AWS re:Invent

# SERVERLESS AT SCALE IS THE NEW NORM

**THOMSON REUTERS**

processes **4,000 requests** per second

**mlbam**

ingests, analyzes and stores **17+ petabytes of data** per season

**FINra**

processes **half a trillion** validations of stock trades daily

**Driver & Vehicle Licensing Agency**

API traffic to register and license more than **47 million driver records** in Great Britain,

**Zillow®**

executes **16 million** requests a month

**Localytics**

processes **tens of billions** of data points monthly

aws

# Running Highly Available Large Scale Systems Is a Lot of Work

Load Balancing

Scaling
Up and
Down

aws

# Handling Failures

aws

# Lambda Handles

**Load Balancing**

**Auto Scaling**

**Handling Failures**

**Security Isolation**

**Managing Utilization**

(and many other things) for you

# Front End Invoke

Orchestrate both synchronous and asynchronous Invokes

# Counting Service

Provides a region wide view of customer concurrency to help enforce set limits

aws

# Worker Manager

Tracks container idle and busy state and schedules incoming invoke requests to available containers
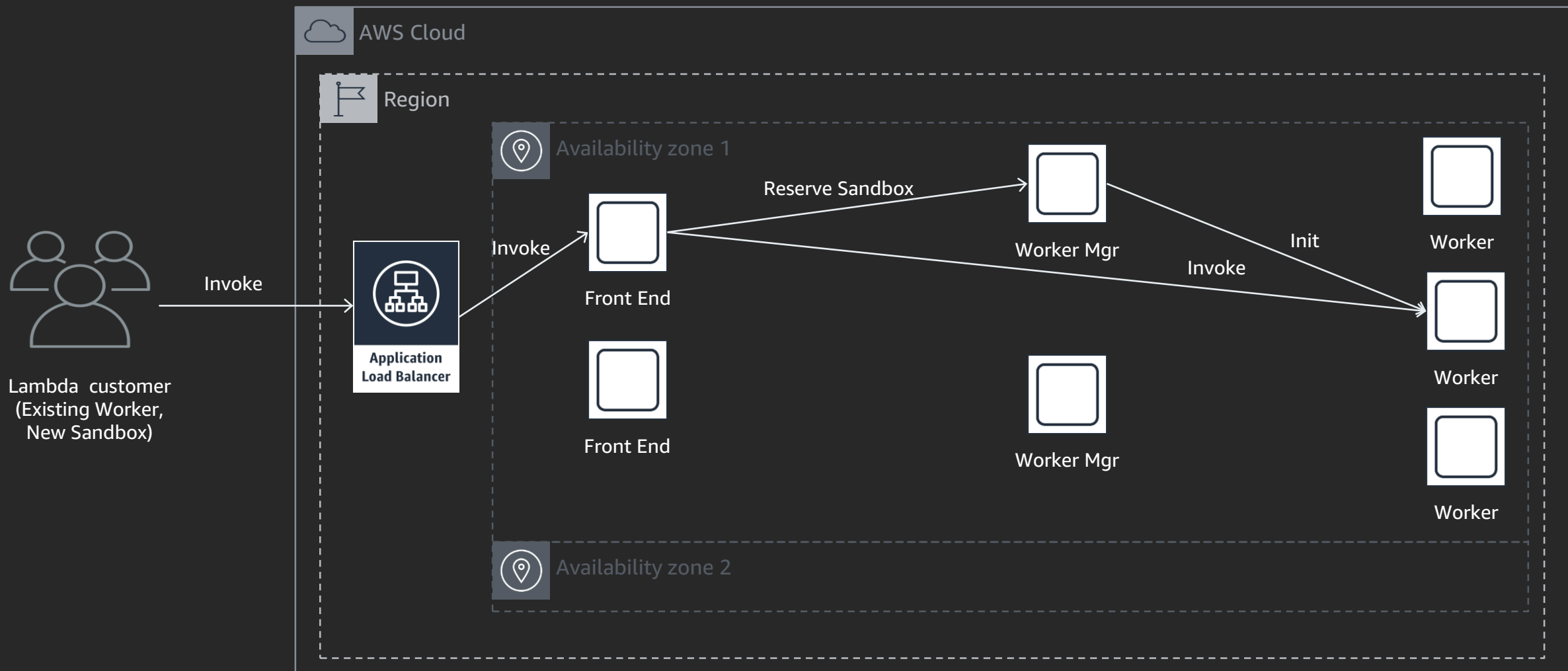
aws

