



15th ANNUAL WORKSHOP 2019

AMAZON ELASTIC FABRIC ADAPTER: ANATOMY, CAPABILITIES, AND THE ROAD AHEAD

Raghu Raja, Sr. SDE
Amazon Web Services

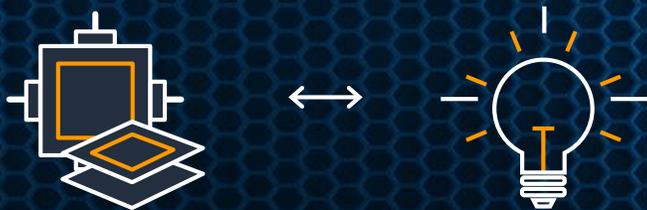


AGENDA

- **Overview of high Performance Computing on AWS**
- **What is EFA?**
- **Deep-dive on EFA**
- **Next steps**



HIGH PERFORMANCE COMPUTING ON AWS



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

BROAD HPC PARTNER COMMUNITY

Application partners



Infrastructure partners



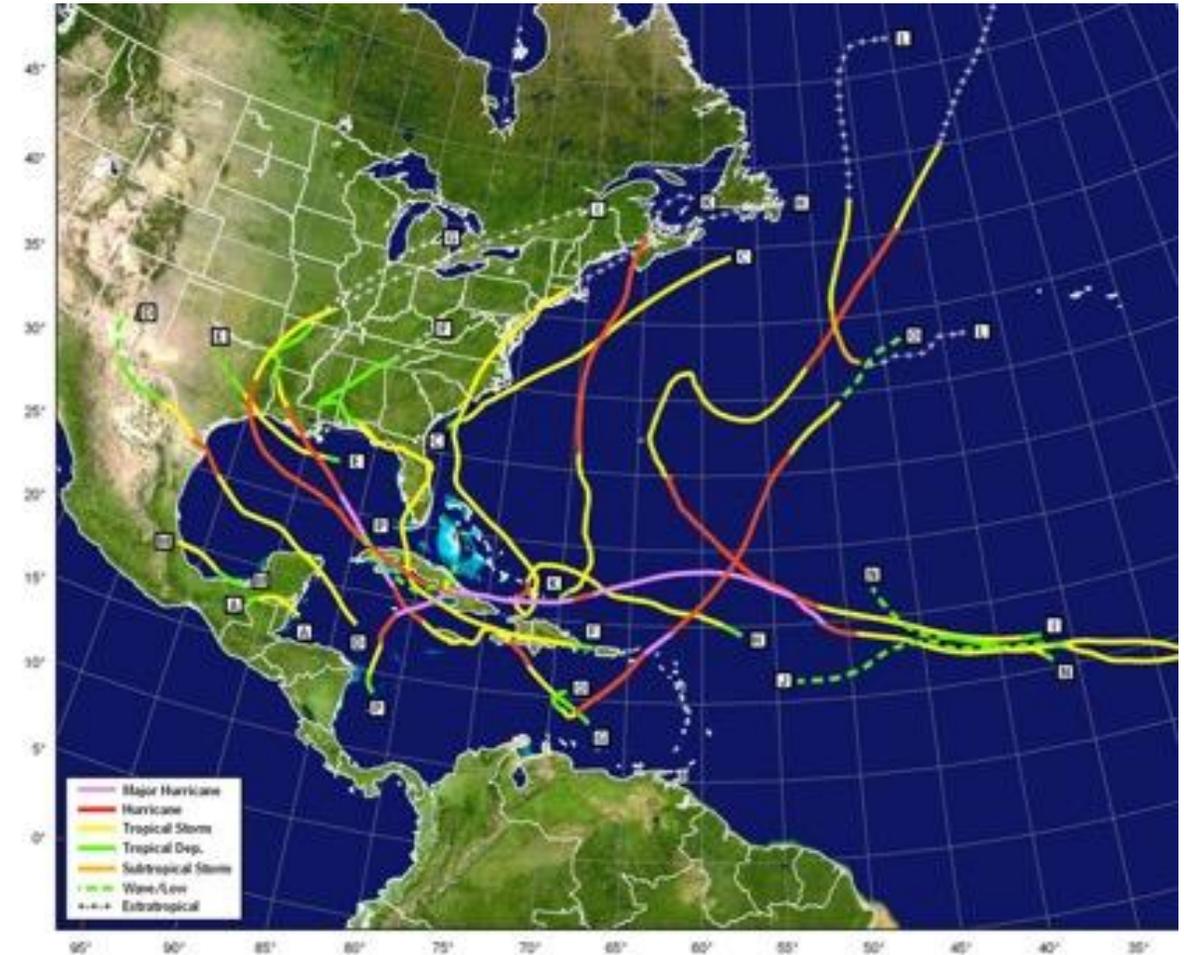
