# Scaling Message Passing on Amazon Web Services with Elastic Fabric Adapter
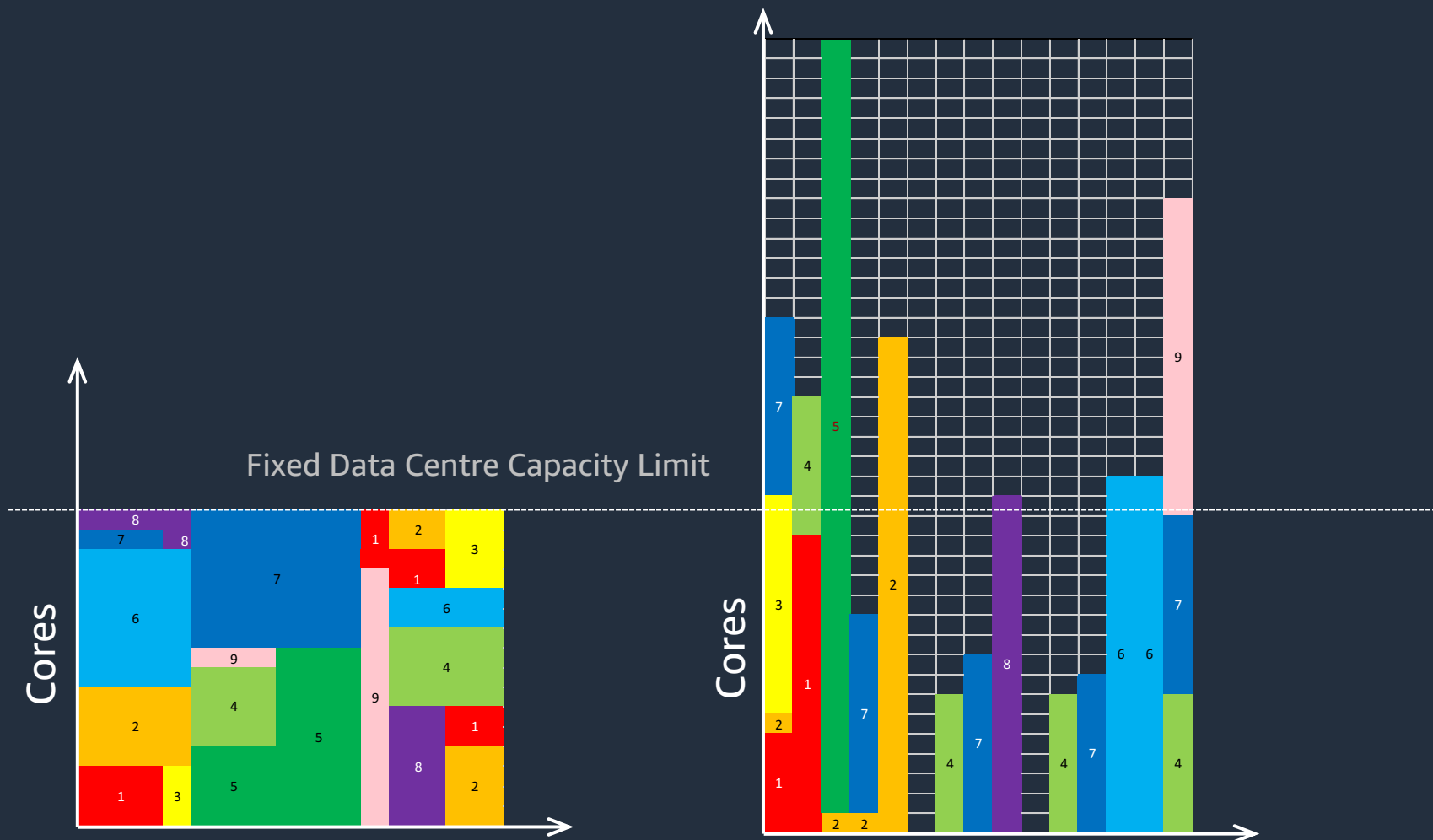
Raghunath Raja Chandrasekar

# Agenda

- Overview of HPC on AWS

- Evolution of networking on AWS

- What is EFA?

- Scalable Reliable Datagrams (SRD) Protocol

- MPI implementations and support for EFA

- Application scaling case-studies

aws

# HPC on AWS

aws

# Time-to-results as the Metric for Success



Finite capacity, usually with long queues to wait in.

Massive capacity when needed to speed up time to results, and agile environment when additional hardware and software experimentation is needed.

# HPC on AWS: Solution Components

## Automation and orchestration

- AWS Batch
- AWS ParallelCluster
- NICE EnginFrame

## Storage

- Amazon EBS
- Amazon FSx for Lustre
- Amazon EFS
- Amazon S3

## Compute

- Amazon EC2 instances (Compute and accelerated)
- Amazon EC2 Spot
- AWS Auto Scaling

## Visualization

- NICE DCV
- Amazon AppStream 2.0

## Networking

- Enhanced networking
- Placement groups
- Elastic Fabric Adapter

aws

# Broad HPC Partner Community
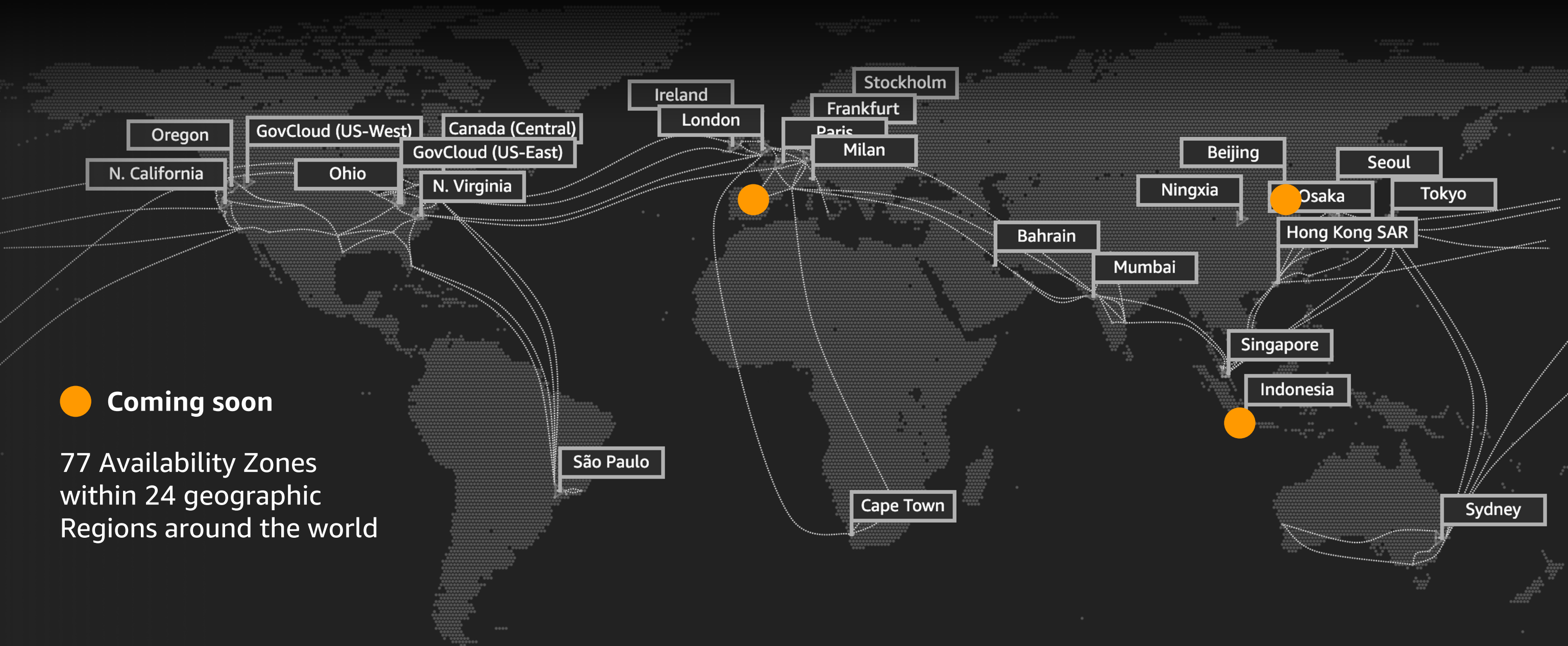
## Application partners
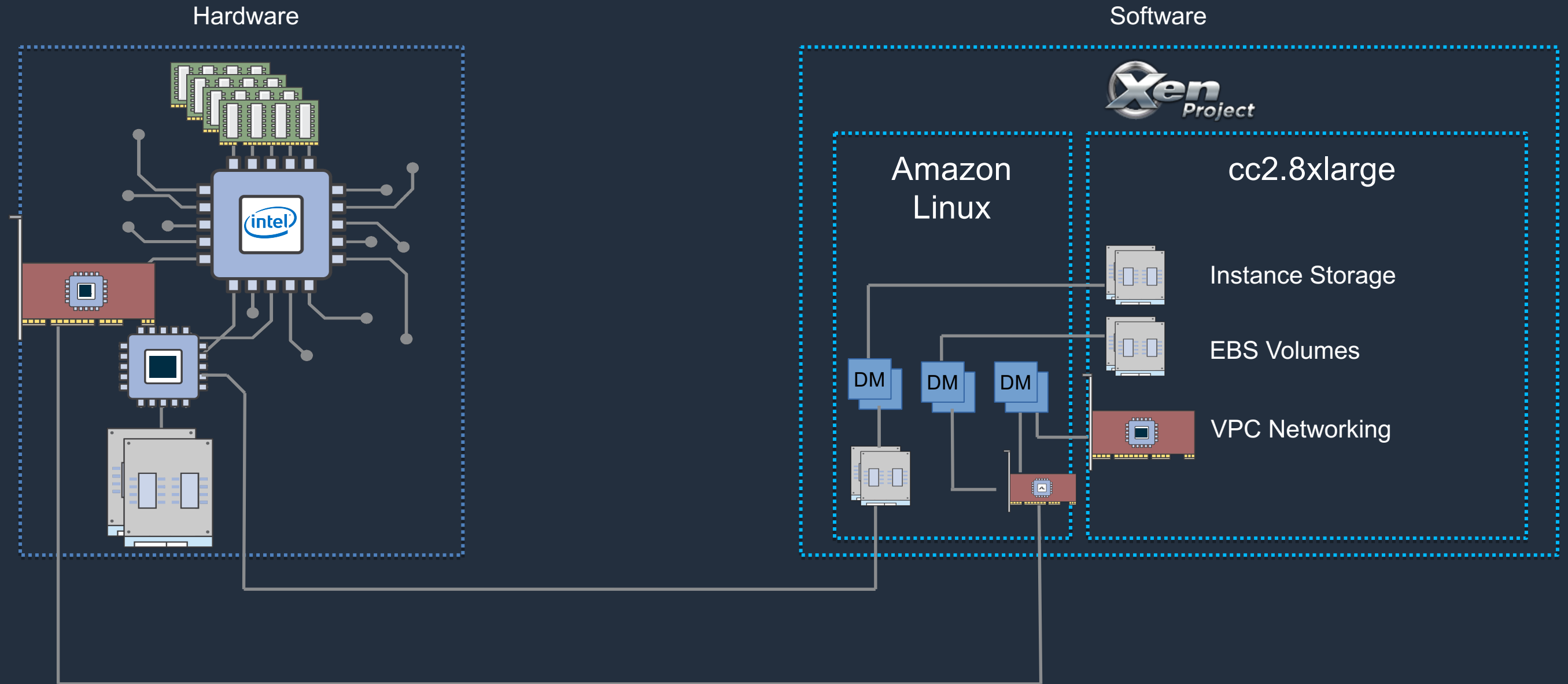