# Worker

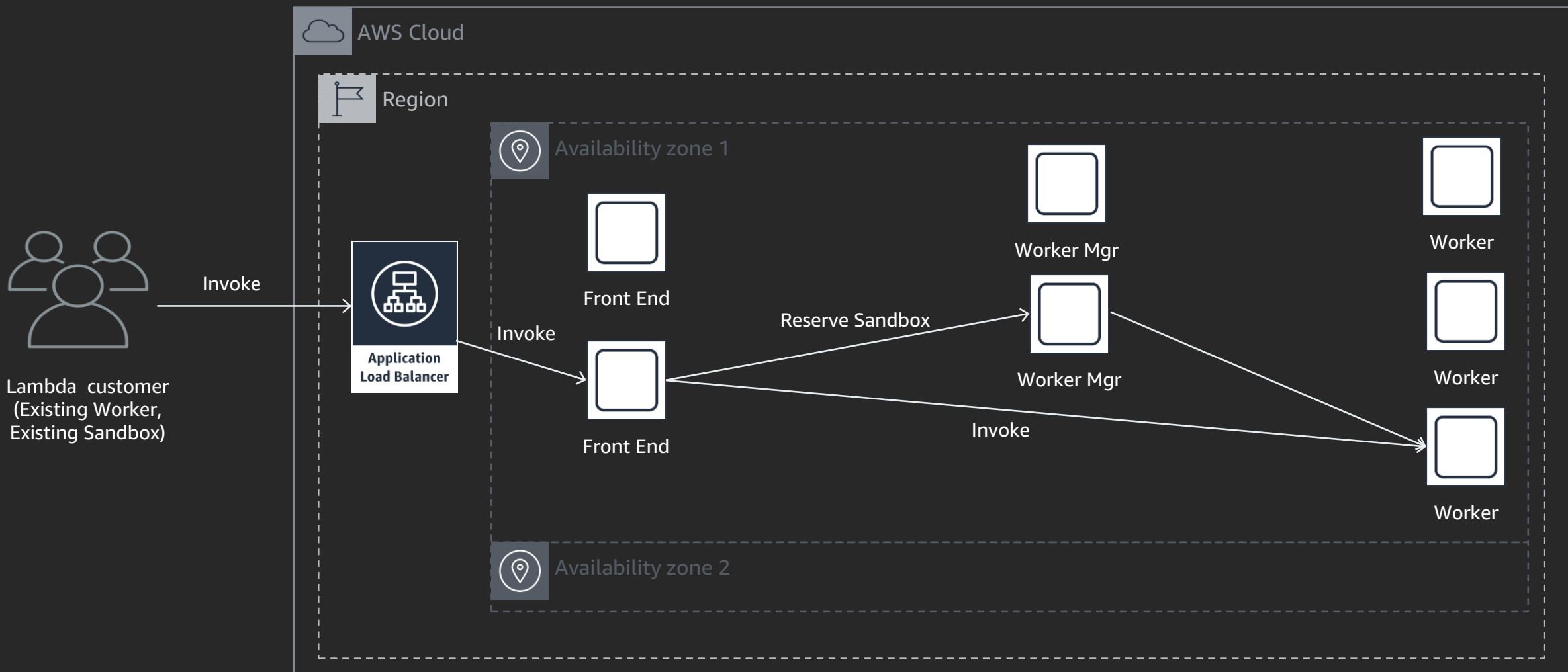Provisions a secure environment for customer code execution

aws

# Placement Service

Places sandboxes on workers to maximize packing density without impacting customer experience or cold-path latency

# Load Balancing

## Routing function traffic across hosts distributed across Availability Zone
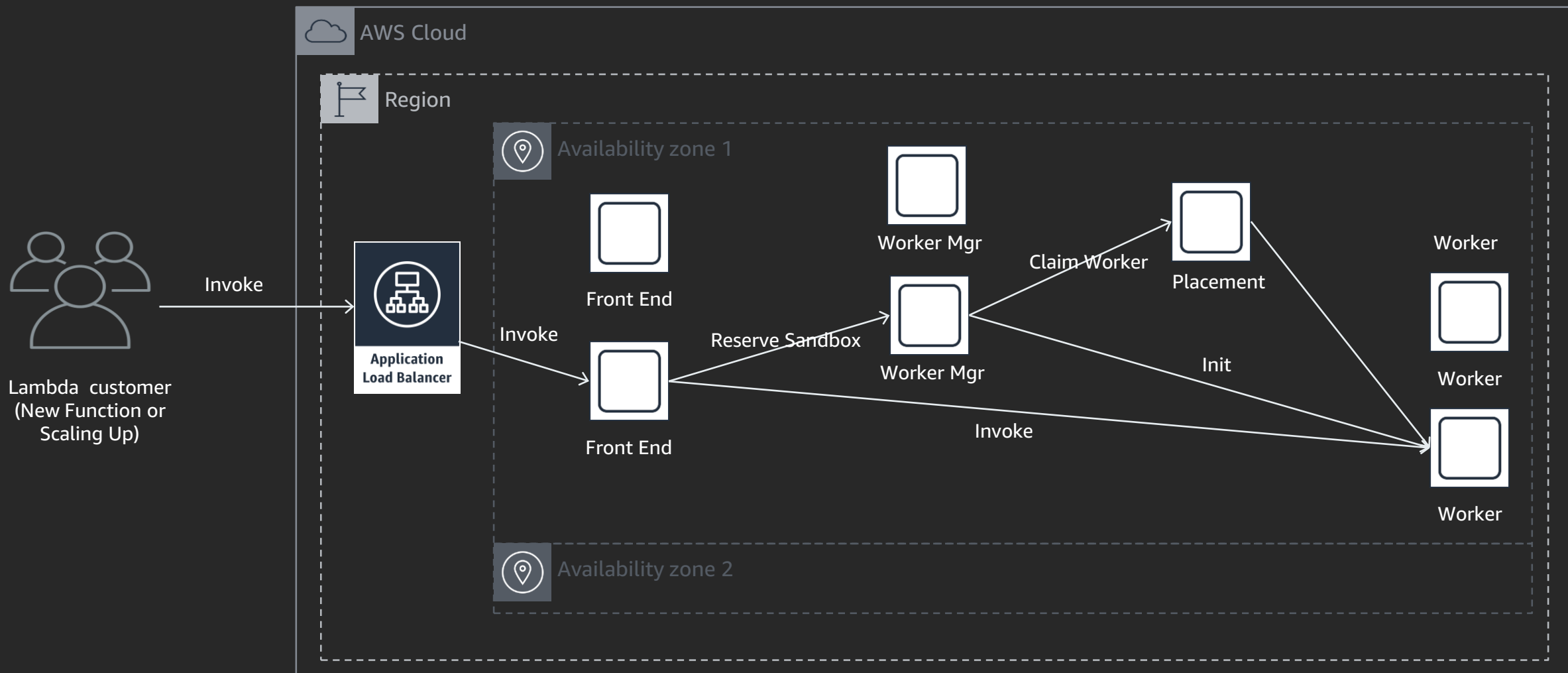
aws