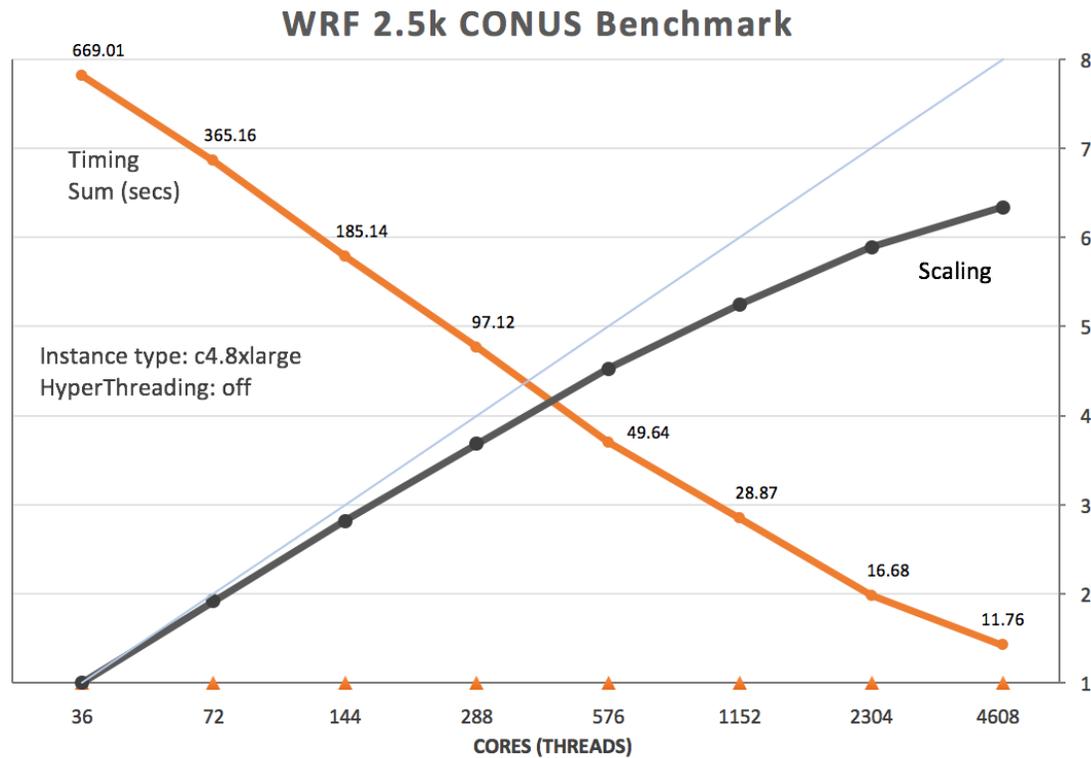
Technology partners



Consulting partners



HPC ON AWS: WEATHER MODELING



HPC ON AWS: DESIGN AND ENGINEERING

- **Boom leverages Rescale and AWS to enable supersonic travel**

- Simulated vortex lift with 200M cell models on 512+ cores
- Increased simulation throughput: 100 jobs in parallel with 6x speedup per job → 600x speedup
- Eliminated IT overhead, including server capital costs & in-house IT and software costs
- Elastic HPC capacity and pay-as-you-go AWS clusters allow business agility & ability to scale

