

XFaaS

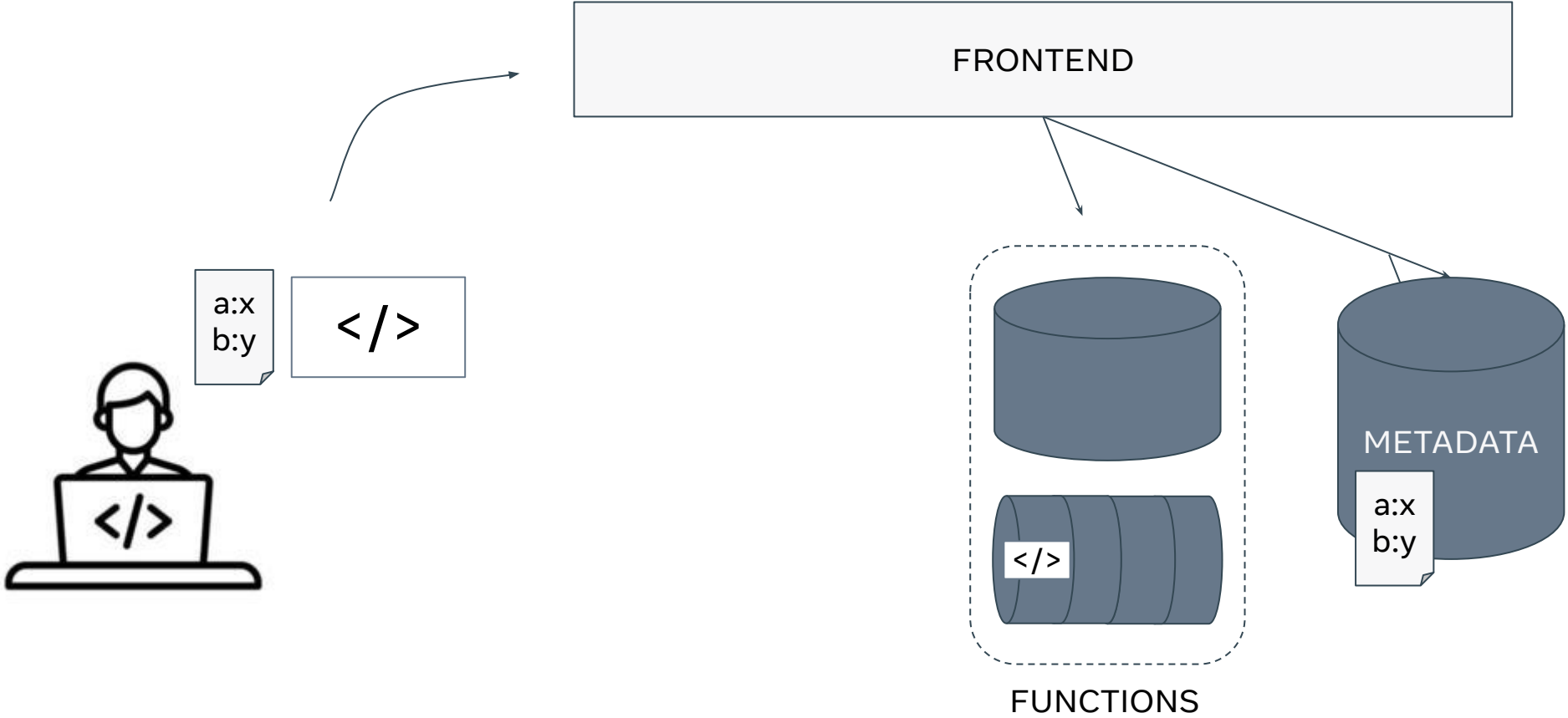
HYPERSCALE AND LOW COST SERVERLESS FUNCTIONS AT META

Alireza Sahraei ₁	Soteris Demetriou ₂	Amirali Sobhgol ₁	Haoran Zhang ₃	Abhigna Nagaraja ₁	Neeraj Pathak ₁	Girish Joshi ₁	Carla Souza ₁	Bo Huang ₁	Wyatt Cook ₁
Andrii Golovei ₁	Pradeep Venkat ₁	Andrew McFague ₁	Dimitrios Skarlatos ₄	Vipul Patel ₁	Ravinder Thind ₁	Ernesto Gonzalez ₁	Yun Jin ₁	Chunqiang Tang ₁	

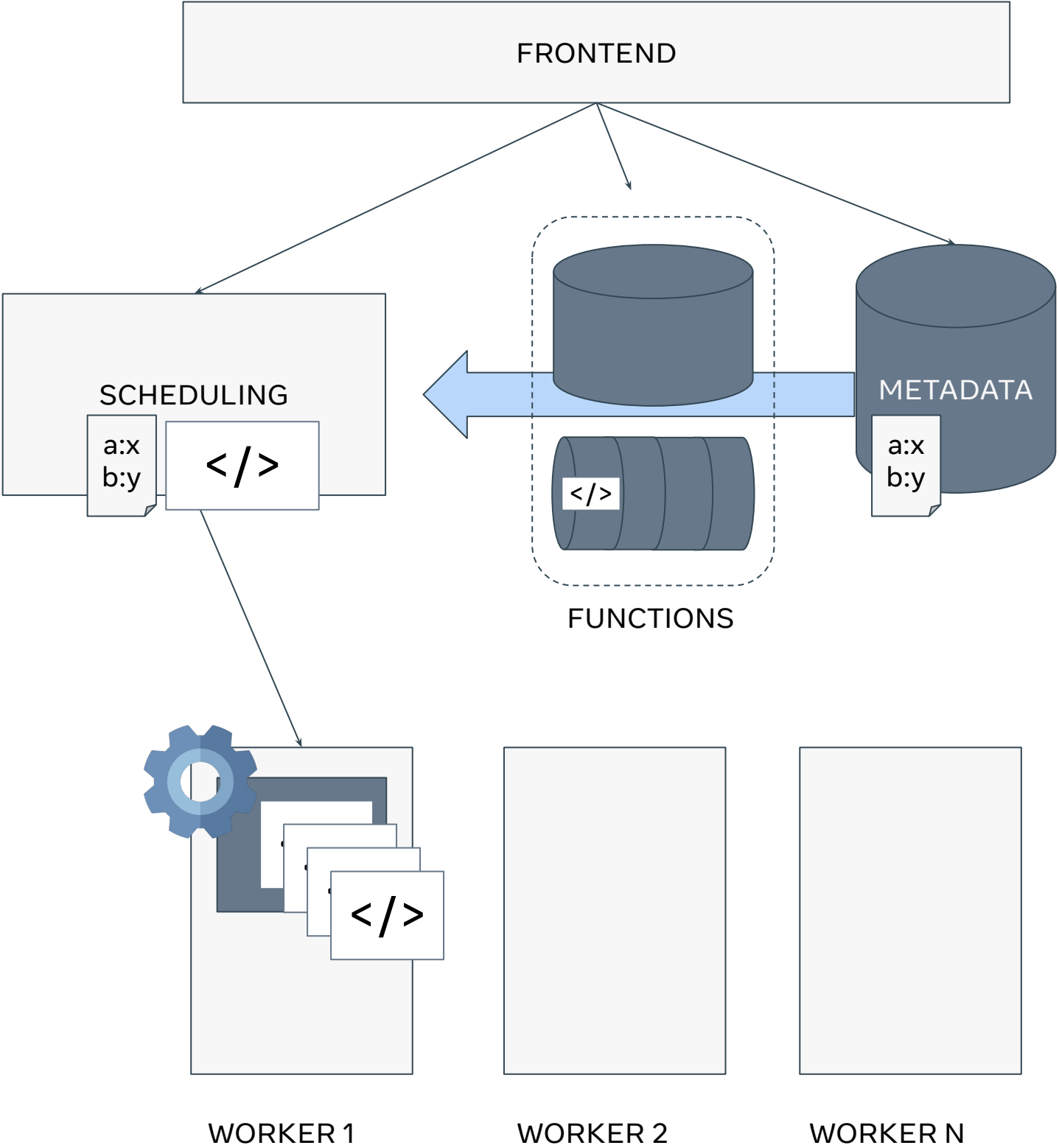


01 BACKGROUND & MOTIVATION

FUNCTION AS A SERVICE



FUNCTION AS A SERVICE



FUNCTION AS A SERVICE

PUBLIC



Brooker, Marc, et al. "On-demand Container Loading in {AWS} Lambda." 2023 USENIX Annual Technical Conference (USENIX ATC 23). 2023

Agache, Alexandru, et al. "Firecracker: Lightweight virtualization for serverless applications." 17th USENIX symposium on networked systems design and implementation (NSDI 20). 2020



Azure Functions

Shahrad et al. Serverless in the wild: Characterizing and optimizing the serverless workload at a large cloud provider. In USENIX Annual Technical Conference, 2020.