# SERVERLESS CUSTOMERS

AWS re:Invent

aws

edmunds

"What took us just a few days to build using a serverless solution based on AWS Lambda would have taken us six months to build from scratch. Our CTO and the rest of the project stakeholders were really happy with how much money and time we saved."
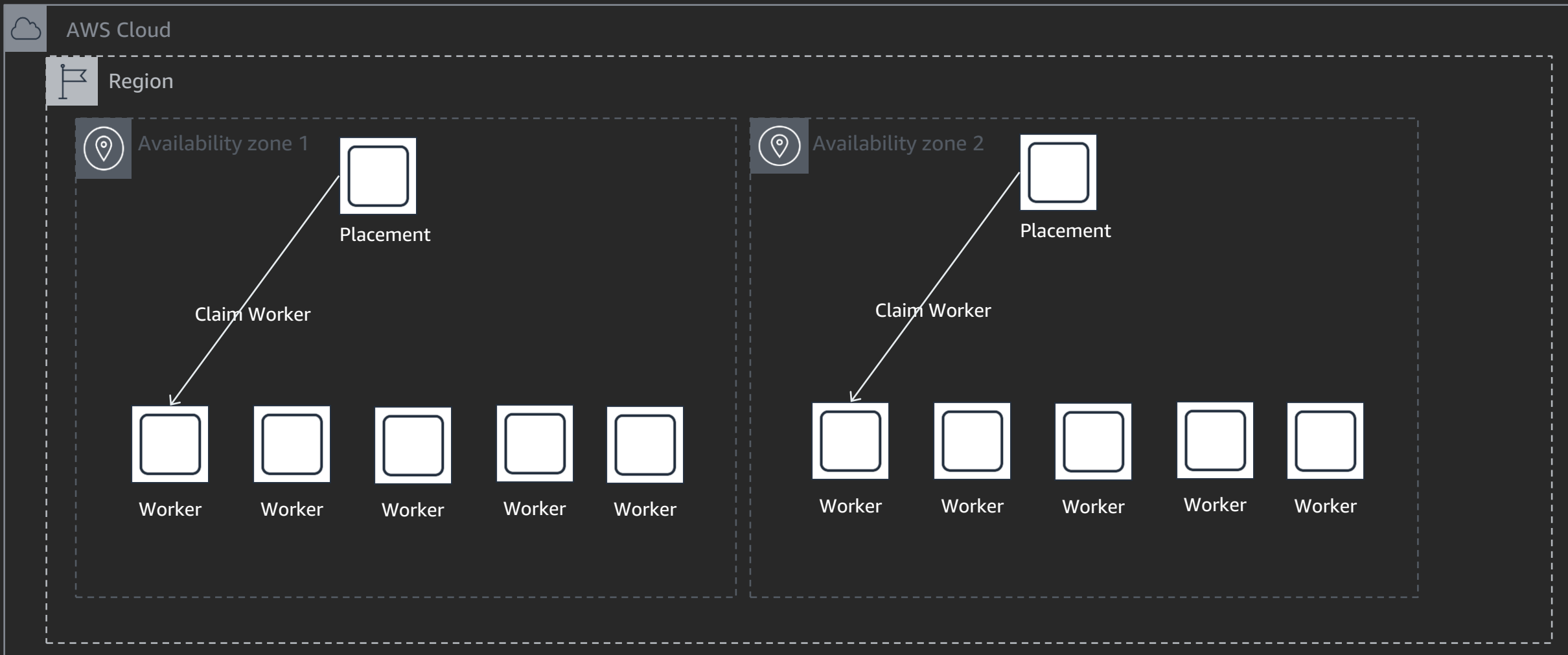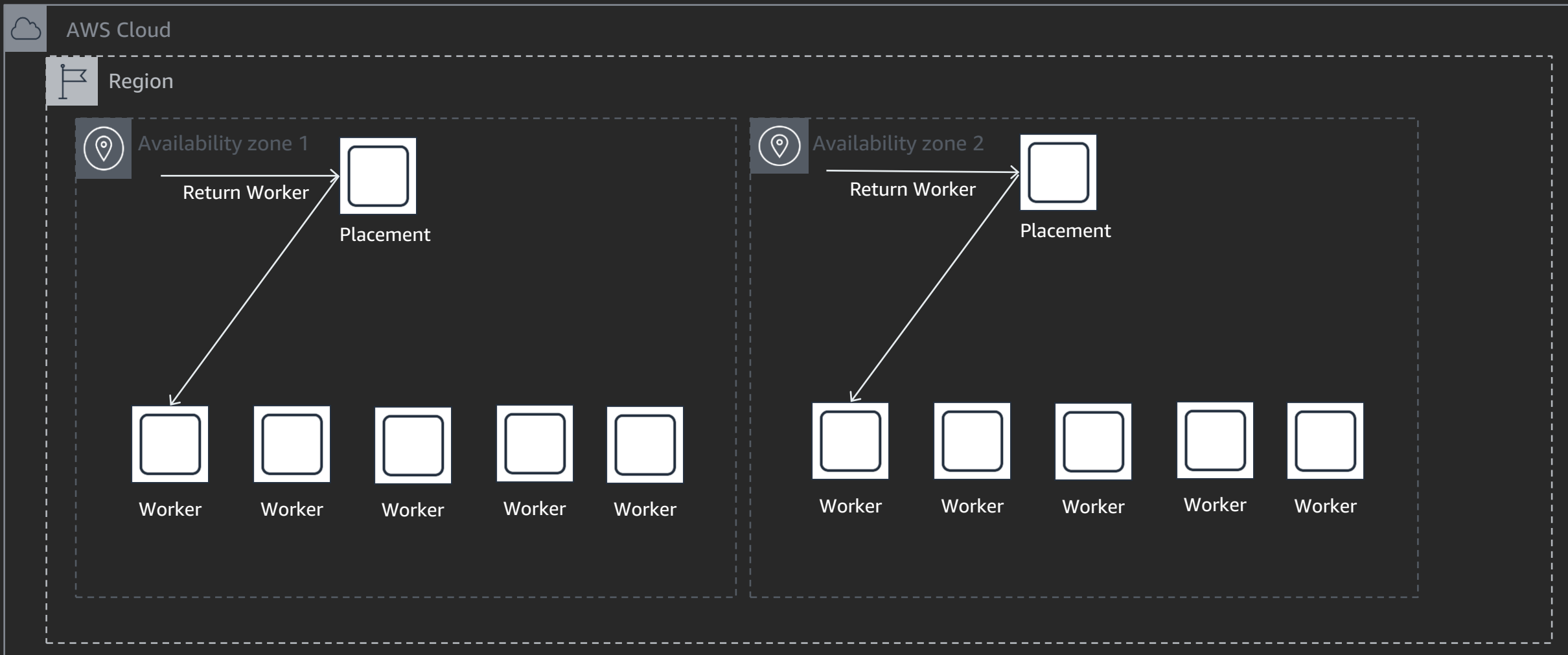
**Nitin Mahajan**
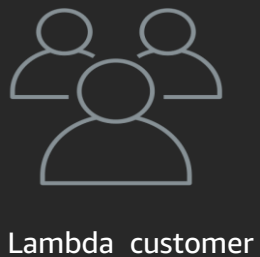Executive Director for Service Engineering