- **“Rescale’s ScaleX cloud platform is a game-changer for engineering. It gives Boom computing resources comparable to building a large on-premise HPC center. Rescale lets us move fast with minimal capital spending and resources overhead.”**

- Josh Krall
- CTO & Co-Founder



HPC ON AWS: MATERIAL SCIENCE

The Western Digital logo, featuring the company name in white, bold, sans-serif font on a black rectangular background.

Over 2.3 million simulation jobs on a **single HPC cluster of 1 million** vCPUs—built using Amazon EC2 Spot Instances.

Time to results: **20** Days → **8** hours

“Storage technology is amazingly complex and we’re constantly pushing the limits of physics and engineering to deliver next-generation capacities and technical innovation. This successful collaboration with AWS shows the extreme scale, power and agility of cloud-based HPC to help us run complex simulations for future storage architecture analysis and materials science explorations. Using AWS to easily shrink simulation time from 20 days to 8 hours allows Western Digital R&D teams to explore new designs and innovations at a pace unimaginable just a short time ago.” – Steve Phillpott, CIO, Western Digital

SAVING KOALAS: GENOME SEQUENCING

Complete sequencing of
3.24 billion base pairs

3 million core-hours of
Amazon EC2 capacity



Australian Museum Research Institute

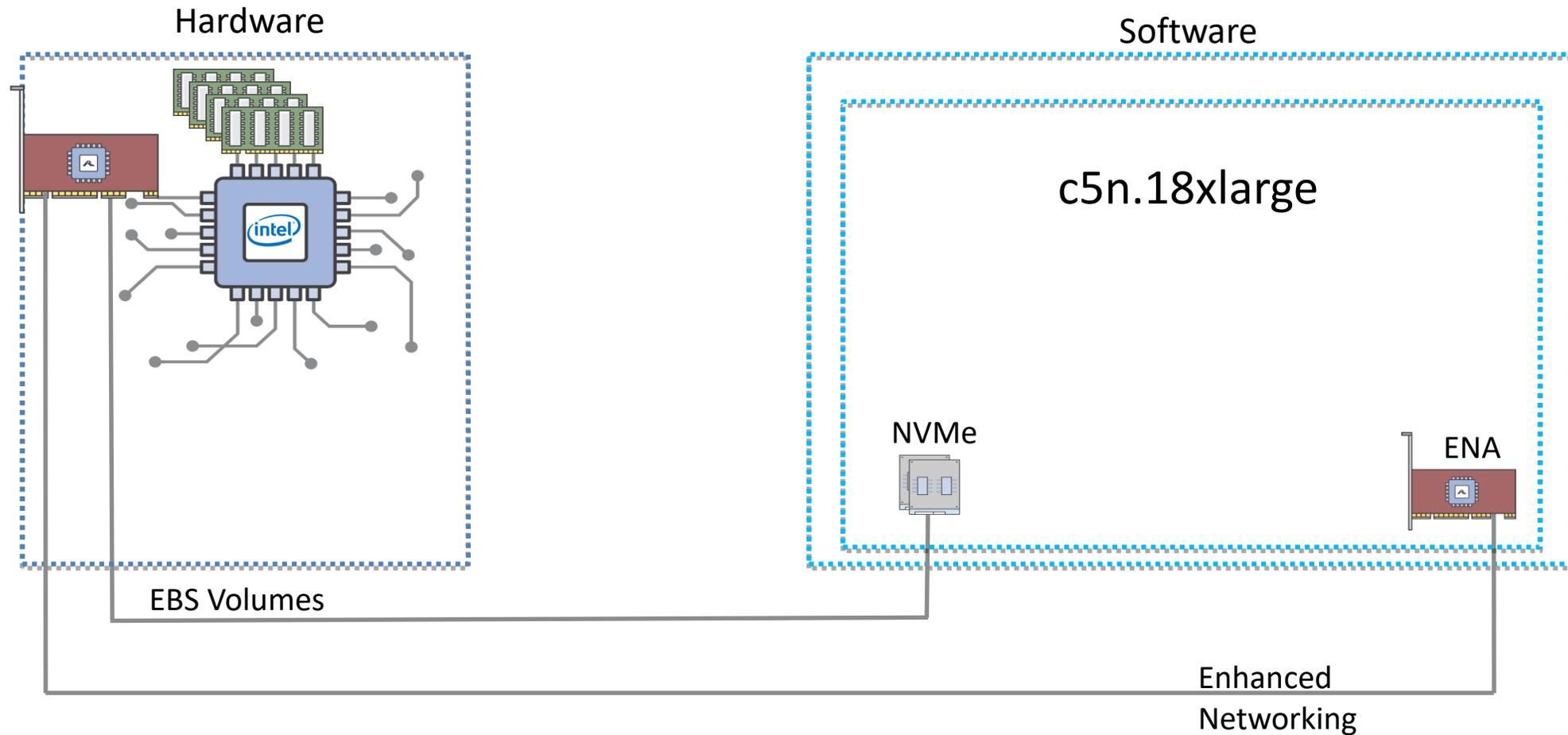
<https://www.nature.com/articles/s41588-018-0153-5>