## Technology partners



## Consulting partners

aws

# Global Infrastructure

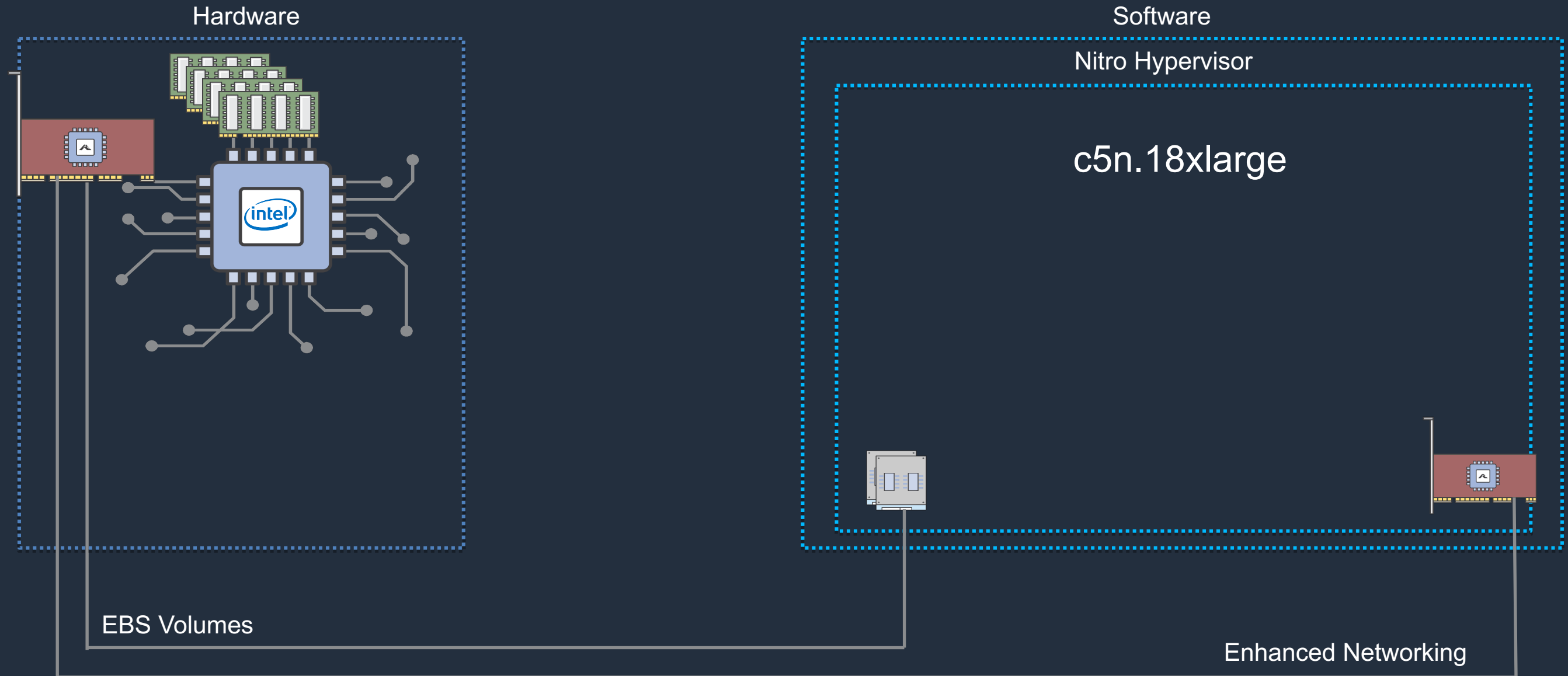We add the equivalent of an entire Fortune 500 company's compute capacity every day

Oregon
N. California
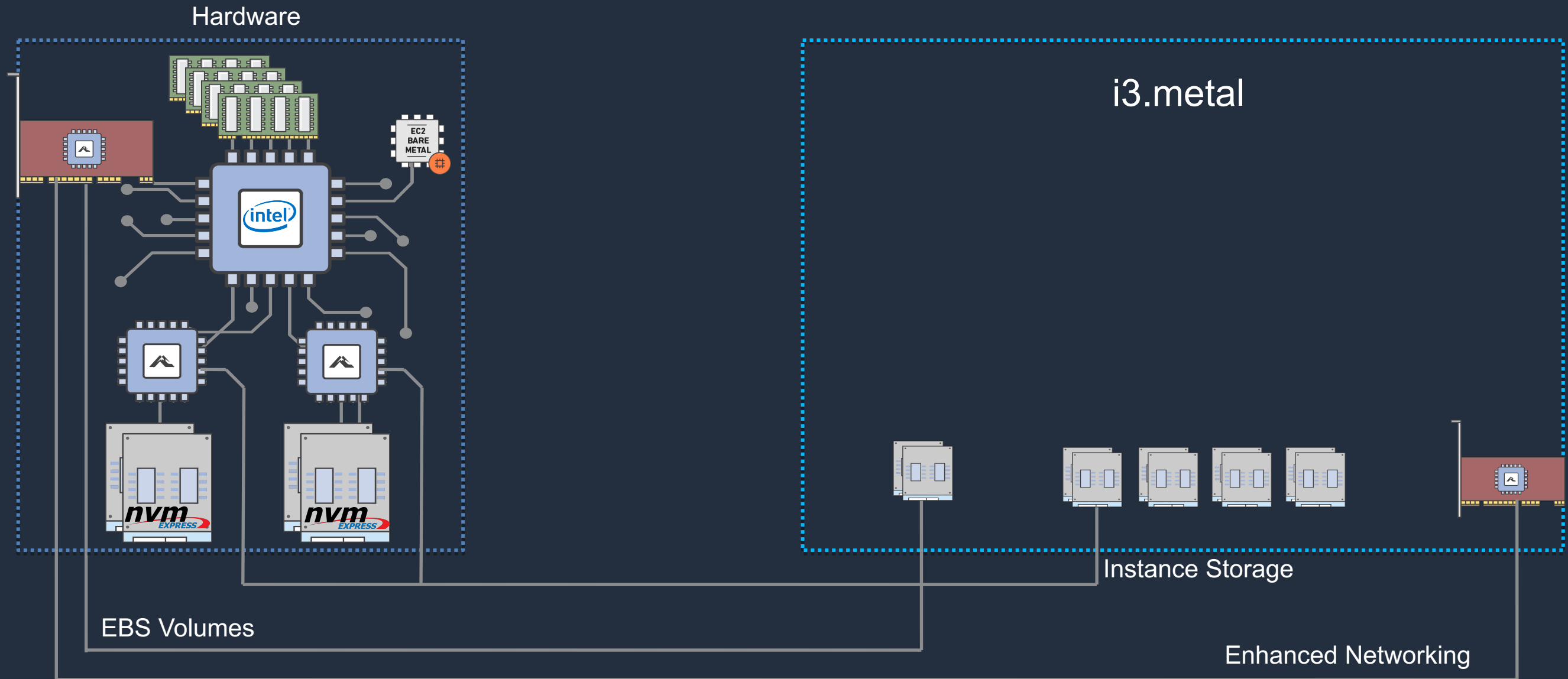GovCloud (US-West)
Ohio
GovCloud (US-East)
Canada (Central)
N. Virginia
Ireland
London
Stockholm
Frankfurt
Paris
Milan
Beijing
Seoul
Ningxia
Osaka
Tokyo
Hong Kong SAR
Bahrain
Mumbai
Singapore
Indonesia
São Paulo
Cape Town
Sydney