Alibaba Cloud

Wang, Ao, et al. "{FaaSNet}: Scalable and fast provisioning of custom serverless container runtimes at alibaba cloud function compute." 2021 USENIX Annual Technical Conference (USENIX ATC 21). 2021



Google Cloud Functions

PRIVATE



This work...




FUNCTION AS A SERVICE AT META

- highly heterogeneous workloads

Workload	Trigger	Calls/ second	CPU (MIPS)	Execution Time (s)	Memory (MB)
Notifications	Data Warehouse	3.4M	65-200	0.55 - 1.1	10 - 90
Morphing Framework	Queue	25K	1.5M - 27M	65 - 155	30 - 230

WHAT ABOUT HARDWARE COSTS?

“81% of the applications are invoked once per minute or less on average. This suggests that the cost of keeping these applications warm, relative to their total execution (billable) time, can be prohibitively high”

 Shahrad et al. Serverless in the wild: Characterizing and optimizing the serverless workload at a large cloud provider. In USENIX Annual Technical Conference, 2020.

02 CHALLENGES

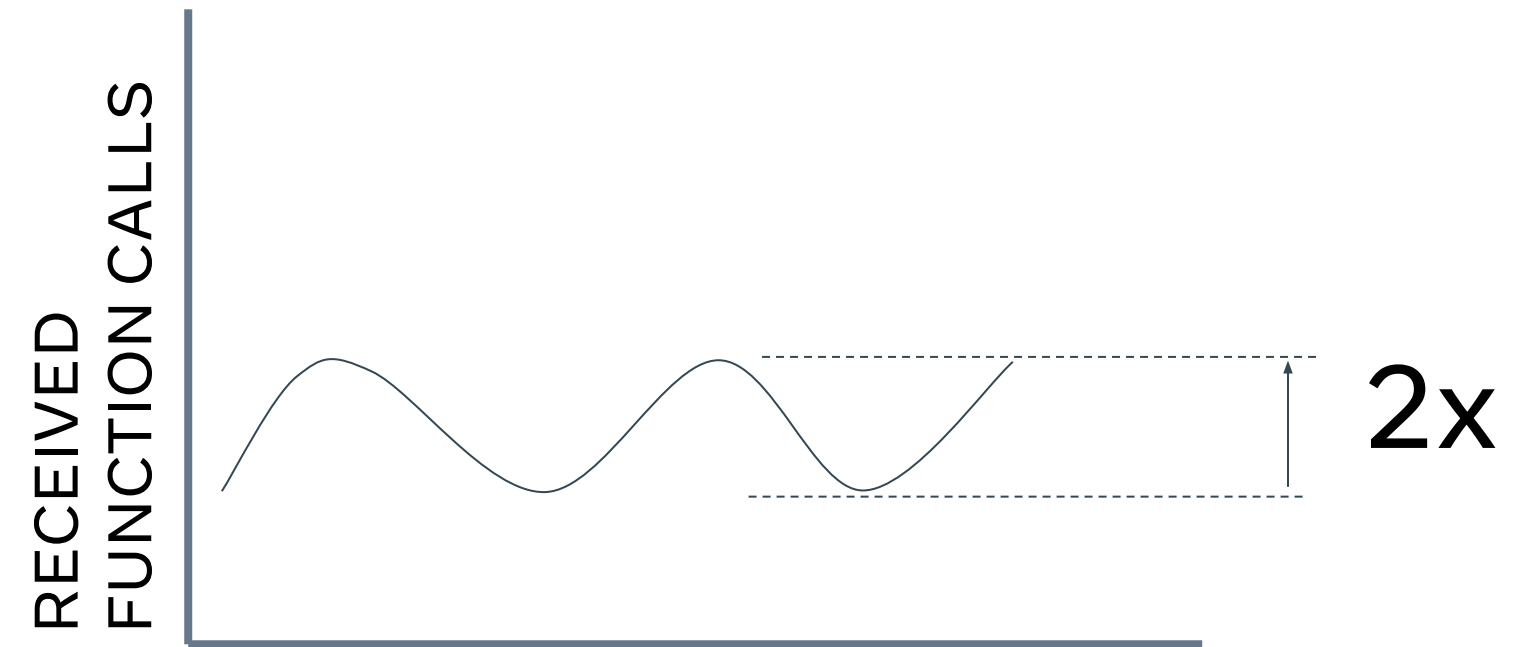
- A Lengthy Cold Start **[NOT COVERED IN THIS TALK]**
- B High Variance of Load
- C Downstream Overloads

High Variance of Load

Problem

1. Previous work reported a high peak-to-trough ratio of function calls
2. At Meta, the ratio can be as high as 4.3

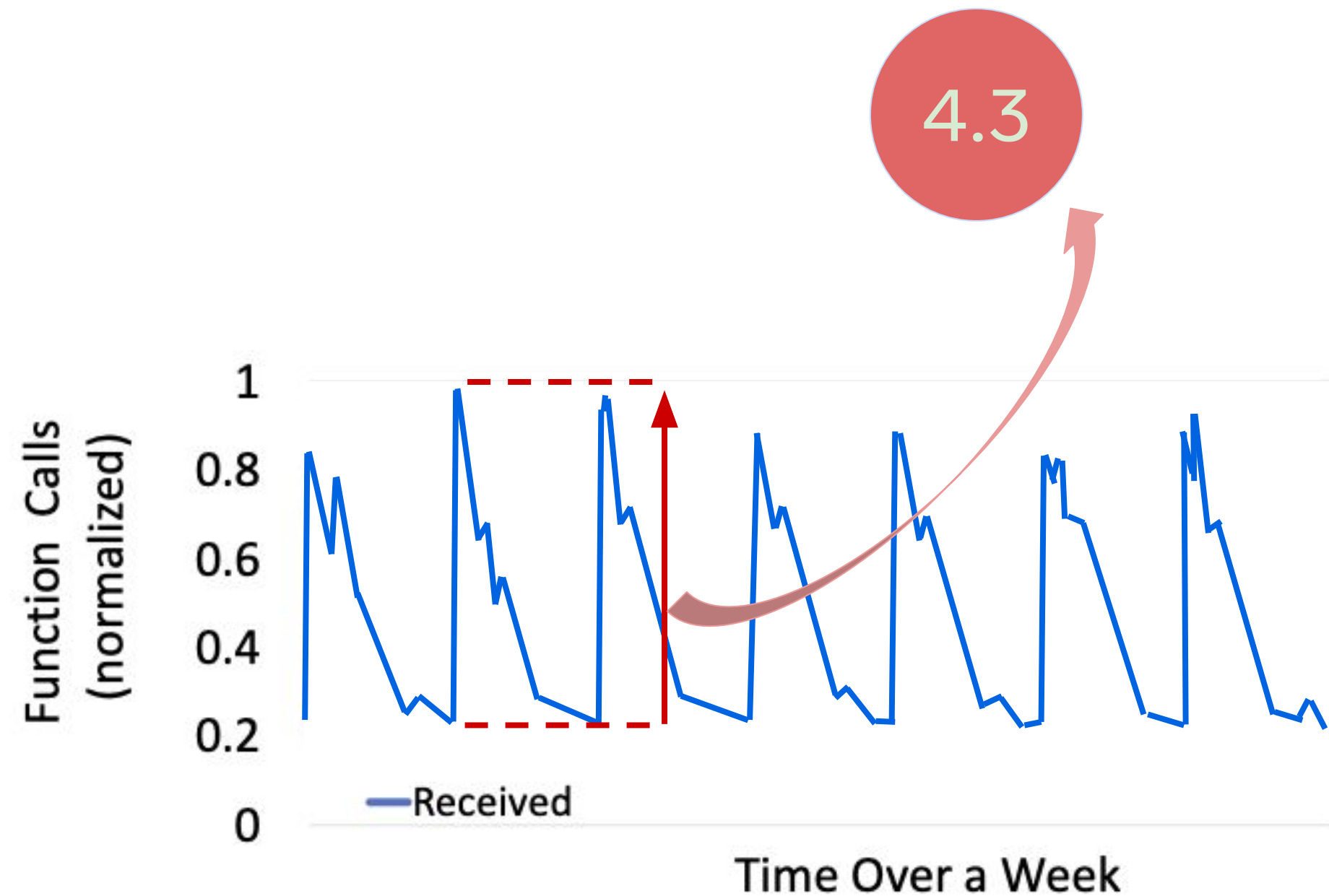
Shahrad et al. Serverless in the wild: Characterizing and optimizing the serverless workload at a large cloud provider. In USENIX Annual Technical Conference (USENIX ATC 20). 2020.