aws

# Auto Scaling

Provision function capacity when needed and releasing when not needed

AWS Cloud

Region

Availability zone 1

Front End

Worker Mgr

Front End

Worker Mgr

Placement

Worker

Worker

Worker

Availability zone 2

Lambda customer
(New Function or
Scaling Up)

Invoke

Invoke

Reserve Sandbox

Claim Worker

Init

Invoke

Application
Load Balancer

AWS re:Invent

AWS Cloud

Region

Availability zone 1

Return Worker

Placement

Worker    Worker    Worker    Worker    Worker

Availability zone 2

Return Worker

Placement

Worker    Worker    Worker    Worker    Worker

# SERVERLESS CUSTOMERS

# Handling Failures

## Handling Host and Availability Zone failure

AWS re:Invent

aws

# With Lambda:
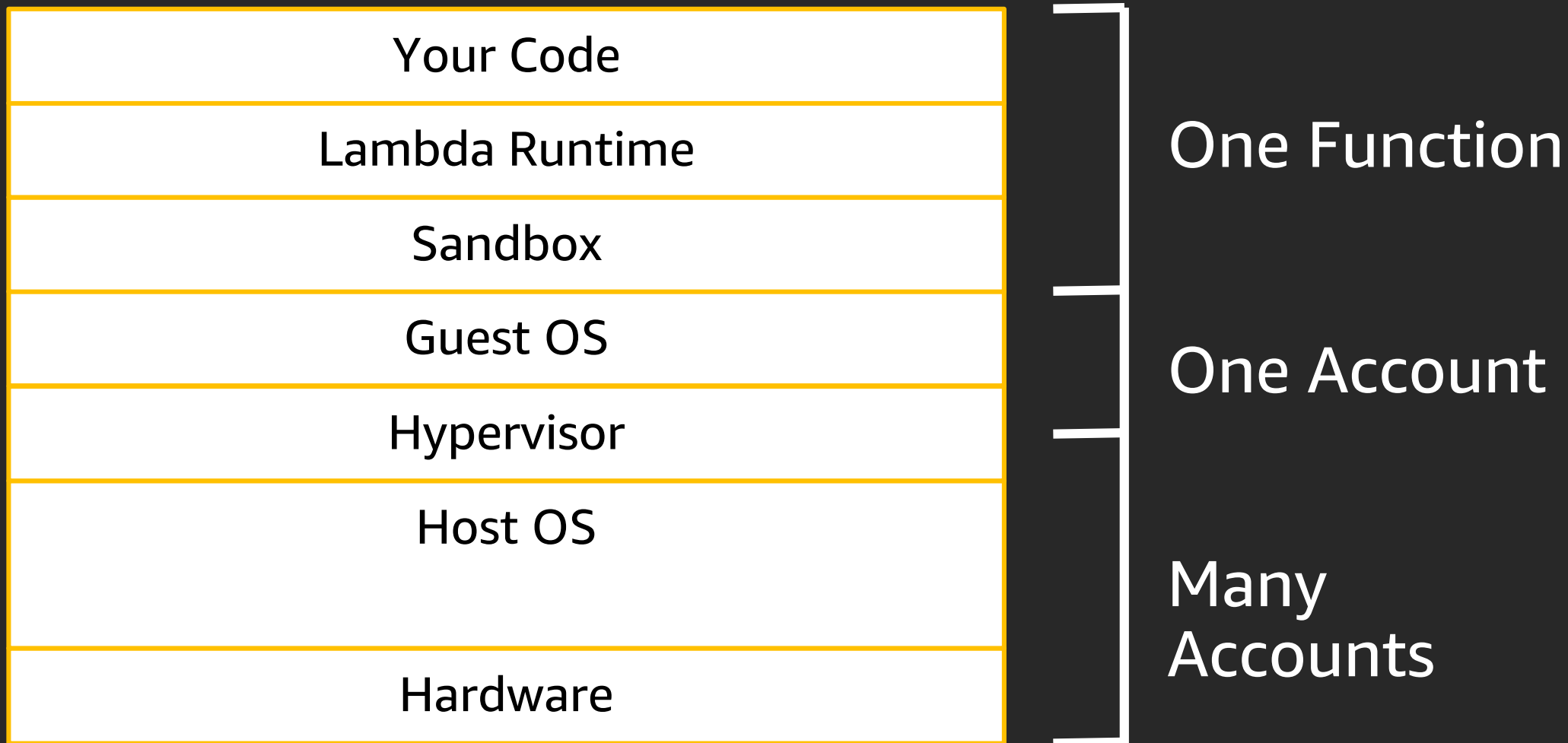# Always have a healthy host

AWS
re:Invent

aws

# Isolation
## Keeping Workloads Safe and Separate

Your Code

Lambda Runtime

Sandbox

Guest OS

Hypervisor

Host OS

Hardware

Your Code

Lambda Runtime
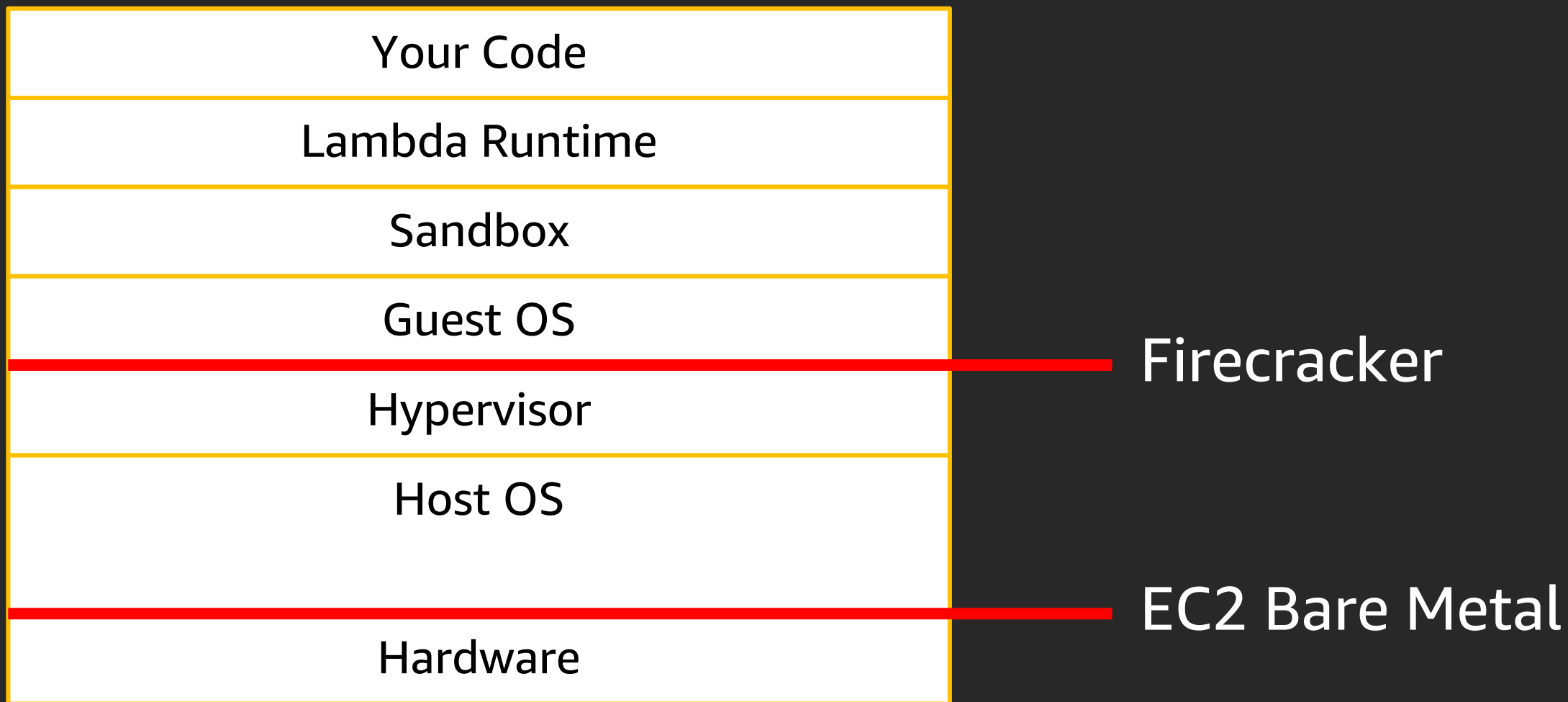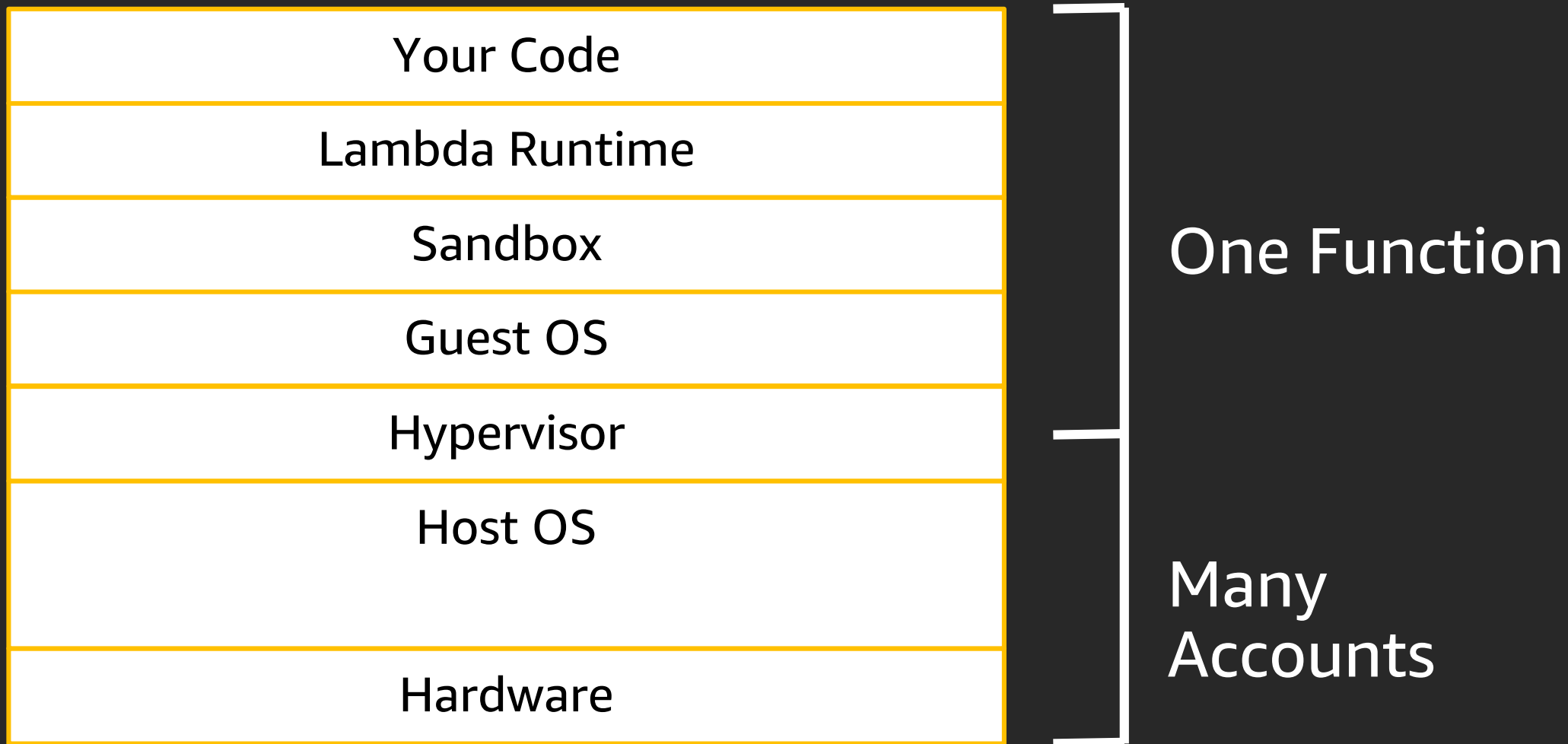
Sandbox