● **Coming soon**

77 Availability Zones
within 24 geographic
Regions around the world

# Evolution of Networking on AWS

aws

# Anatomy of an EC2 Instance – Circa 2011



Hardware

Software

Amazon Linux

cc2.8xlarge

Instance Storage

EBS Volumes

VPC Networking

DM   DM   DM

# What EC2 Instances Look Like Today



Hardware

Software

Nitro Hypervisor

c5n.18xlarge

EBS Volumes

Enhanced Networking

aws

# What EC2 Looks Like Today… when you don't want hypervisors



Hardware

EC2 BARE METAL

i3.metal

nvm EXPRESS

nvm EXPRESS

Instance Storage

EBS Volumes

Enhanced Networking

aws

# What AWS Looks Like Today... with EFA



Hardware

Software

c5n.18xlarge

NVMe

EFA

ENA

EBS Volumes

Enhanced Networking

aws

# EFA – For HPC and Machine Learning

Scale **tightly-coupled**
HPC applications on AWS

| C5n | P3dn | I3en | M5(d)n | R5(d)n | G4dn Metal* | C5a(d)n Metal* |

intel  nVIDIA  AMD

**\* Coming soon**

## EFA

Elastic Fabric Adapter, best for large HPC workloads

High data throughput

100 Gbps network bandwidth

Congestion control for cloud scale and rapid packet loss recovery.

Lower latency for message passing and more effective application-layer comms.
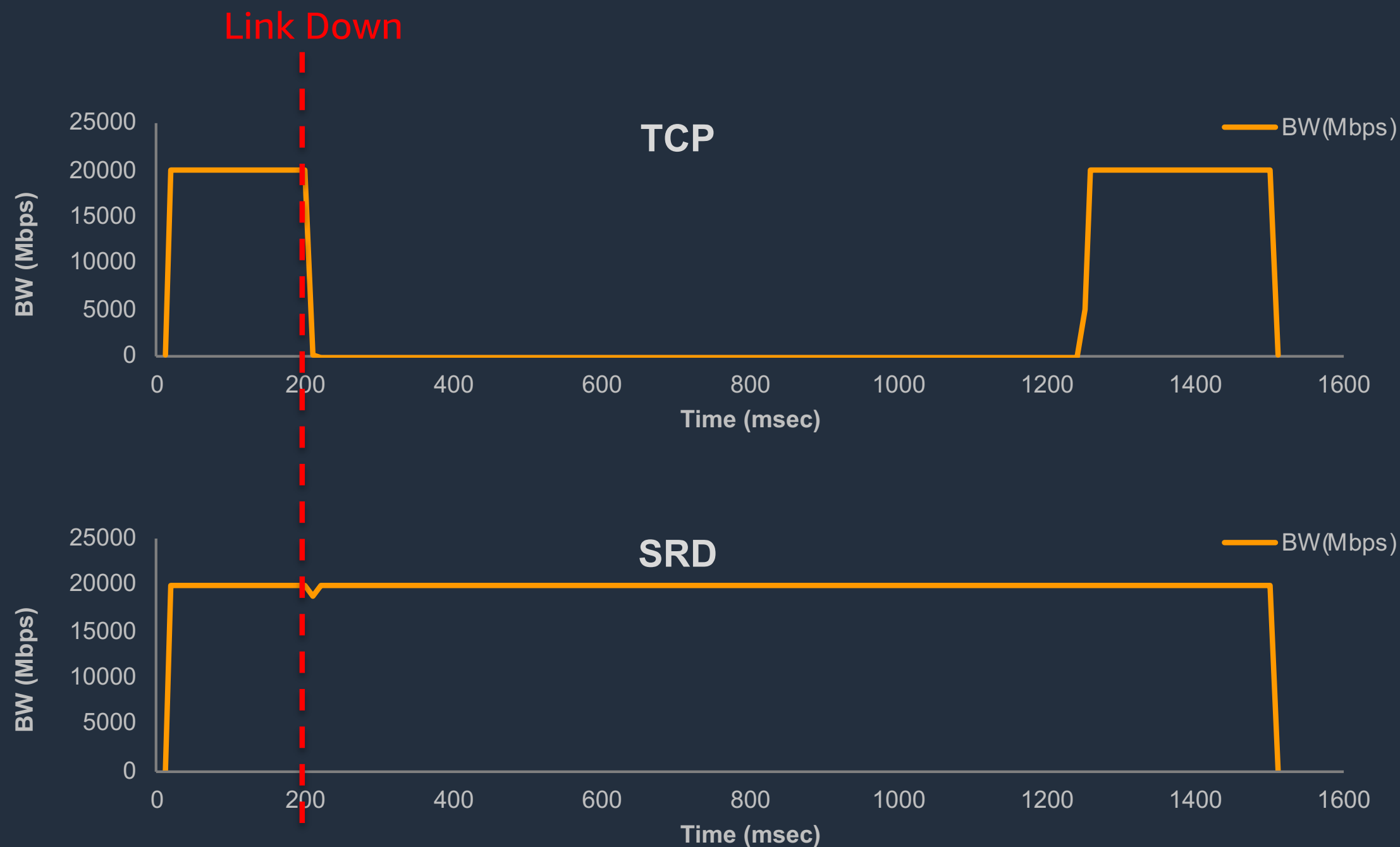
aws