High Variance of Load

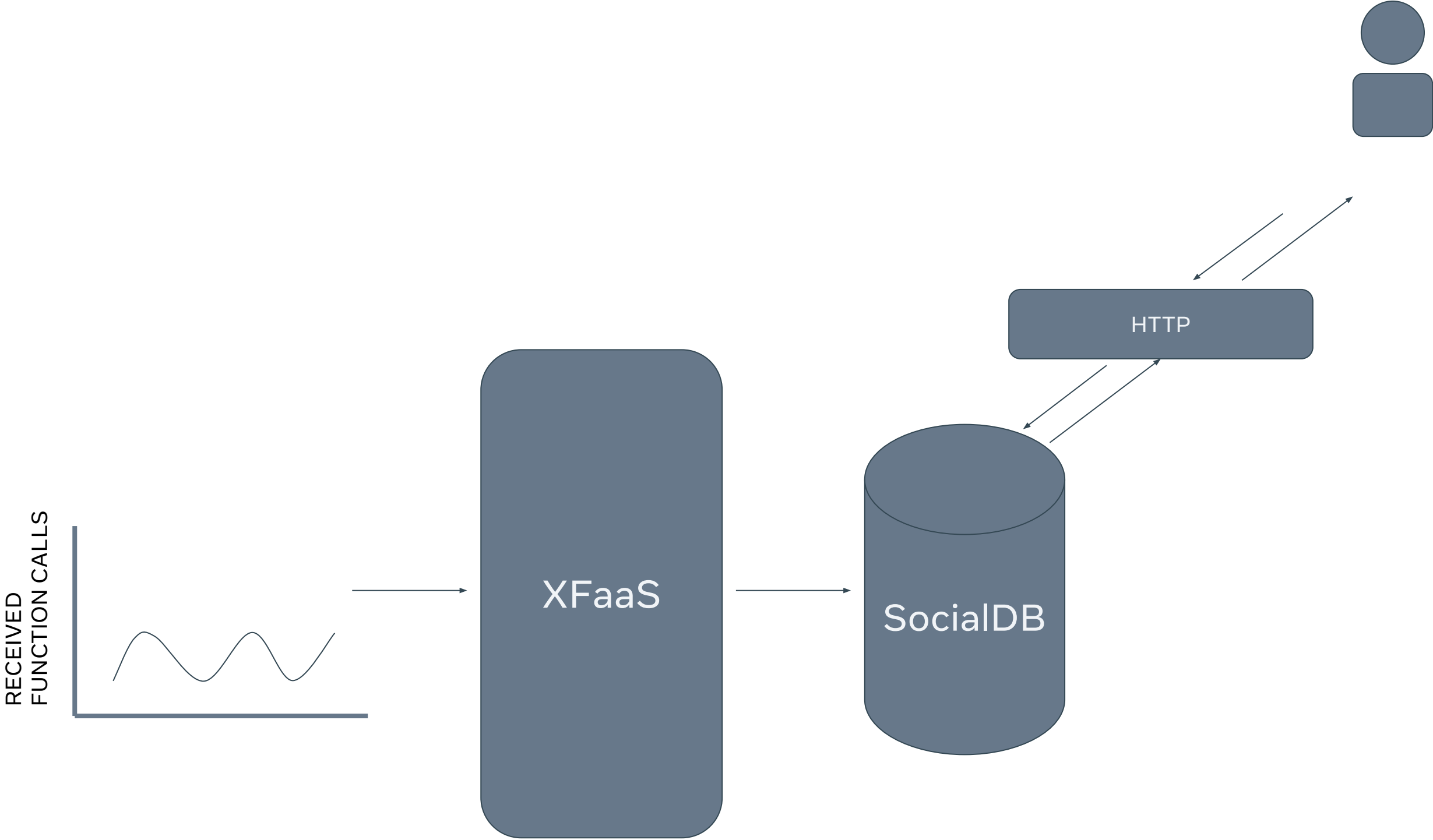
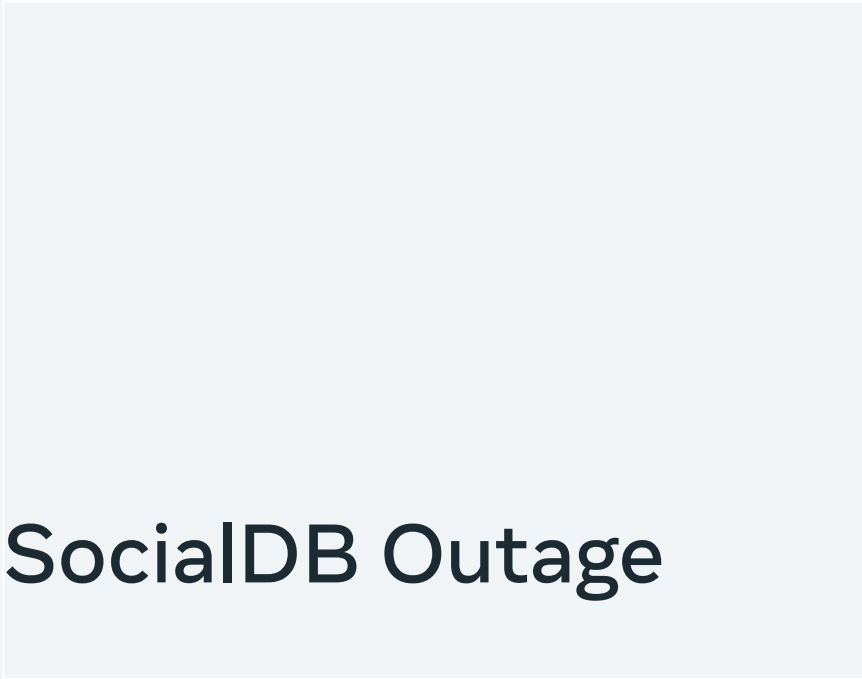
Problem