Guest OS

Hypervisor

Host OS

Hardware

One Function

One Account

Many
Accounts

| Your Code |
| Lambda Runtime |
| Sandbox |
| Guest OS |
| Hypervisor |
| Host OS |
| Hardware |

cgroups,
namespaces,
seccomp,
iptables,
& chroot

aws

| |
|---|
| Your Code |
| Lambda Runtime |
| Sandbox |
| Guest OS |
| Hypervisor |
| Host OS |
| Hardware |

virtualization & device emulation

aws

| |
|---|
| Your Code |
| Lambda Runtime |
| Sandbox |
| Guest OS |
| Hypervisor |
| Host OS |
| Hardware |

Firecracker

EC2 Bare Metal

aws

# Innovating to Improve Isolation



Guest OS

Virtual Devices

Hypervisor

Host OS

Device Emulation

Hardware

Physical Devices

# Innovating to Improve Isolation

Guest OS

*virtio drivers*

Hypervisor

Host OS

*virtio* host in Firecracker

Hardware

Physical Devices

# Innovating to Improve Isolation

# Utilization

## Keeping Servers Busy

# % of Resources Doing Useful Work

## (vs. idle or waste)

# With Lambda:
# Pay only for useful work.

# Inside Lambda: Optimize To Keep Servers Busy

aws

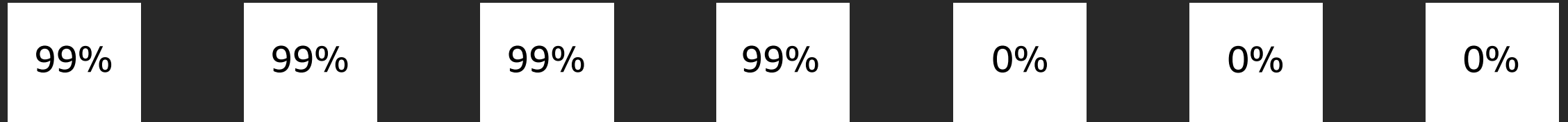# Available Sandboxes For a Function
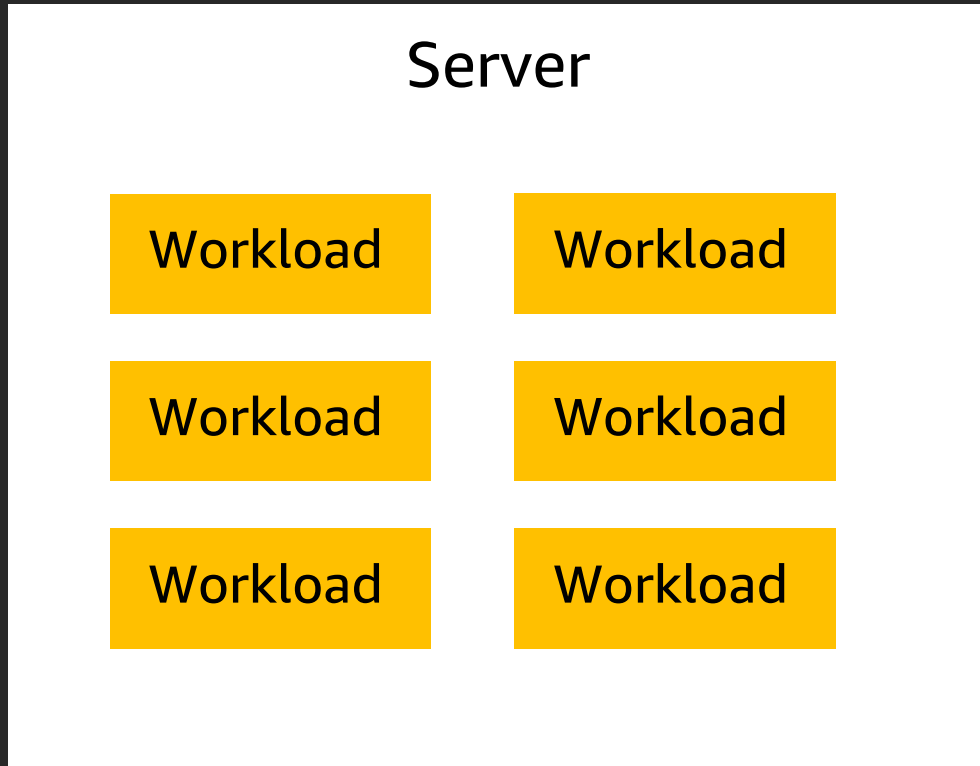
# Bad: Balance The Load