# Scalable Reliable Datagrams (SRD)

aws

# Scalable Reliable Datagram Protocol
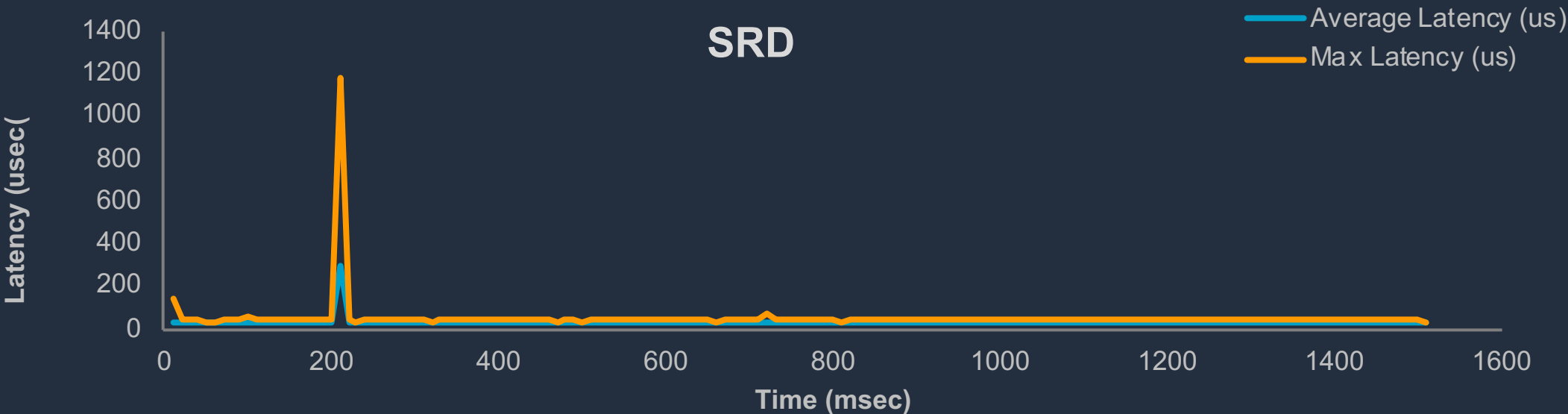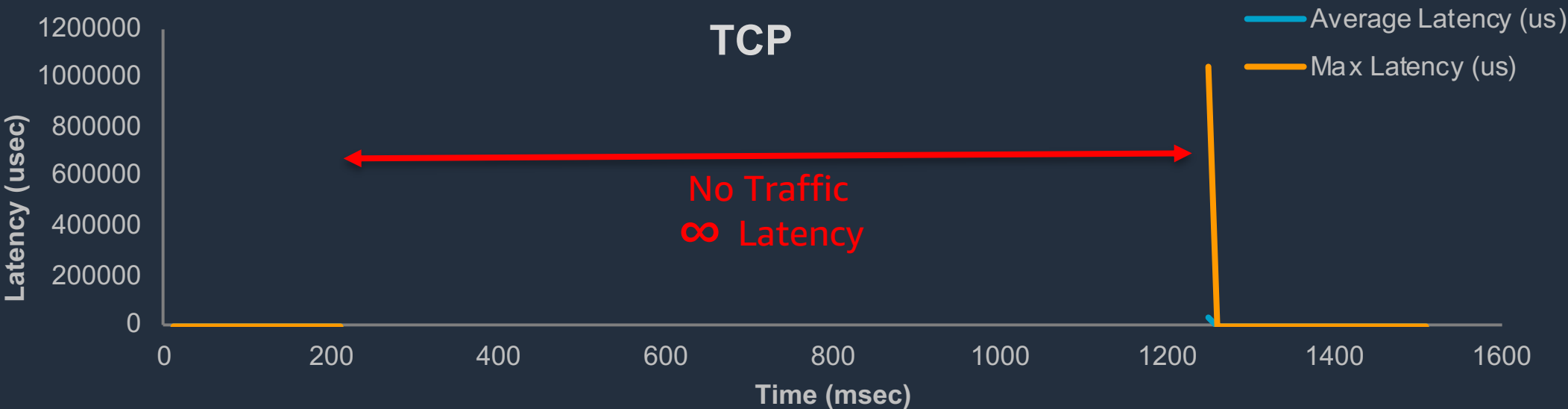


*Courtesy: Wikipedia*

- New protocol designed for AWS's unique datacenter network

- Implemented as part of our 3<sup>rd</sup> generation Nitro chip

- EFA exposes SRD as a reliable datagram interface

- Inspired by IB Reliable Datagram, without the drawbacks

- Packet spraying over multiple ECMP paths

- Out-of-order delivery – no head-of-line blocking

- Congestion and flow control designed for large-scale cloud

L. Shalev, H. Ayoub, N. Bshara and E. Sabbag, "Supercomputing on Nitro in AWS Cloud," in IEEE Micro 2020
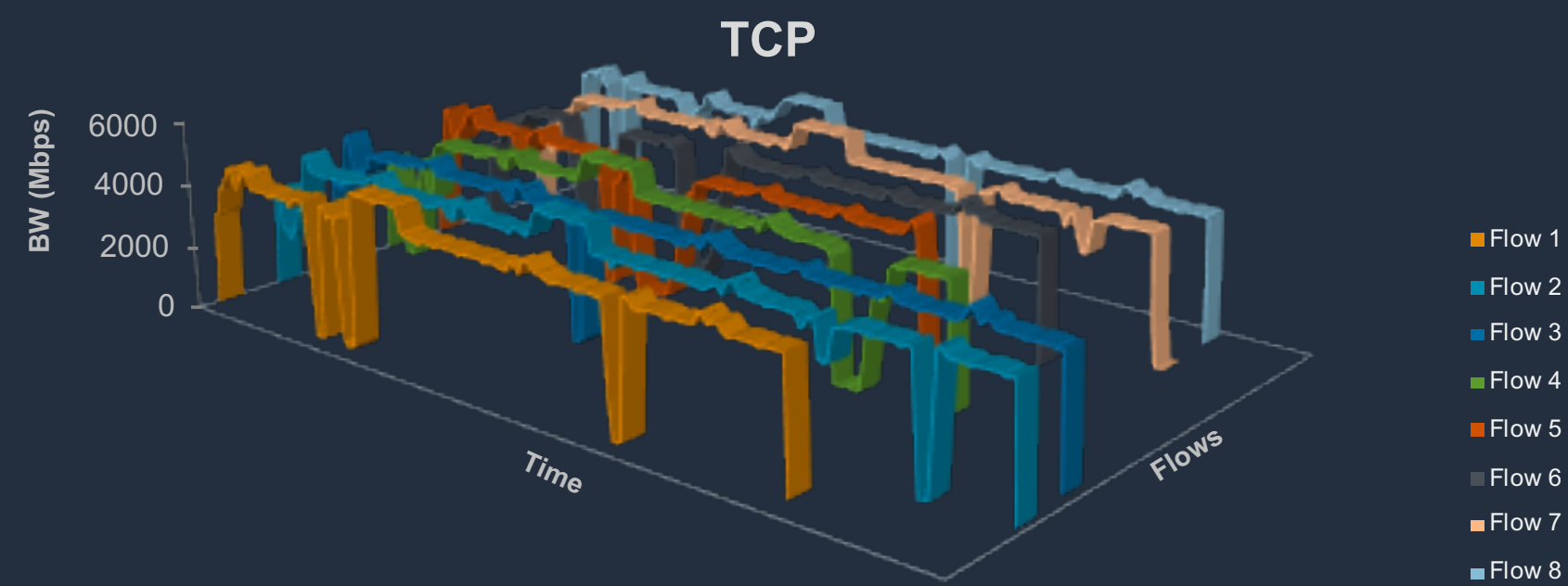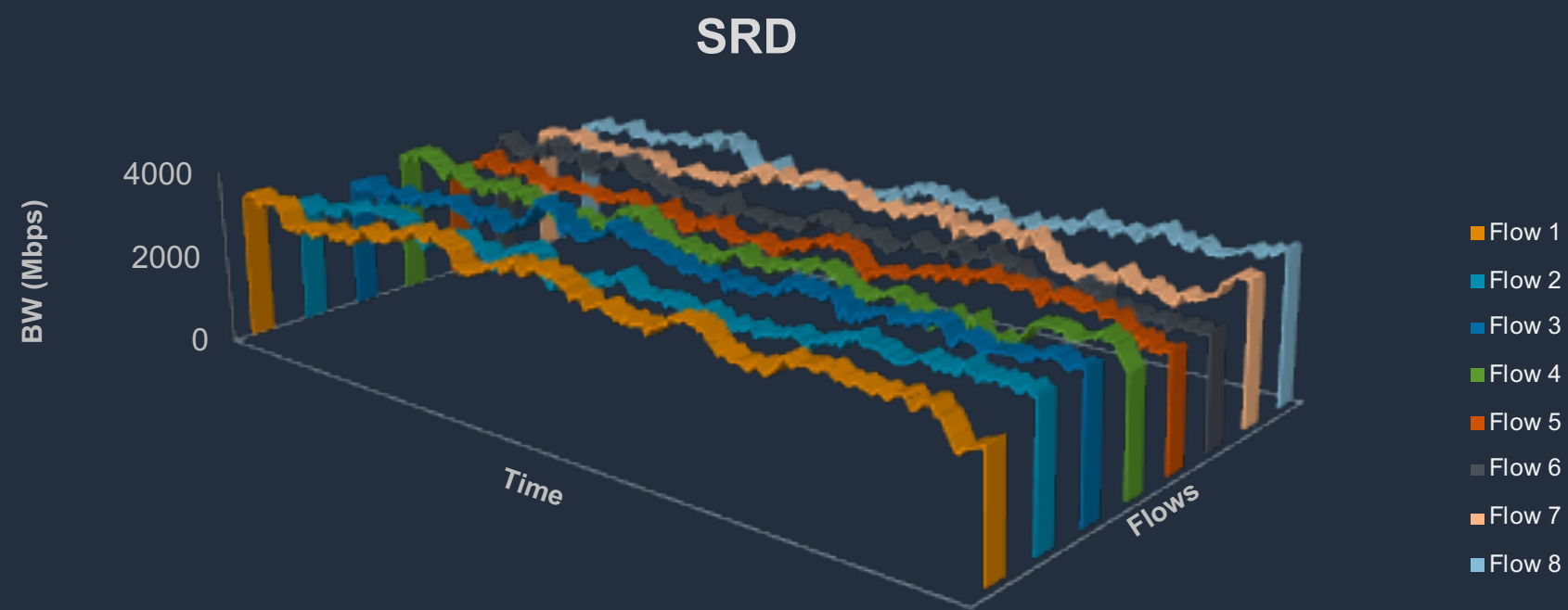
aws