1. Previous work reported a high peak-to-trough ratio of function calls
2. At Meta, the ratio can be as high as 4.3



DOWNSTREAM OVERLOADS

Problem

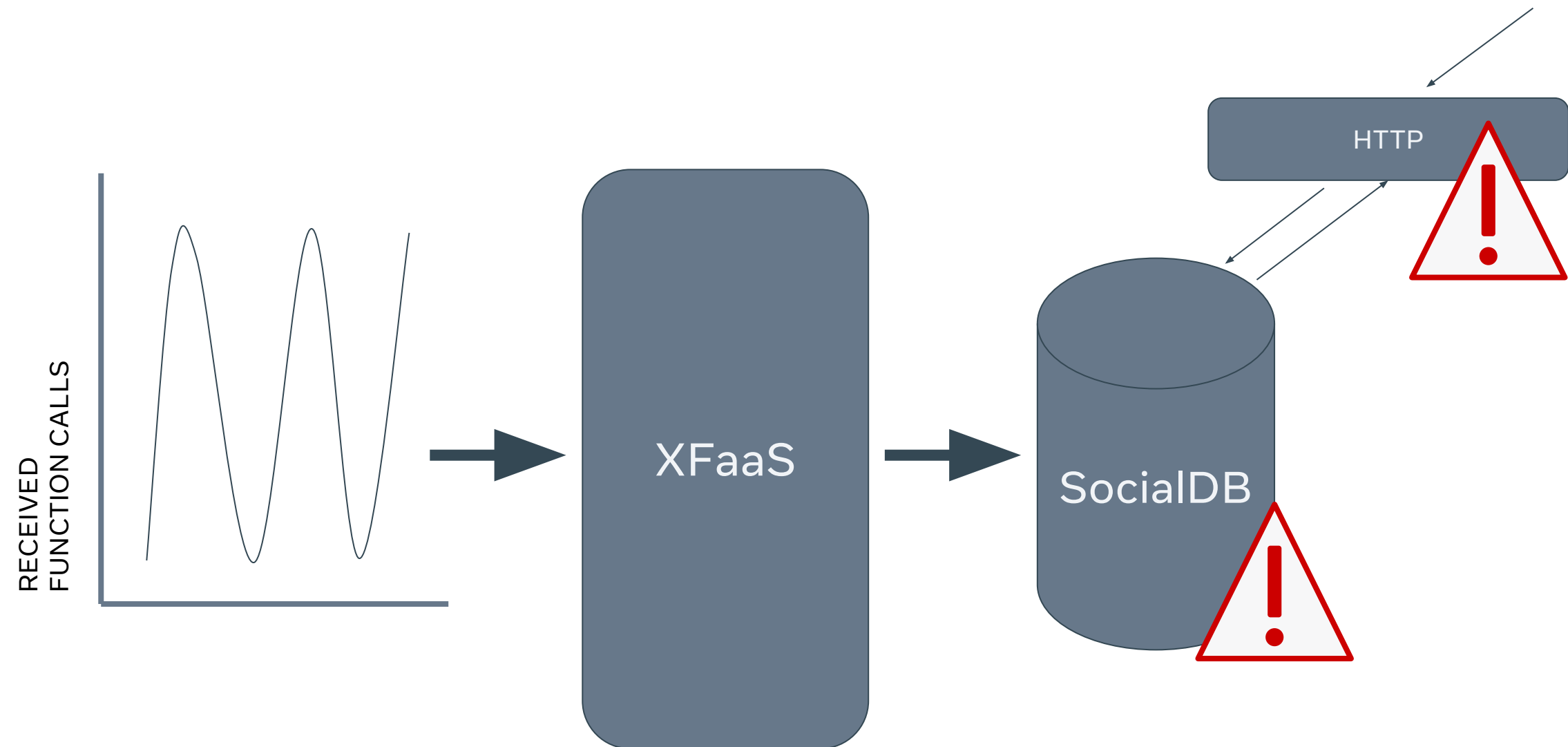


DOWNSTREAM OVERLOADS

Problem

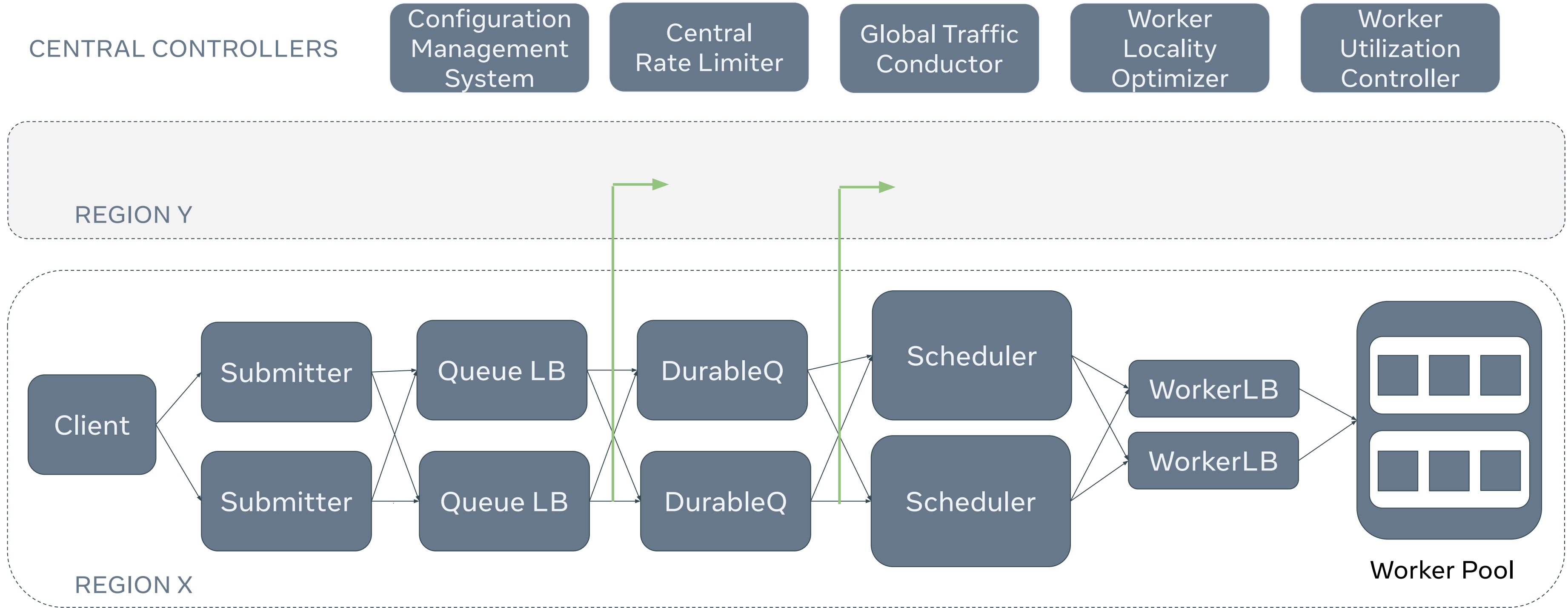
SocialDB Outage

- manual resolution
- several hours to resolve
- coarse-grained



03 SYSTEM OVERVIEW

03 SYSTEM OVERVIEW



Next...