| 60% | 60% | 60% | 60% | 60% | 60% | 60% |

AWS re:Invent

aws

# Good: Concentrate The Load

| 99% | 99% | 99% | 99% | 0% | 0% | 0% |
|-----|-----|-----|-----|-----|-----|-----|

# Cache Locality
# Ability to Autoscale

aws

# Bad: Pack Server With One Workload

Server

Workload    Workload

Workload    Workload

Workload    Workload

aws

# Better: Pack With Many Loads

**Server**

Workload
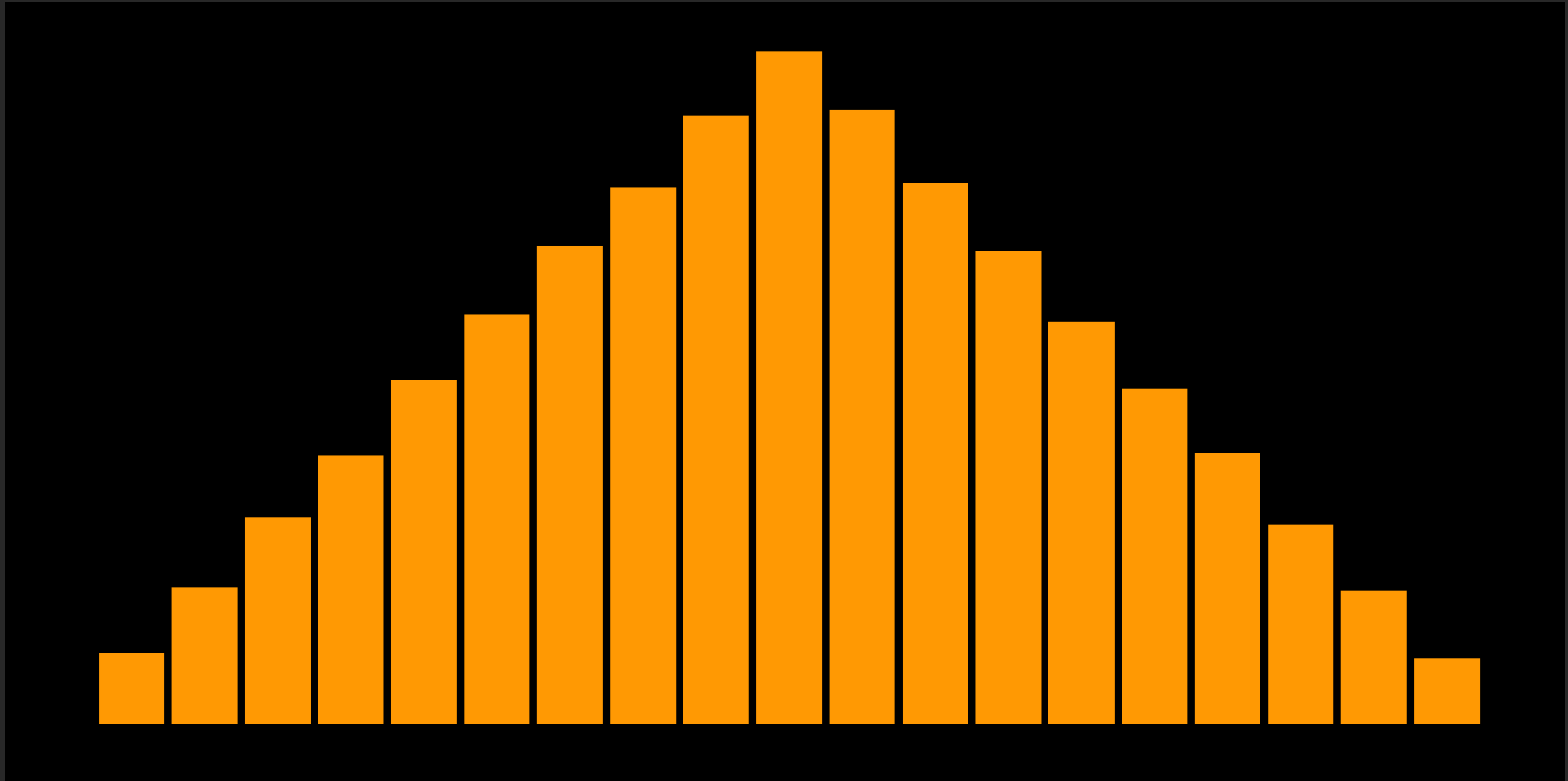Workload
Workload
Workload
Workload
Workload