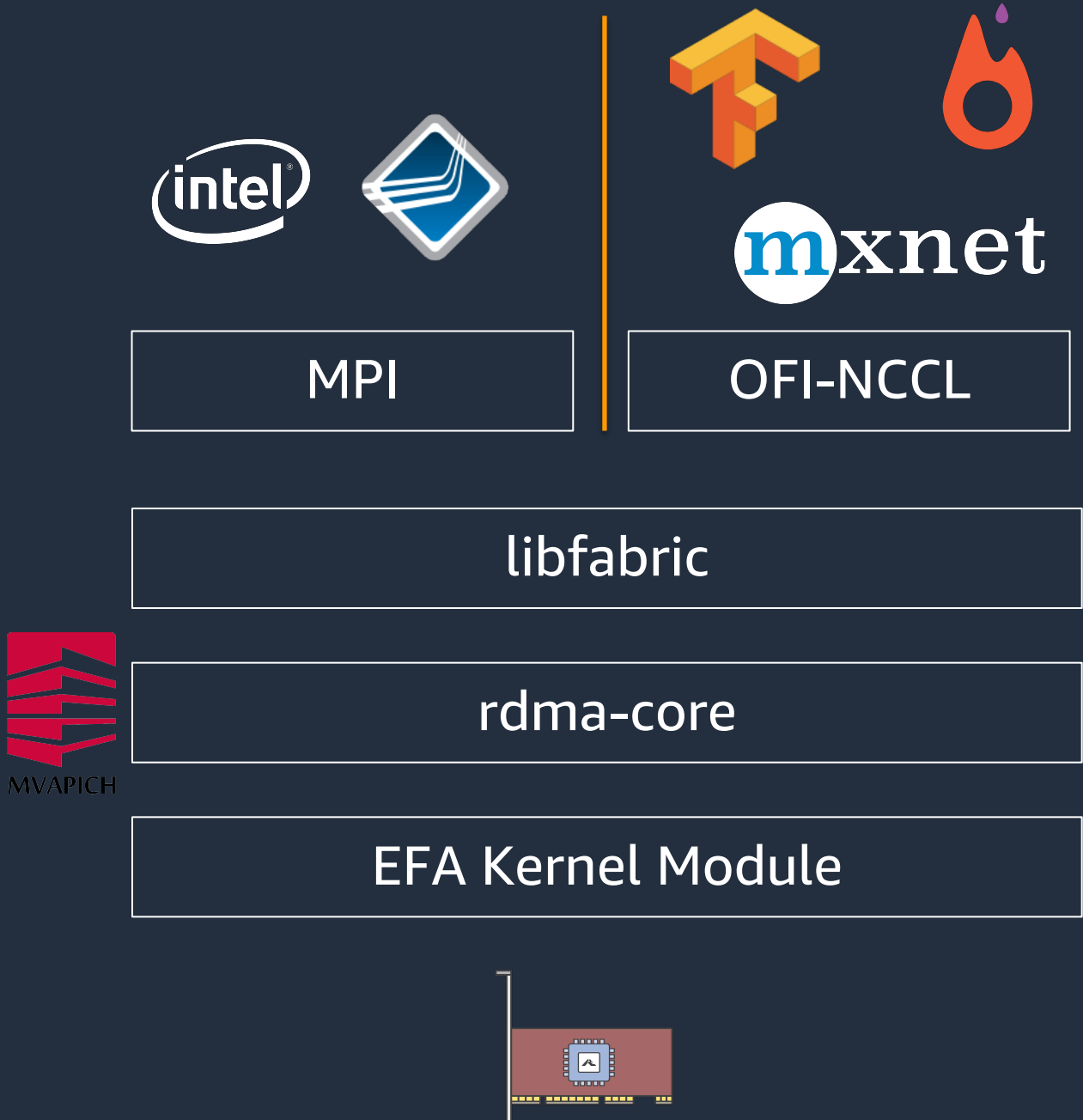# SRD Link Failure Handling - Throughput

# SRD Link Failure Handling - Latency
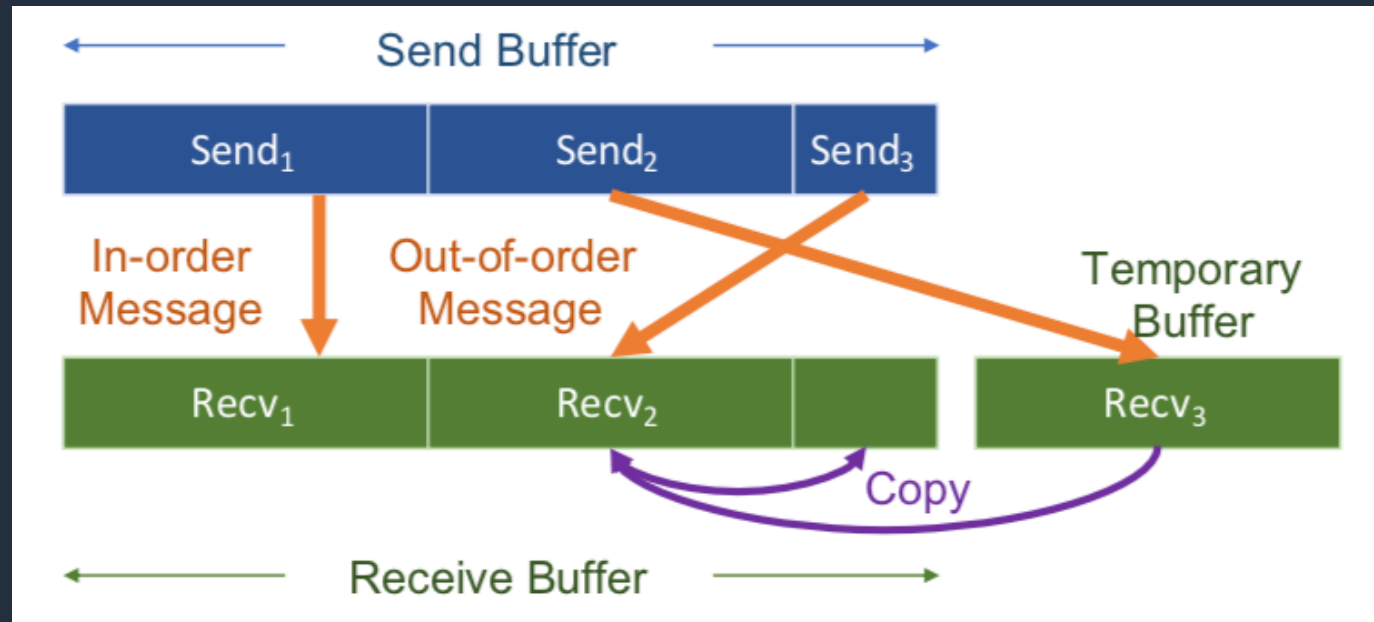
# SRD Congestion Control on Incast



© 2020, Amazon Web Services, Inc. or its Affiliates.

# EFA Software Ecosystem

aws

# EFA Software Ecosystem

MPI

OFI-NCCL

libfabric

rdma-core

MVAPICH
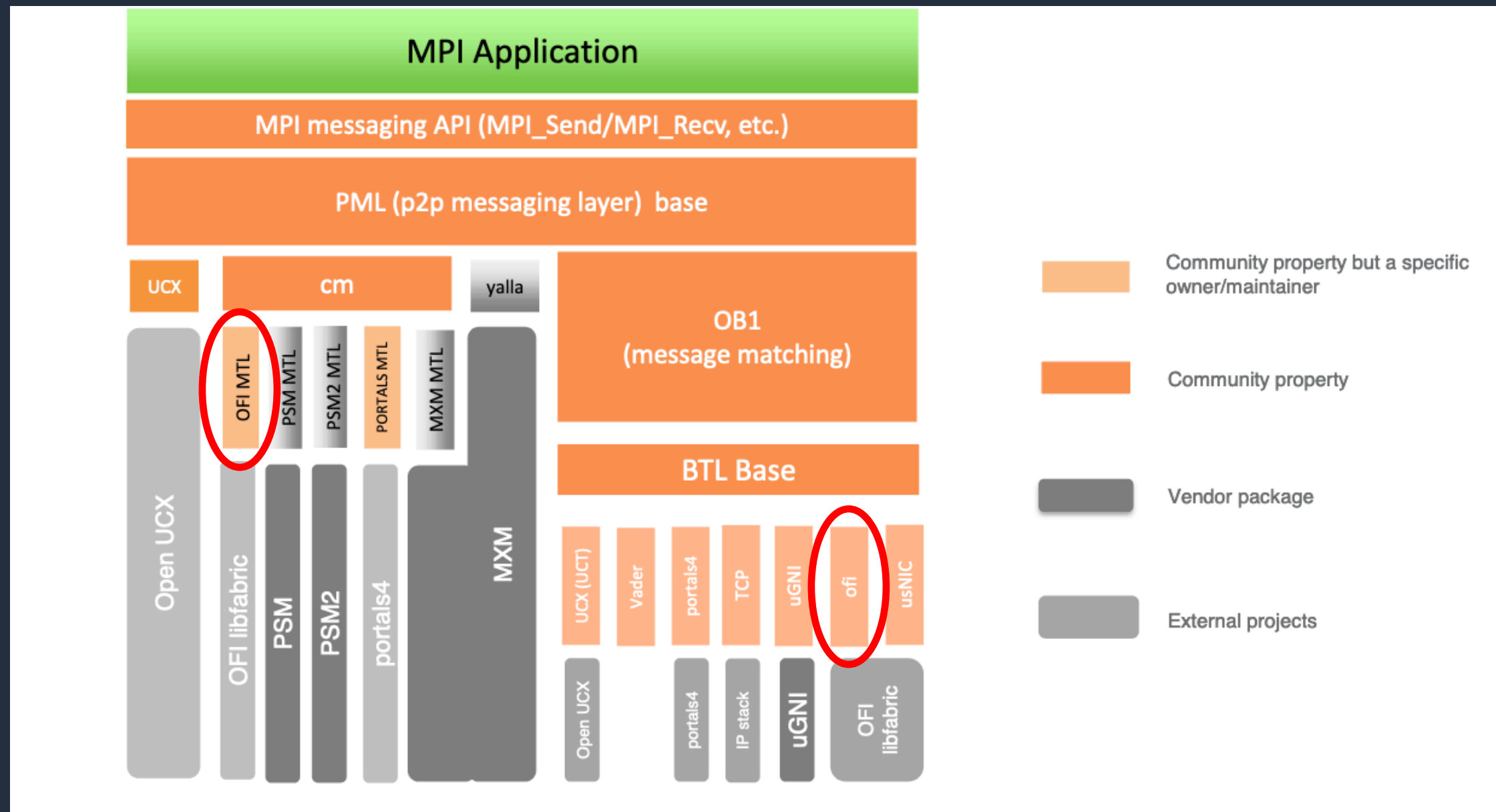
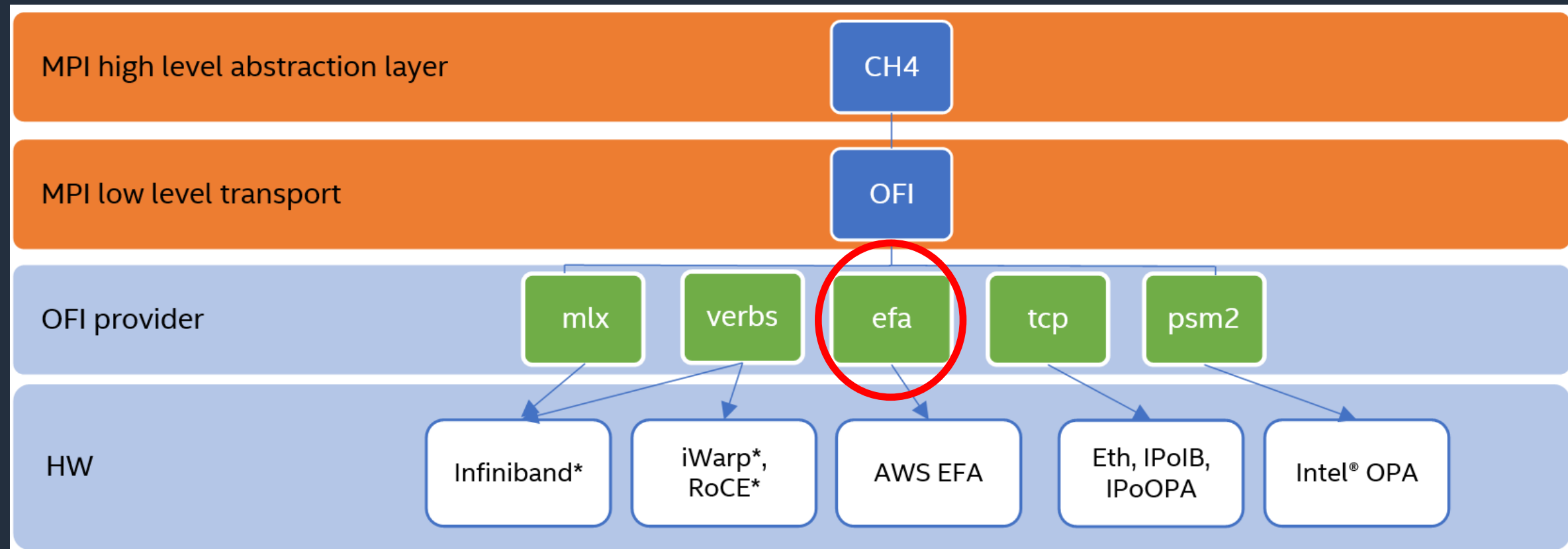EFA Kernel Module