04 DEFERRED COMPUTE - DESIGN & EVALUATION

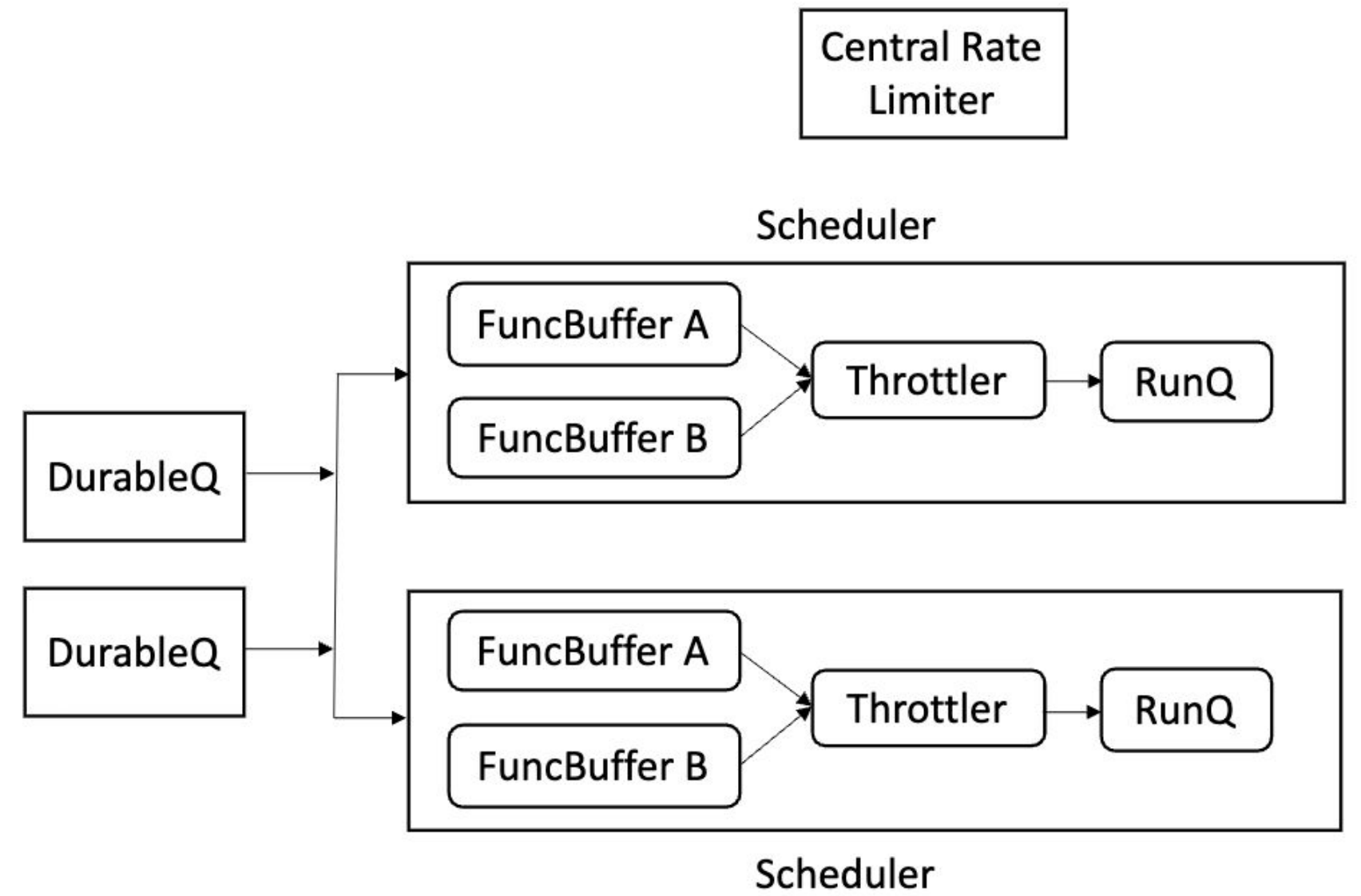
05 DOWNSTREAM PROTECTION - DESIGN & EVALUATION

04 DEFERRED COMPUTE - DESIGN & EVALUATION

1. Reserved Quota

- CPU cycles a function can consume
- Transformed to RPS for enforcement

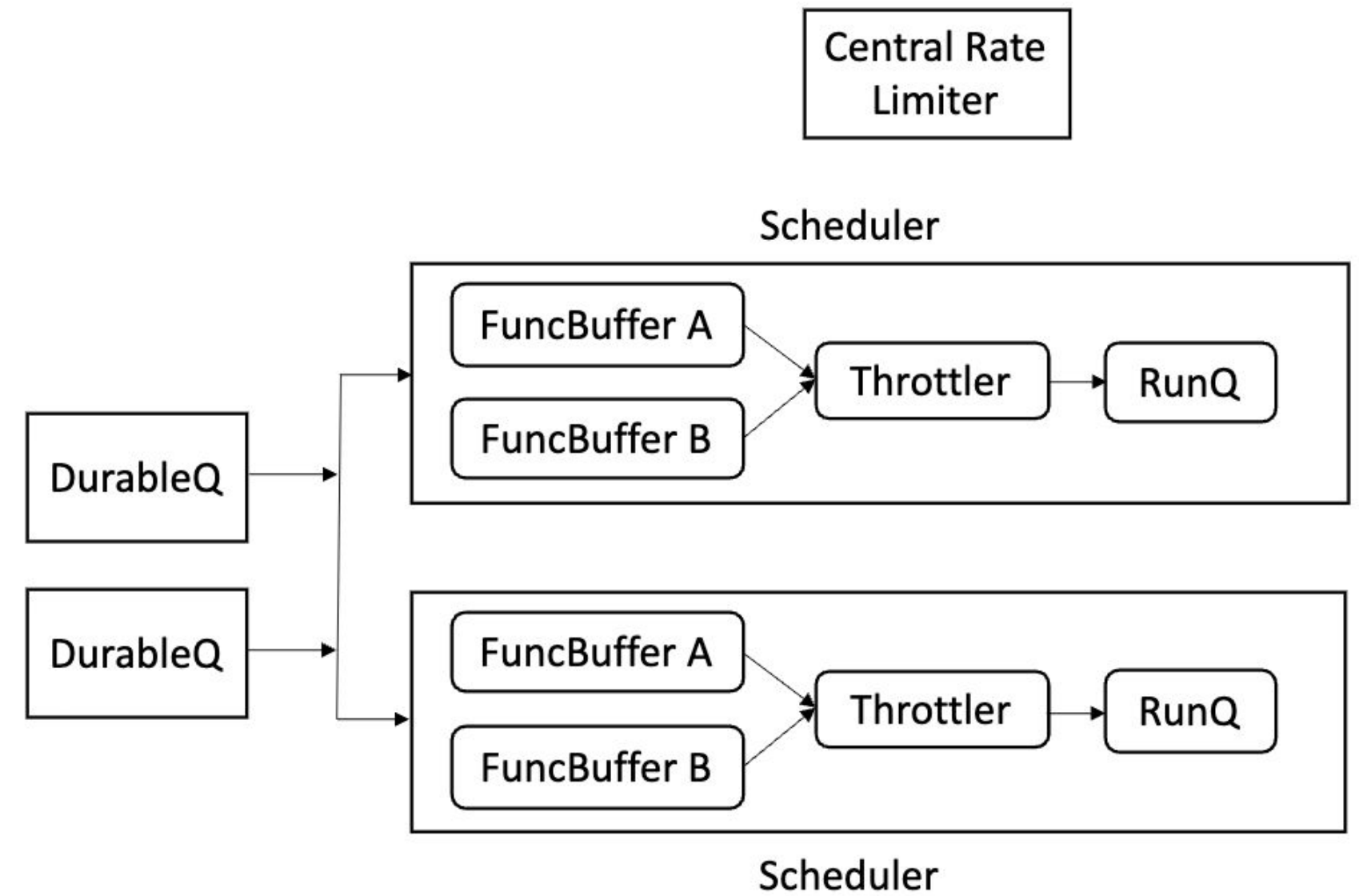
2. Opportunistic Quota



1. Reserved Quota

2. Opportunistic Quota

- Dynamically adjusted based on worker utilization
- Deferred to off-peak hours
- SLA of 24 hrs