<https://aws.amazon.com/blogs/aws/saving-koalas-using-genomics-research-and-cloud-computing/>

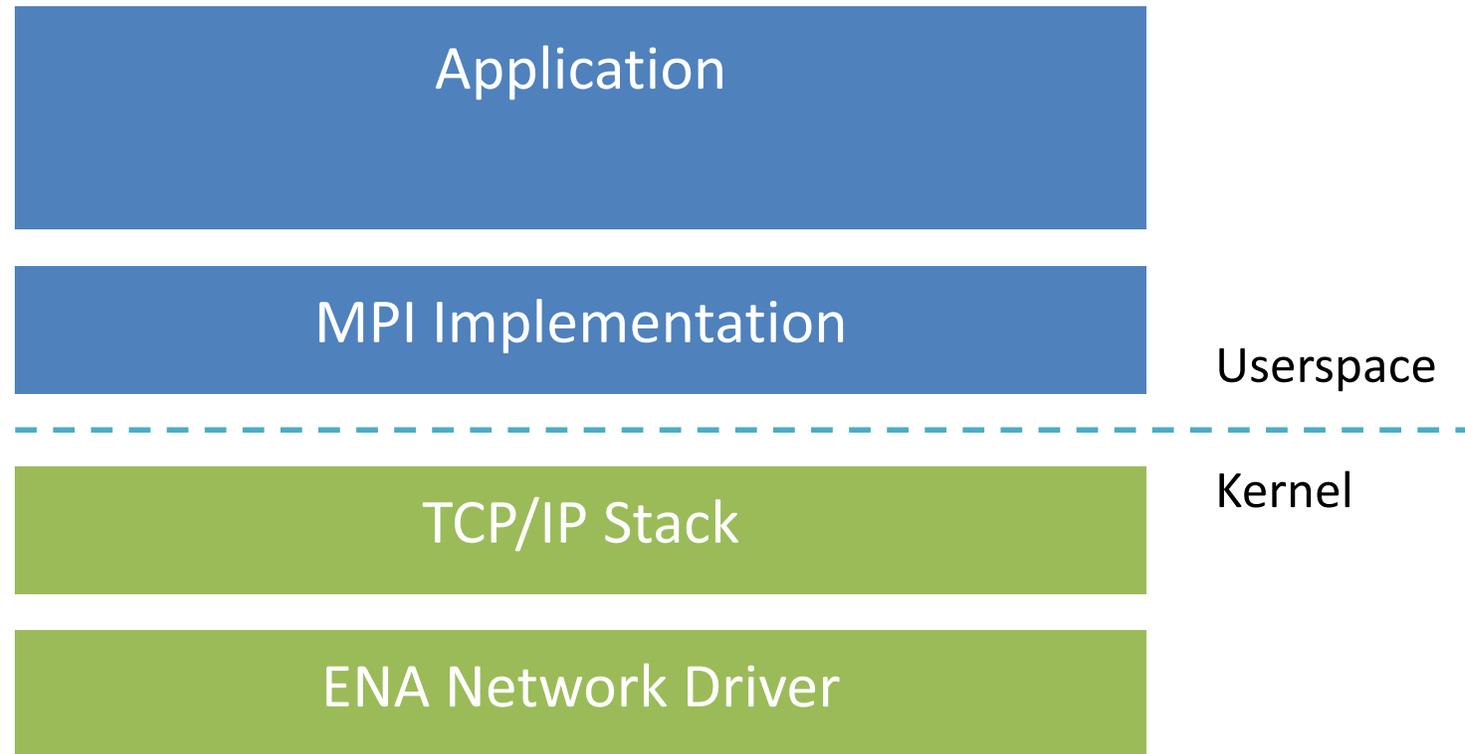


WHAT IS ELASTIC FABRIC ADAPTER (EFA)?