aws

# MVAPICH2-X-AWS



- Directly programs to rdma-core
- Reordering with copy-out
- Use of immediate data for seq ID
- Long message packetization
- Tag matching
- Intra-node path with XPMEM

S. Chakraborty, S. Xu, H. Subramoni and D. K. Panda, Designing Scalable and High-Performance MPI Libraries on Amazon Elastic Adapter, Hot Interconnect, 2019

aws

# Open MPI and EFA

# Intel MPI and EFA



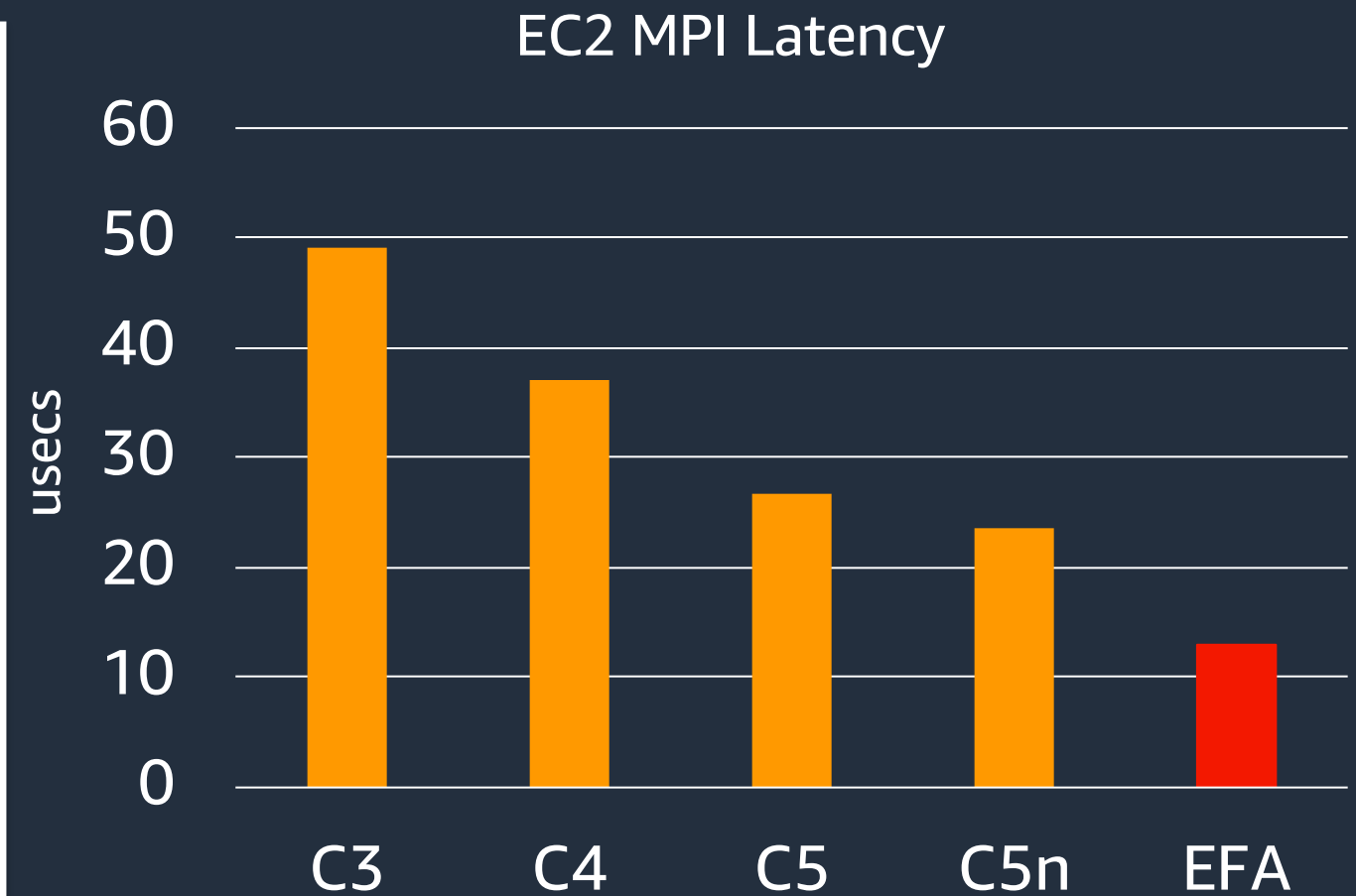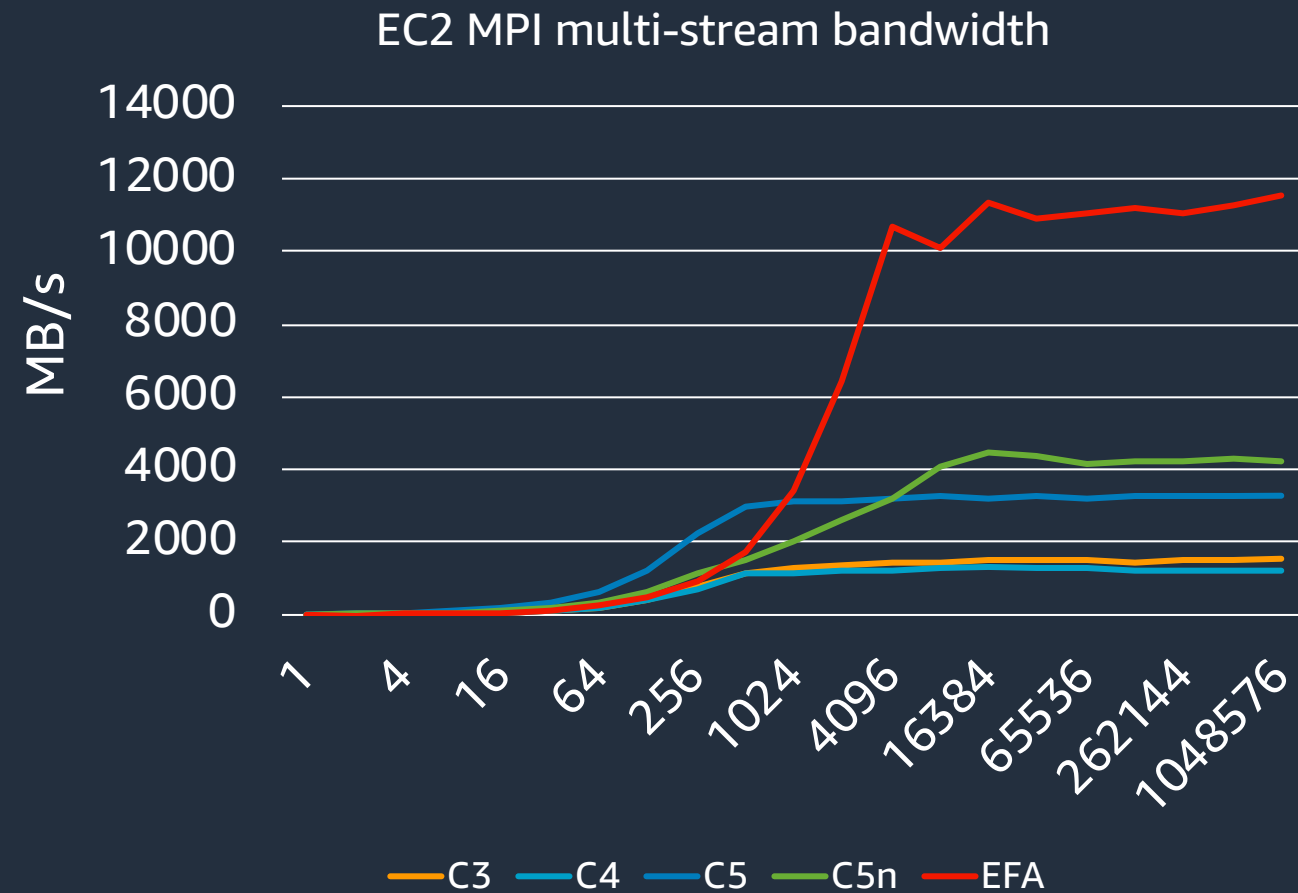https://software.intel.com/content/www/us/en/develop/articles/intel-mpi-library-2019-over-libfabric.html20