$$throttling_rate = base_rate_from_quota * S$$

High Worker Utilization: $S \searrow$

Low Worker Utilization: $S \nearrow$

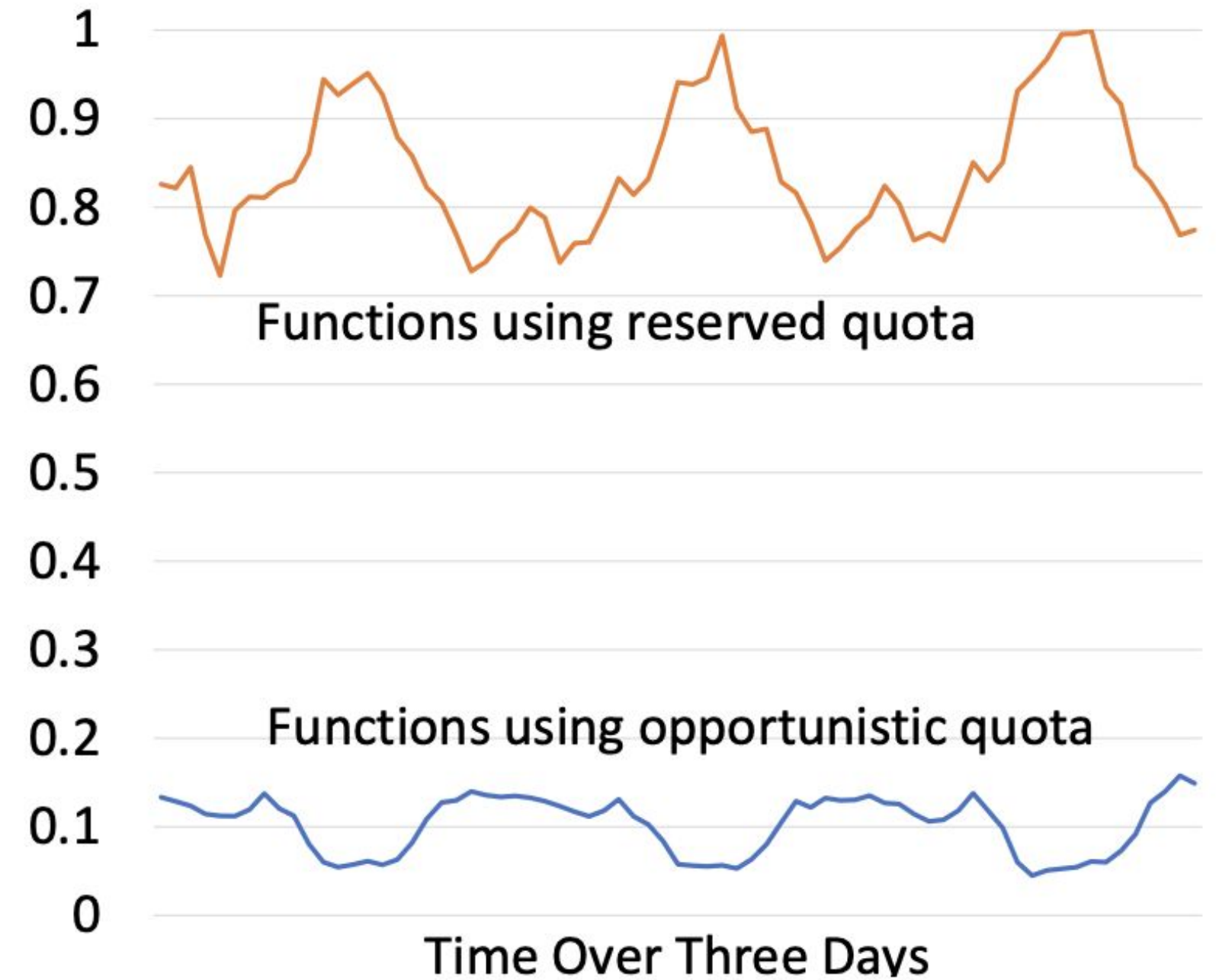
1. Reserved Quota
2. Opportunistic Quota

3. Per function criticality level
4. Explicit future execution start time

[NOT COVERED IN THIS TALK]

- Daily Peak Pattern
- Opportunistic Functions are Throttled during Peak

Total CPU Cycles Consumed by Functions



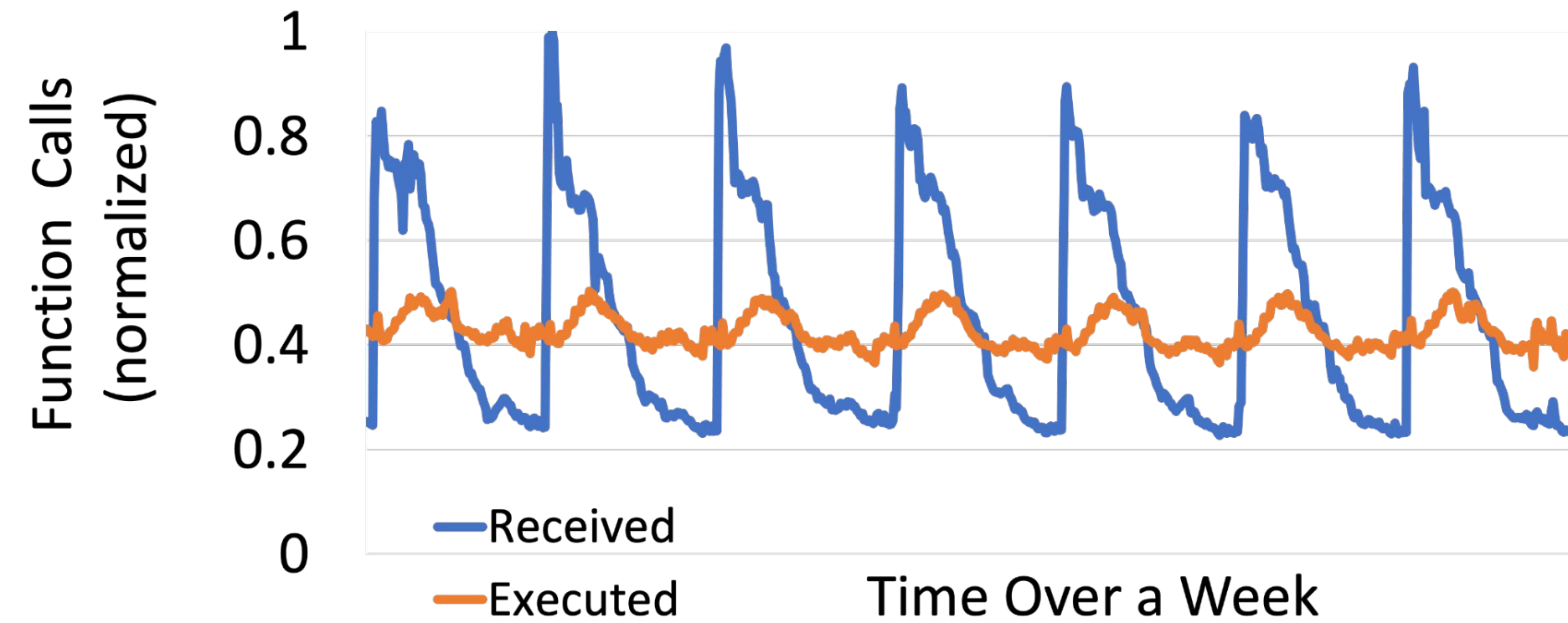
All Deferred Compute Features at Work

- Reserved Quota
- Opportunistic Quota
- Per Function Criticality
- Explicit future execution time

Cross Regional Load Balancing

Results:

