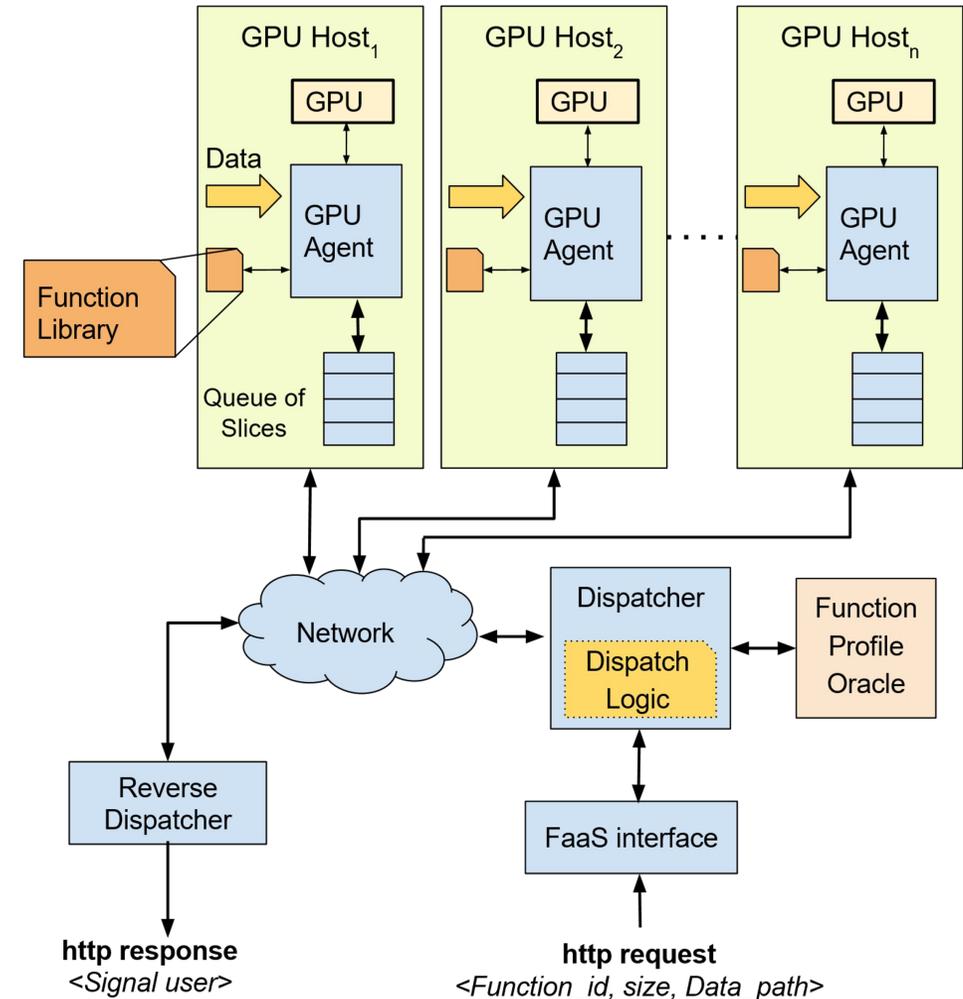
- GPU kernel slicing

## • Solution components

- #slices and #GPUs for each function
- which types of GPUs, and #GPUs per type?
- online heterogeneous GPU scheduling

### Goals:

Matching task types with appropriate accelerators  
Minimize average turn-around time of functions  
Priority based allocation of acceleration services



# Does it all add up?

Slicing-based and accelerator appropriate mapping improved resource efficiency and performance

**Better than or equal to No slicing approach under various conditions**