# EFA Installer

```
$ curl -O https://s3-us-west-2.amazonaws.com/aws-efa-installer/aws-efa-installer-latest.tar.gz
$ tar -xf aws-efa-installer-1.4.1.tar.gz
$ cd aws-efa-installer
$ sudo ./efa_installer.sh -y
= Starting Amazon Elastic Fabric Adapter Installation Script =
= EFA Installer Version: 1.4.1 =

== Installing EFA dependencies ==
<snip>
== Writing EFA profile.d configuration ==
== Configuring system limits for EFA ==
Limits for Elastic Fabric Adapter configured.
== Testing EFA device ==
Starting server...
Starting client...
<snip>
=========================================================
EFA installation complete.
- Please logout/login to complete the installation.
- Libfabric was installed in /opt/amazon/efa
- Open MPI was installed in /opt/amazon/openmpi
=========================================================
```
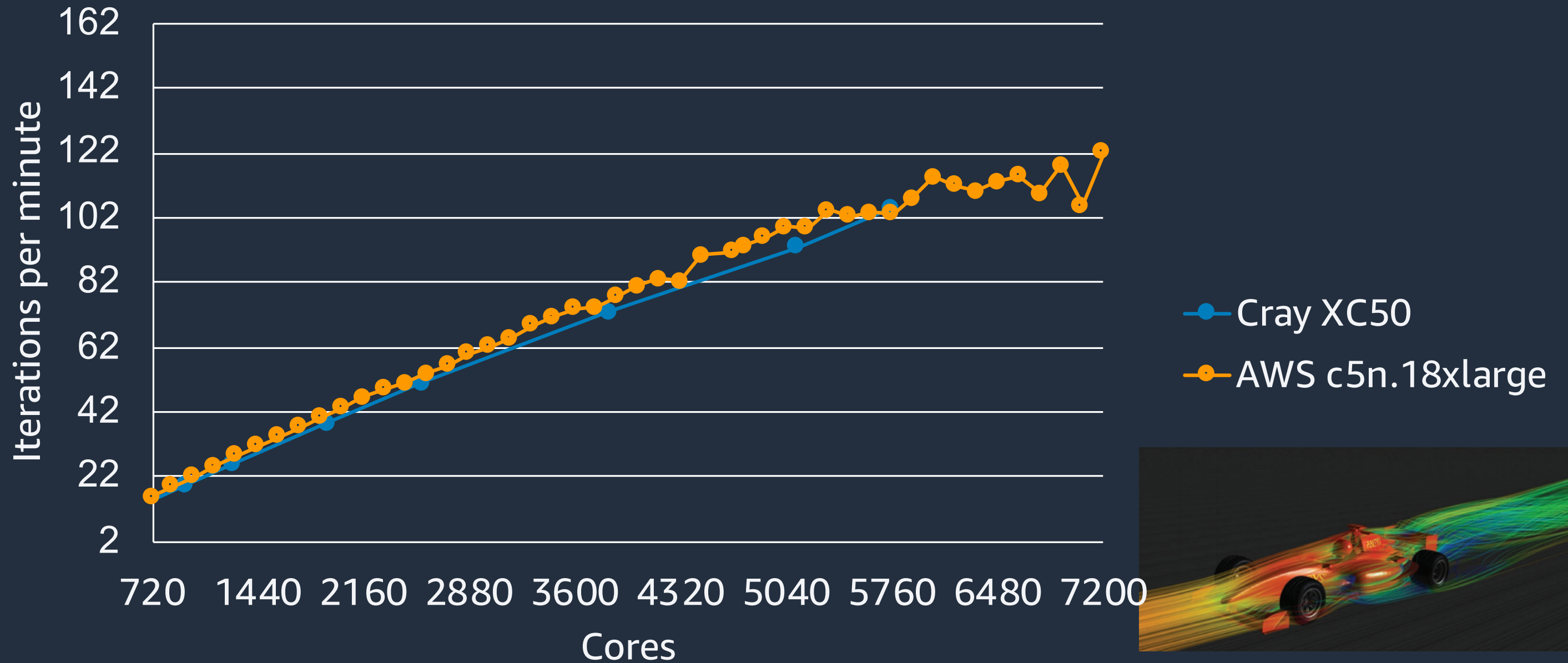
aws

# Performance Results

aws

# MPI Benchmarks



EC2 MPI multi-stream bandwidth

Legend: C3, C4, C5, C5n, EFA

EC2 MPI Latency

aws

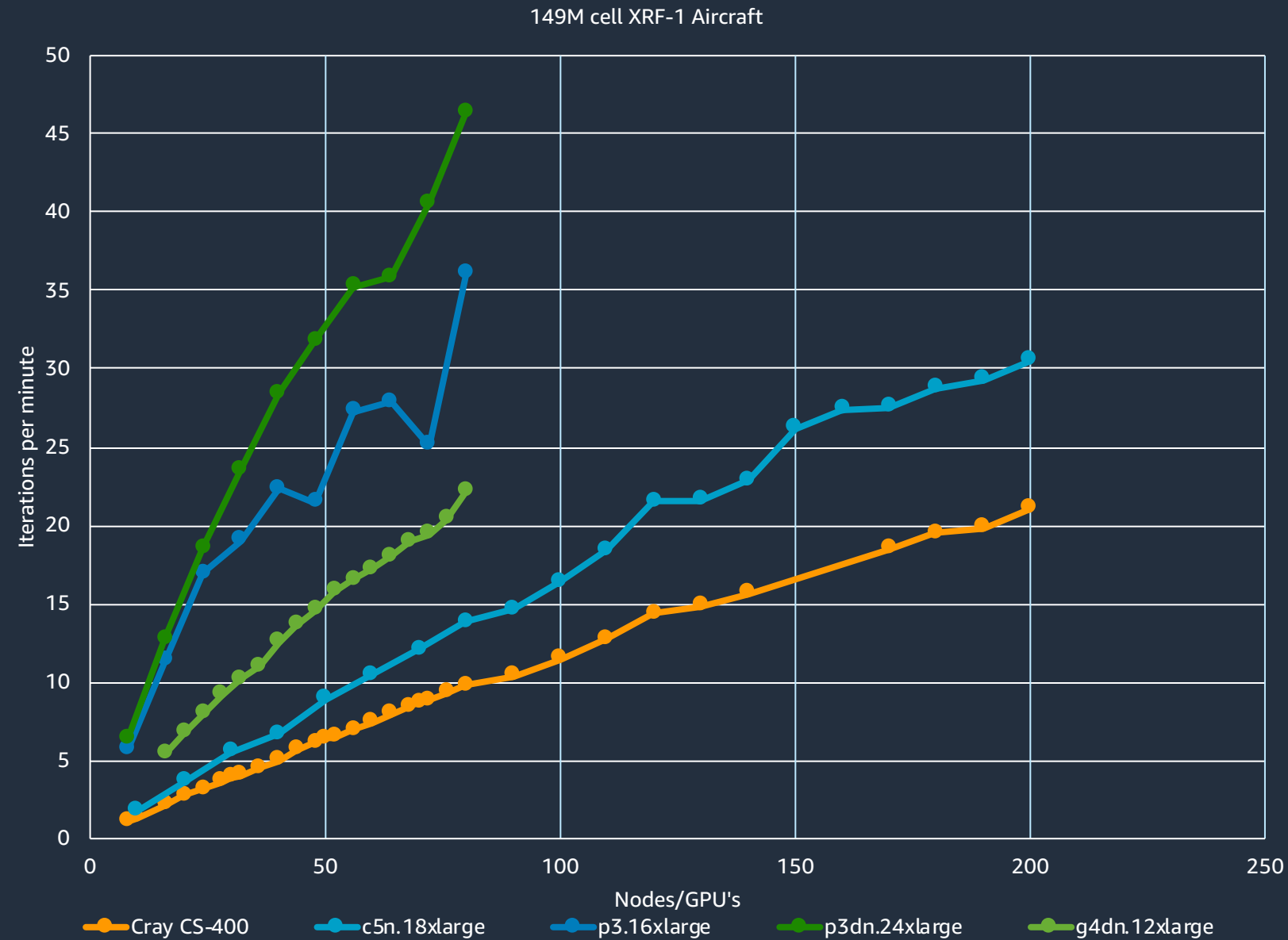# Scaling with Applications: Fluent

ANSYS Fluent 19.5 – F1 (**140M cells**) – IntelMPI 2019.5 – AL2 – PC2.5.1

# Scaling with Applications: zCFD

**149M cell XRF-1 Aircraft**



- Recent work with Zenotech (zCFD)