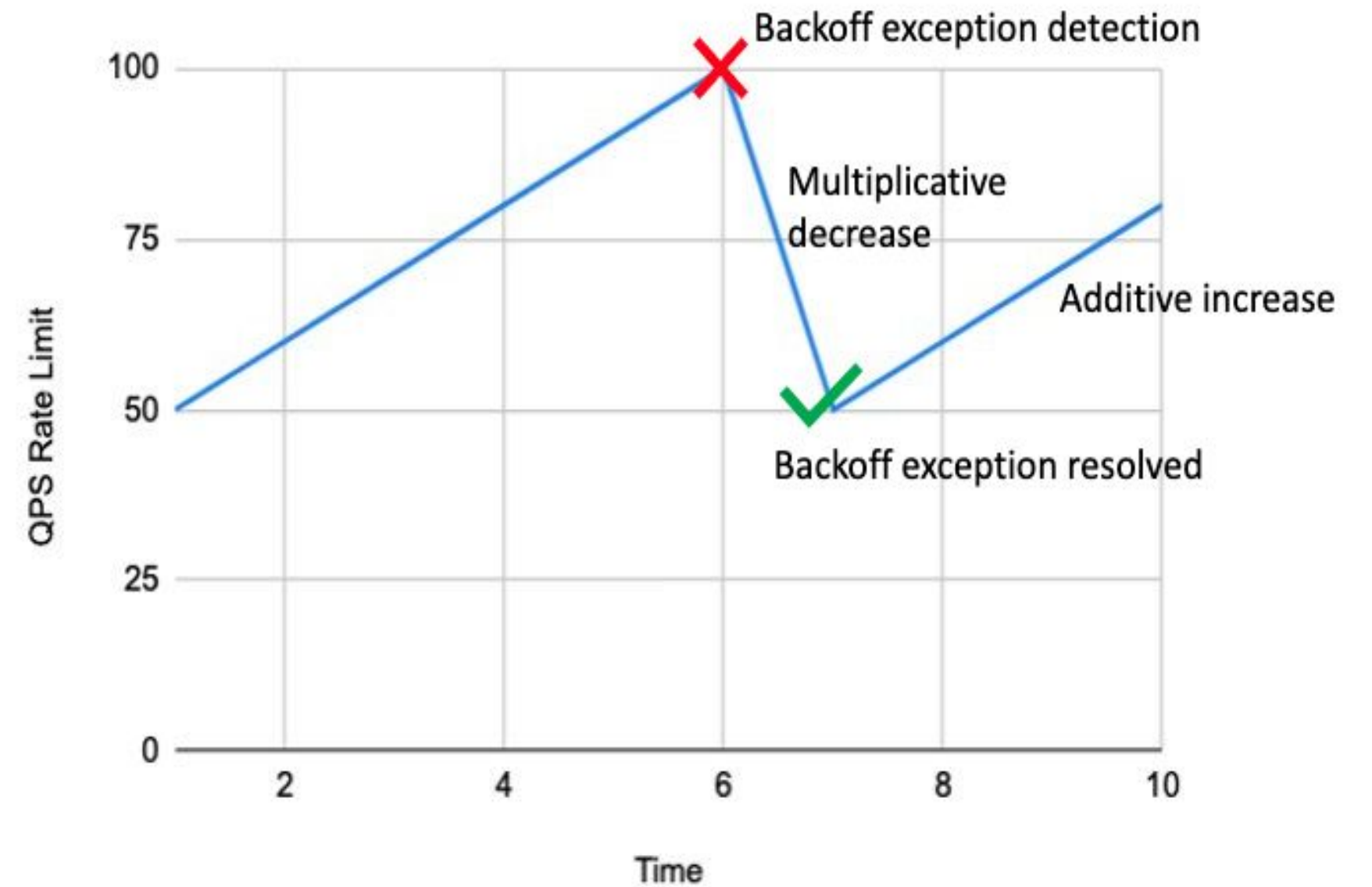
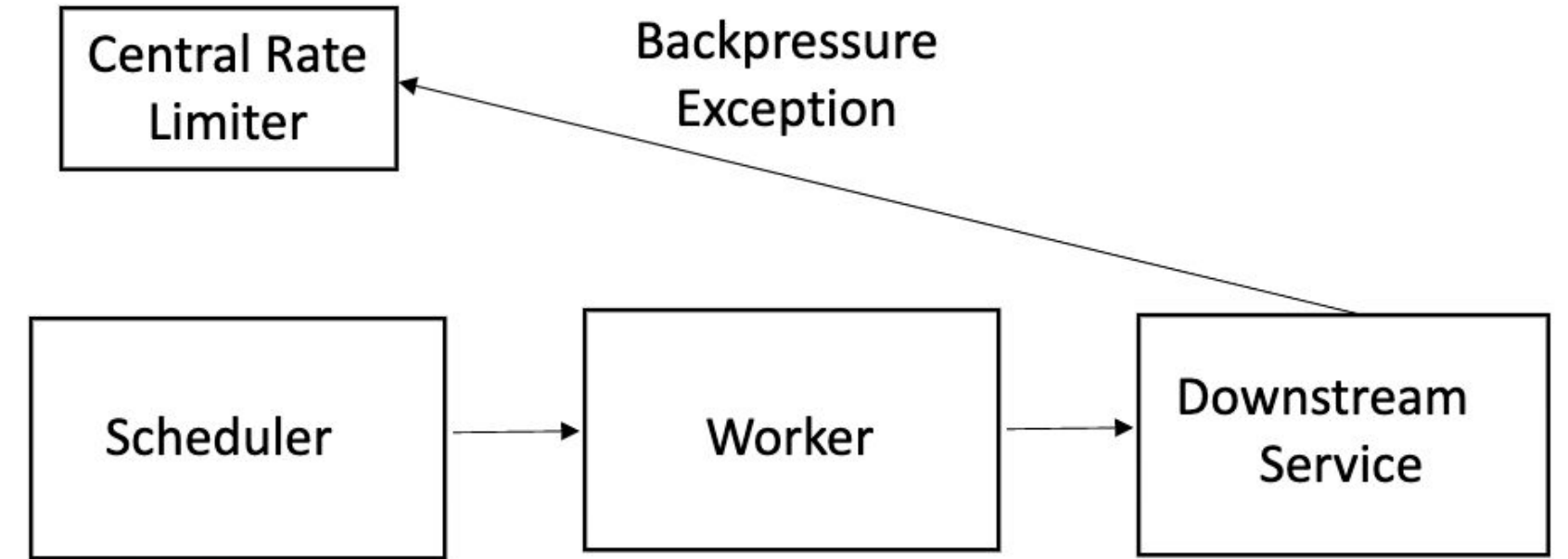
- PeaktoTrough reduced from 4.3x to 1.4x
- 66% Daily Average CPU Utilization



05 DOWNSTREAM PROTECTION - DESIGN & EVALUATION

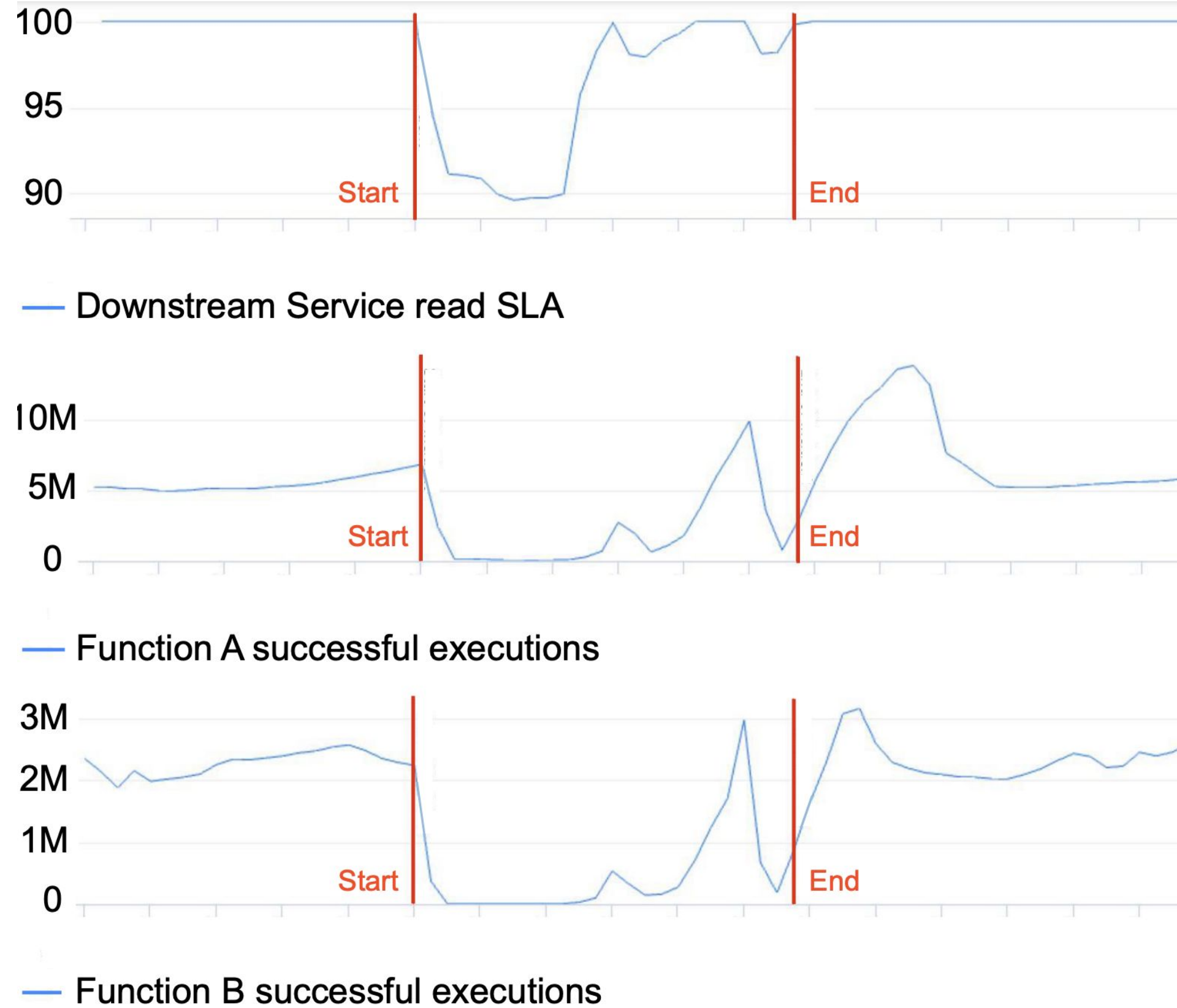
Backpressure Handling

- Responds to Downstream Backpressure Exceptions
- Throttling rate is set by AIMD algorithm