AMAZON ELASTIC COMPUTE CLOUD (EC2): 101

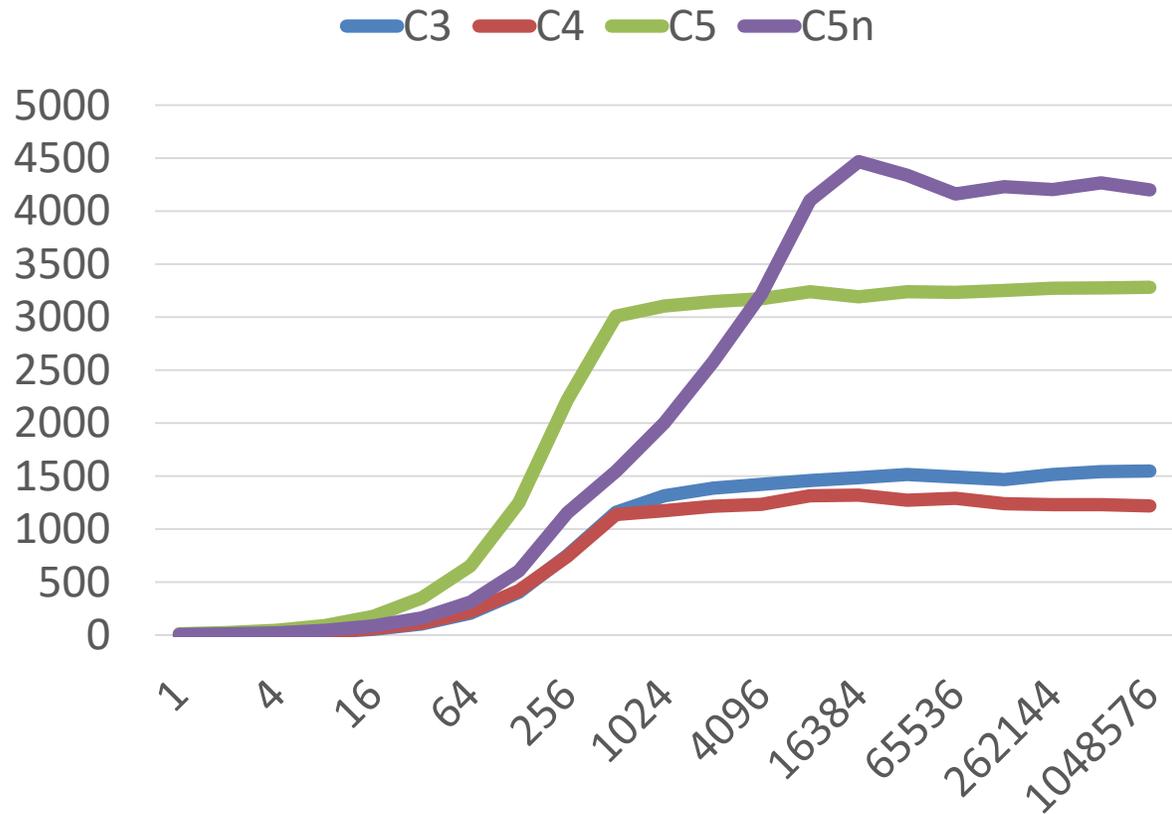


HPC SOFTWARE STACK ON EC2

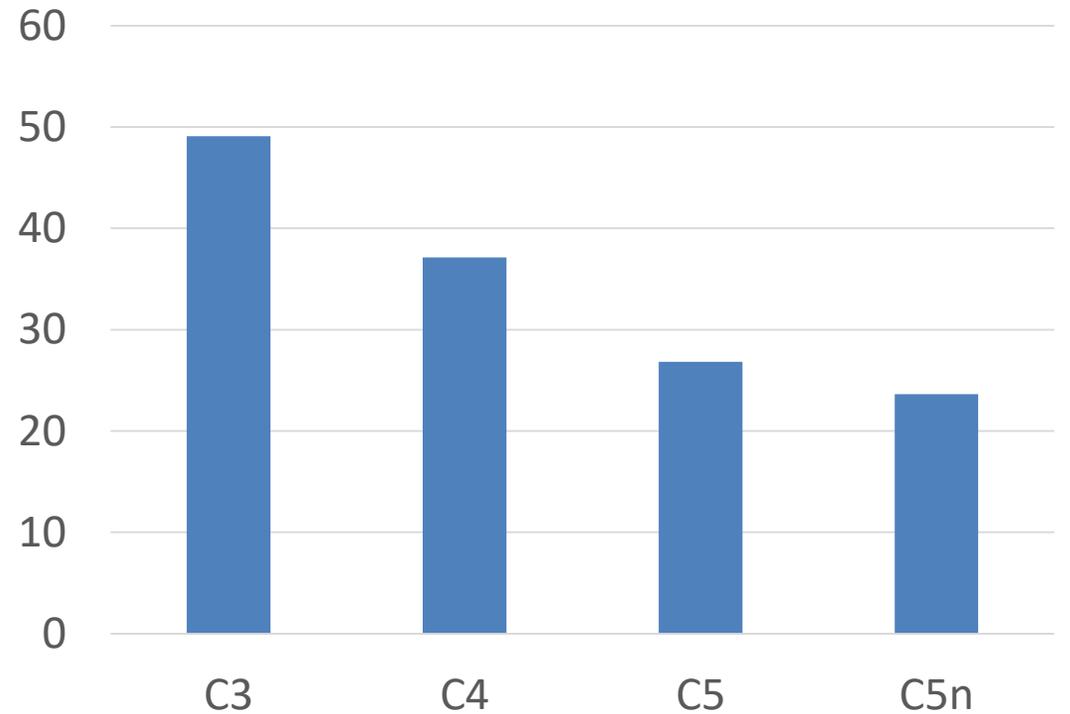


HPC NETWORK PERFORMANCE

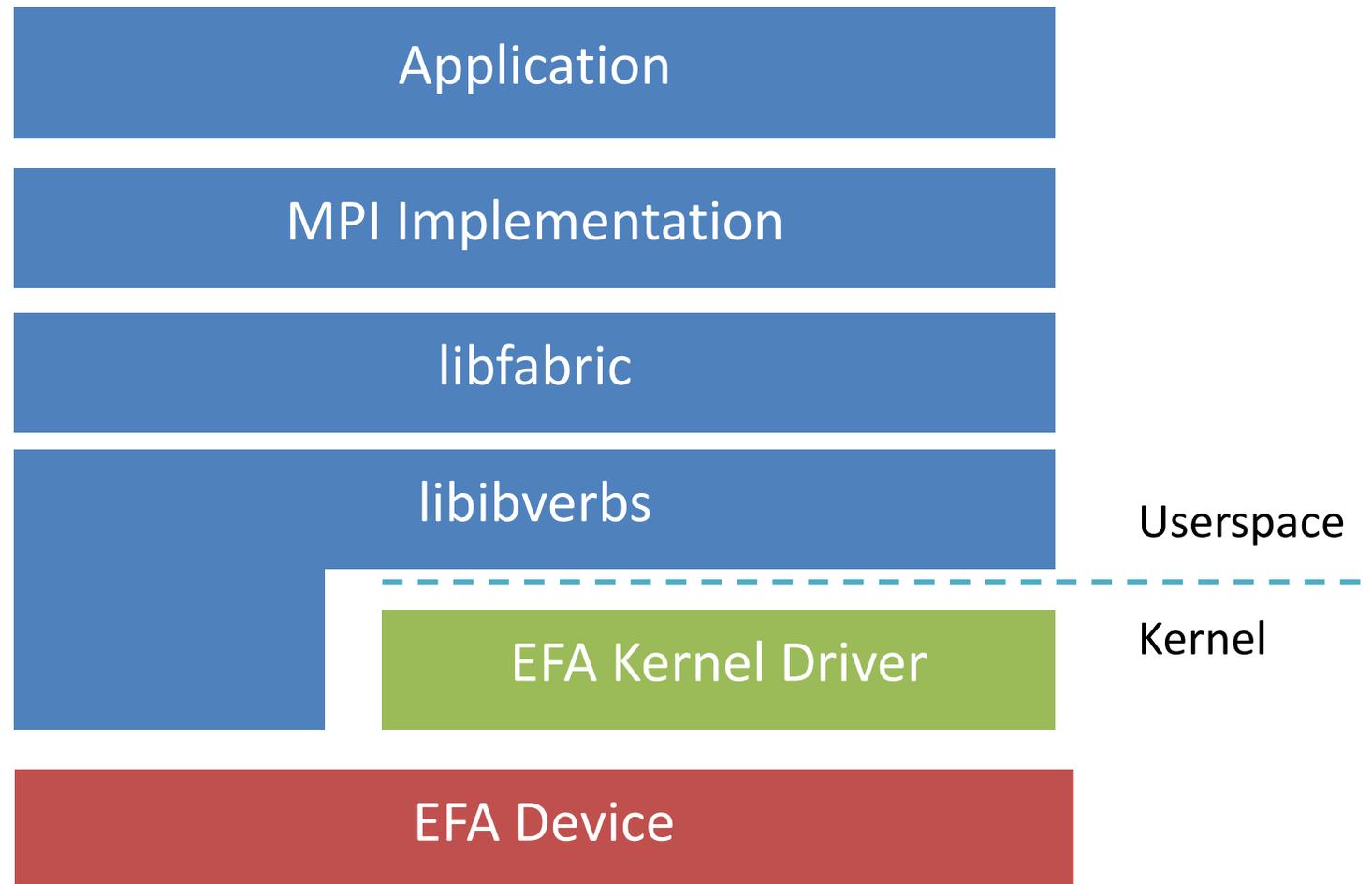
EC2 MPI multi-stream bandwidth



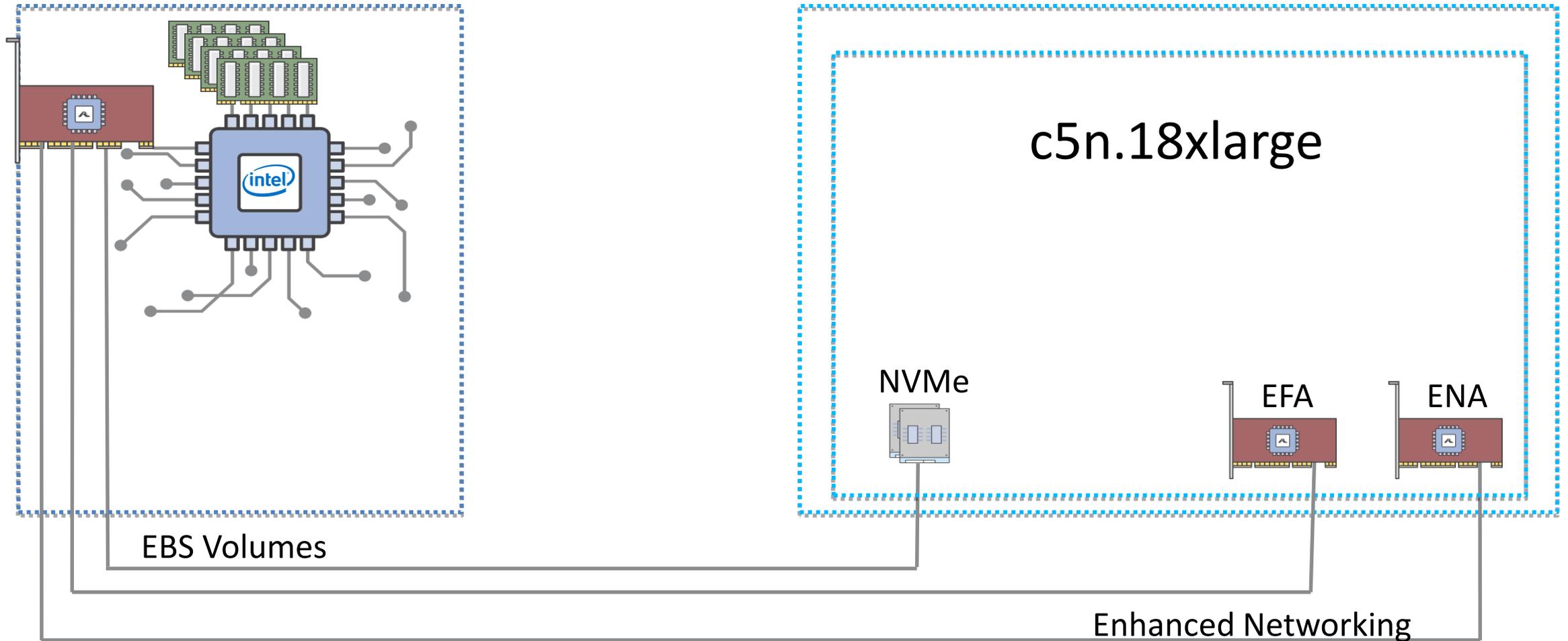
EC2 MPI Latency



HPC SOFTWARE STACK WITH EFA