- GPUs (p3.24xlarge and g4dn.16xlarge Amazon EC2 instances i.e Nvidia v100 and T4's) deliver results faster for a lower cost
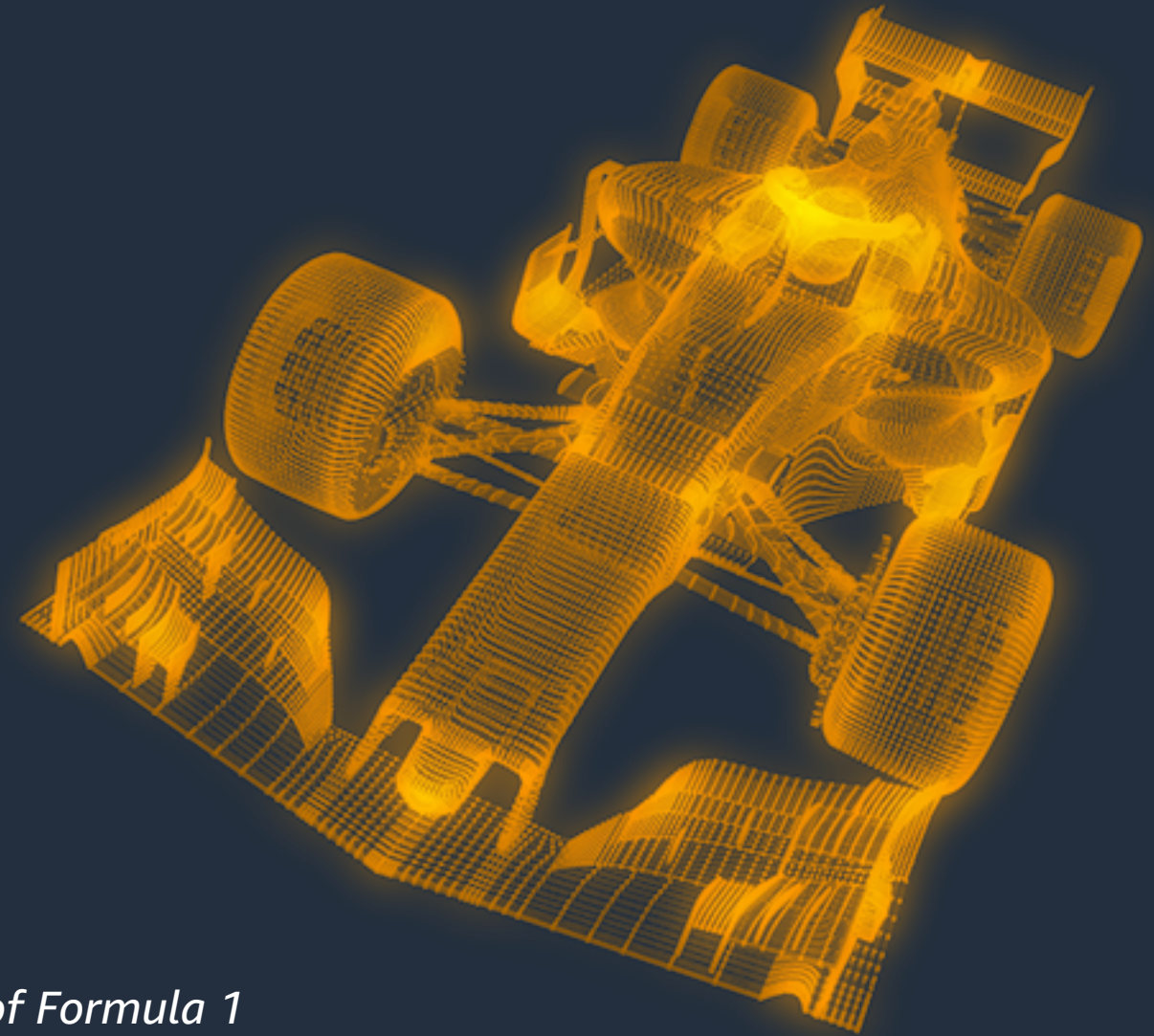
# Customer case study: Formula 1 on AWS

- No access to wind-tunnels for 12-24 months, only CFD

- **60hrs** to **10hrs** thanks to AWS.

- 192 cores to 1440 cores (EC2 c5n.18xlarge+ AWS ParallelCluster)

- On-demand + now smaller spot jobs (exploring other instance types)

*"This project with AWS was one of the most revolutionary in the history of Formula 1 aerodynamics,"* said Pat Symonds, Chief Technical Officer of Formula 1
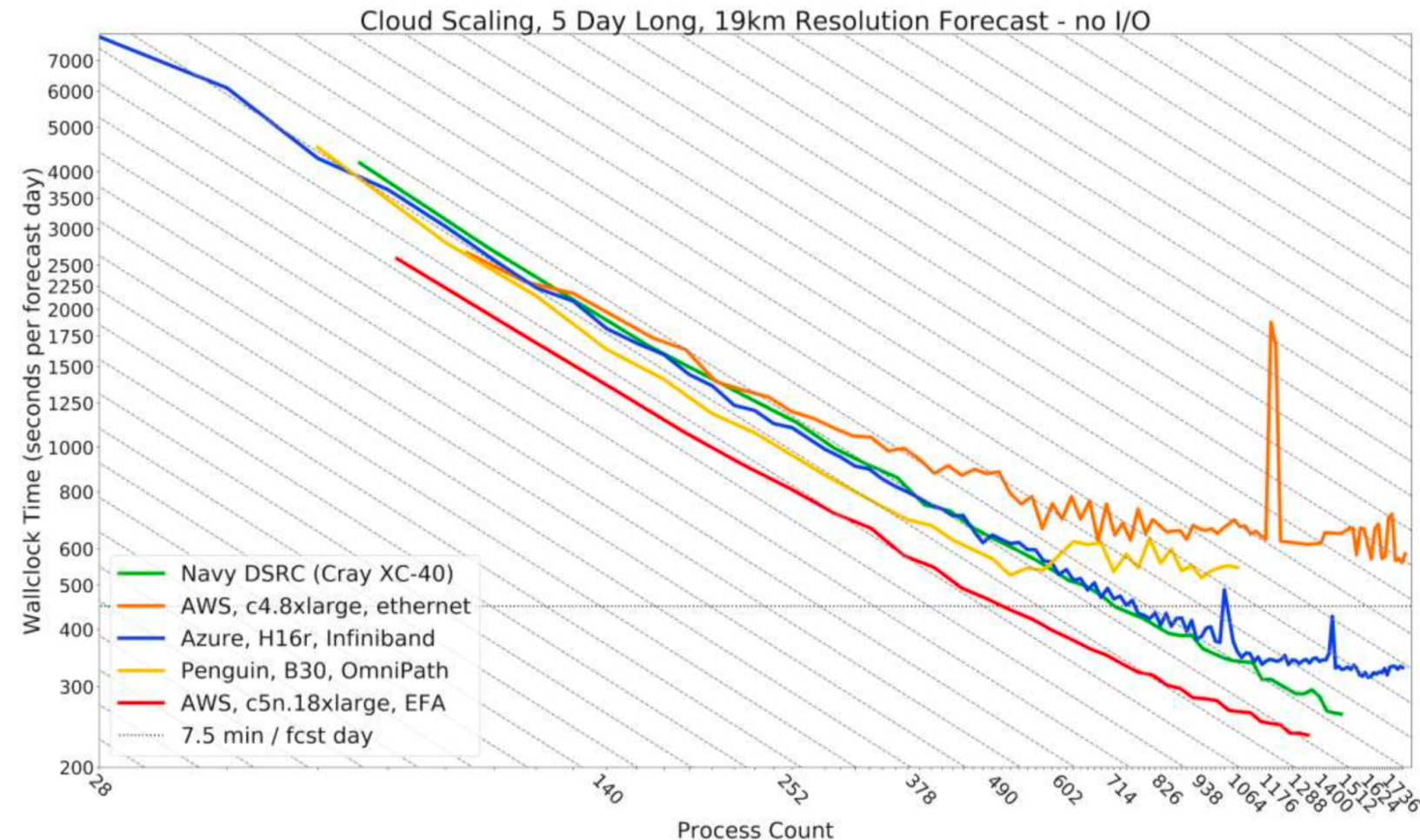
aws

# Scaling with Applications: NAVGEM



*Courtesy: U.S. Naval Research Laboratory*

https://www.youtube.com/watch?v=GTHWf0OVGrw&t=1177s

# Thank you!

Raghu Raja (craghun@amazon.com)

🐦 @rrcsraghu

aws