05 DOWNSTREAM PROTECTION - DESIGN & EVALUATION

- Real incident during overload of WTCache in front of Social DB (TAO¹)
- Recovery was complete in two hours without any engineering intervention



¹ Nathan Bronson, et al. "TAO: Facebook's Distributed Data Store for the Social Graph." In Proceedings of the 2013 USENIX Annual Technical Conference, 2013

XFaaS

HYPERSCALE AND LOW COST SERVELESS FUNCTIONS AT META

06 Summary

$O(10^{12})$

)
Function
Calls/
Day

$O(10^5)$

Servers

>10

Regions

- XFaaS utilizes the concept of universal workers to eliminate cold start [NOT COVERED IN THIS TALK]
- Even if we eliminate cold start, we will still be often underutilized with need to autoscale almost instantaneously by 4x
- XFaaS embodies several methods to smoothen out the function execution curve => **daily avg CPU utilization at 66%**
- Ensures protection of downstream services



FUNCTION AS A SERVICE AT META

- function types
- workload examples

Triggers	Functions	Function Calls	Compute Usage
Queue-triggered	89%	15%	86%
Event-triggered	8%	85%	14%
Timer-triggered	3%	<1%	<1%

FUNCTION AS A SERVICE AT META

- highly heterogeneous workloads
- workload examples

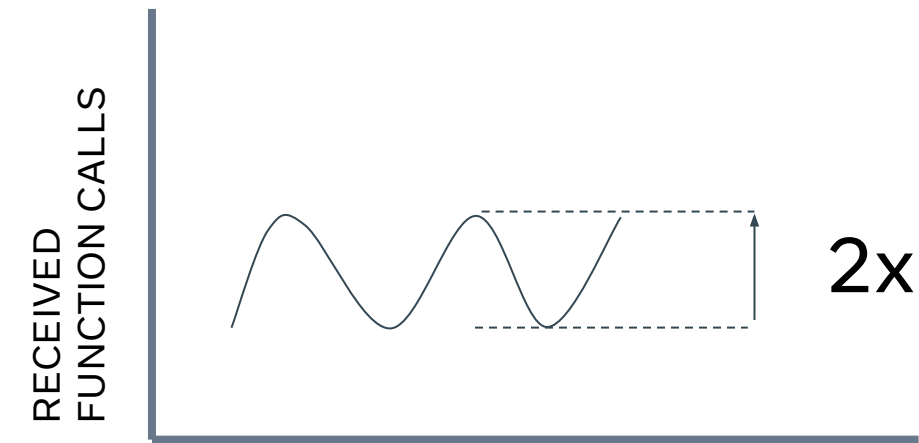
CPU Usage (MIPS)				
Function Type	P10	P50	P90	P99
Queue-triggered	20.40	221.80	7,611	1,064,280
Event-triggered	0.54	11,36	189	2,981
Timer-triggered	0.37	576.00	44,839	369,282

High Variance of Load

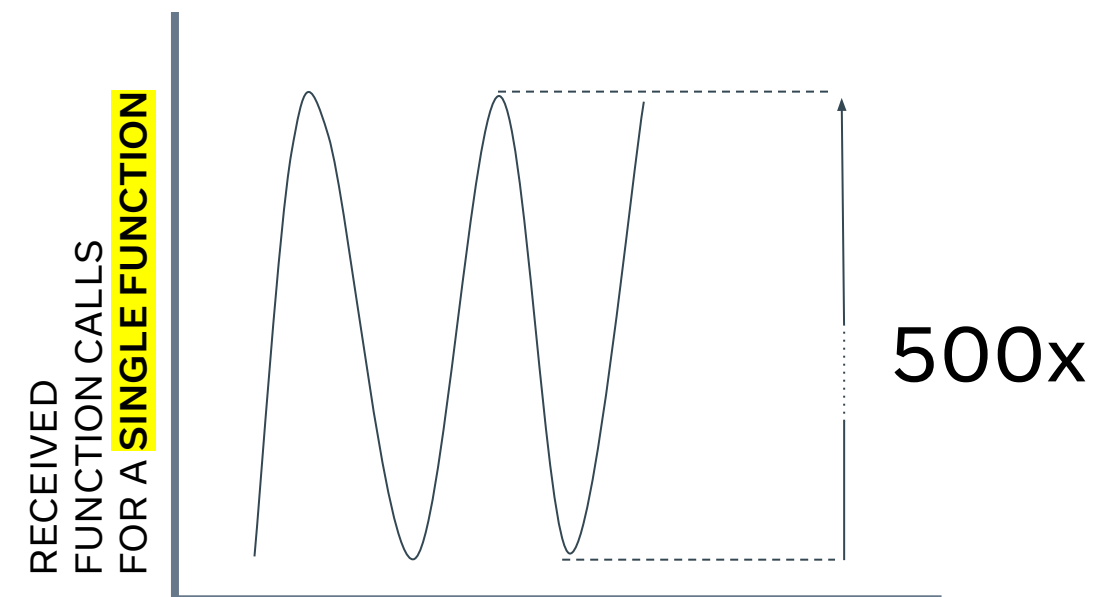
Problem

1. Previous work reported a high peak-to-trough ratio of function calls
2. At Meta, the ratio can be as high as 4.3

Shahrad et al. Serverless in the wild: Characterizing and optimizing the serverless workload at a large cloud provider. In USENIX Annual Technical Conference (USENIX ATC 20). 2020.



Wang, Ao, et al. "{FaaSNet}: Scalable and fast provisioning of custom serverless container runtimes at alibaba cloud function compute." In USENIX Annual Technical Conference (USENIX ATC 21). 2021.



Lengthy Cold Start

Problem

1. Cloud Provider overhead costs
2. Actual work (only step 8)

[NOT COVERED IN THIS TALK]

INITIALIZATION PHASE:

- (1) Start the VM.
- (2) Fetch the container image and the function's code.
- (3) Initialize the container.
- (4) Start the language runtime such as Python or PHP.
- (5) Load common libraries into memory.
- (6) Load the function code into memory.
- (7) Optionally, do JIT compilation.

INVOCATION PHASE:

- (8) Invoke the function multiple times as needed.

SHUTDOWN PHASE:

- (9) Stop the container if it receives no requests for X minutes (X=10/20/10 minutes for AWS/Azure/OpenWhisk respectively).
- (10) Optionally, stop the VM.

Lengthy Cold Start

XFaaS Solution: Universal Worker

1. Proactive code deployment (skip 1-5)
2. Function process colocation (skip)
3. Cooperative JIT compilation (skip 6-7 for regularly invoked functions)
4. Locality Groups (helps skip 6-7 by improving JIT code cache hit rate)

INITIALIZATION PHASE:

- (1) Start the VM.
- (2) Fetch the container image and the function's code.
- (3) Initialize the container.
- (4) Start the language runtime such as Python or PHP.
- (5) Load common libraries into memory.
- (6) Load the function code into memory.
- (7) Optionally, do JIT compilation.

INVOCATION PHASE:

- (8) Invoke the function multiple times as needed.

SHUTDOWN PHASE:

- (9) Stop the container if it receives no requests for X minutes (X=10/20/10 minutes for AWS/Azure/OpenWhisk respectively).
- (10) Optionally, stop the VM.