Take advantage of *Statistical Multiplexing*
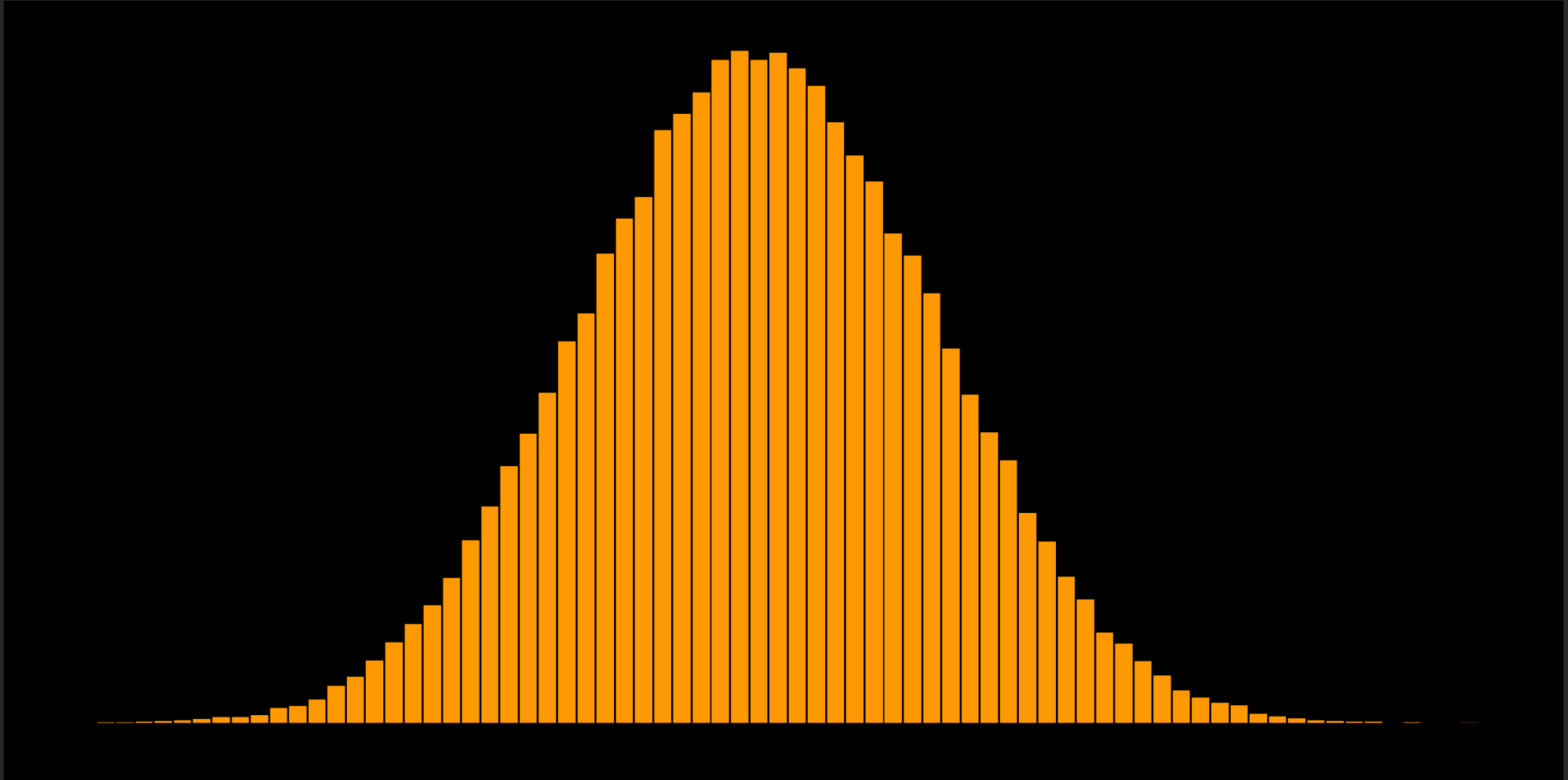
# Throwing a d20

aws

# Throwing a d10, twice

aws

# Throwing a d10, ten times

aws

# Best: Placement Optimization

**Server**

| | |
|---|---|
| Workload | Workload |
| Workload | Workload |
| Workload | Workload |

Pick workloads that pack together well.

Minimize contention.

aws

# Improving VPC start-up and scaling: now

Worker

Lambda
Function

Local
NAT

ENI in
your VPC

Your VPC

# Improving VPC start-up and scaling: 2019

**Secure Tunnel**

Worker

Lambda
Function

Remote
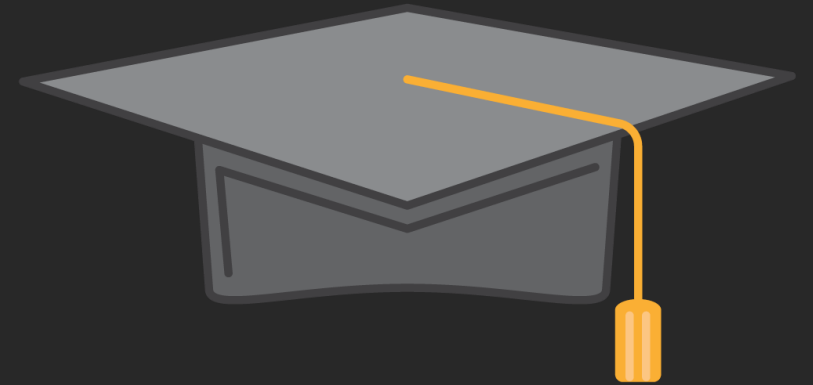NAT

ENI in
your VPC

Your VPC

# Improving VPC start-up and scaling: benefits



Faster
Scaling

Lower
Latency

Easier To
Use

# Firecracker Hypervisor vs. Others

↓ Startup time

↓ Memory overhead

= Performance

↑ Flexibility

aws

# Firecracker Unlocks Higher Utilization and Scale

# In Conclusion

AWS
re:Invent

aws

# Thank you!

Marc Brooker – Senior Principal Engineer, Amazon Serverless
Holly Mesrobian – Director of Engineering, Amazon Lambda

aws

Please complete the session survey in the mobile app.