EFA DEVICE



EFA DEVICE

00:05.0 Ethernet controller: Amazon.com, Inc. Elastic Network Adapter (ENA)

Subsystem: Amazon.com, Inc. Elastic Network Adapter (ENA)

Physical Slot: 5

Flags: bus master, fast devsel, latency 0

Memory at fe814000 (32-bit, non-prefetchable) [size=16K]

Memory at f8400000 (32-bit, prefetchable) [size=4M]

Capabilities: [70] Express Endpoint, MSI 00

Capabilities: [b0] MSI-X: Enable+ Count=33 Masked-

Kernel driver in use: ena

Kernel modules: ena

00:06.0 Ethernet controller: Amazon.com, Inc. Elastic Fabric Adapter (EFA)

Subsystem: Amazon.com, Inc. Elastic Fabric Adapter (EFA)

Physical Slot: 6

Flags: bus master, fast devsel, latency 0

Memory at fe818000 (32-bit, non-prefetchable) [size=16K]

Memory at f0000000 (64-bit, prefetchable) [size=128M]

Memory at fe000000 (64-bit, non-prefetchable) [size=8M]

Capabilities: [70] Express Endpoint, MSI 00

Capabilities: [b0] MSI-X: Enable+ Count=129 Masked-

Kernel driver in use: efa

SCALABLE RELIABLE DATAGRAM (SRD)

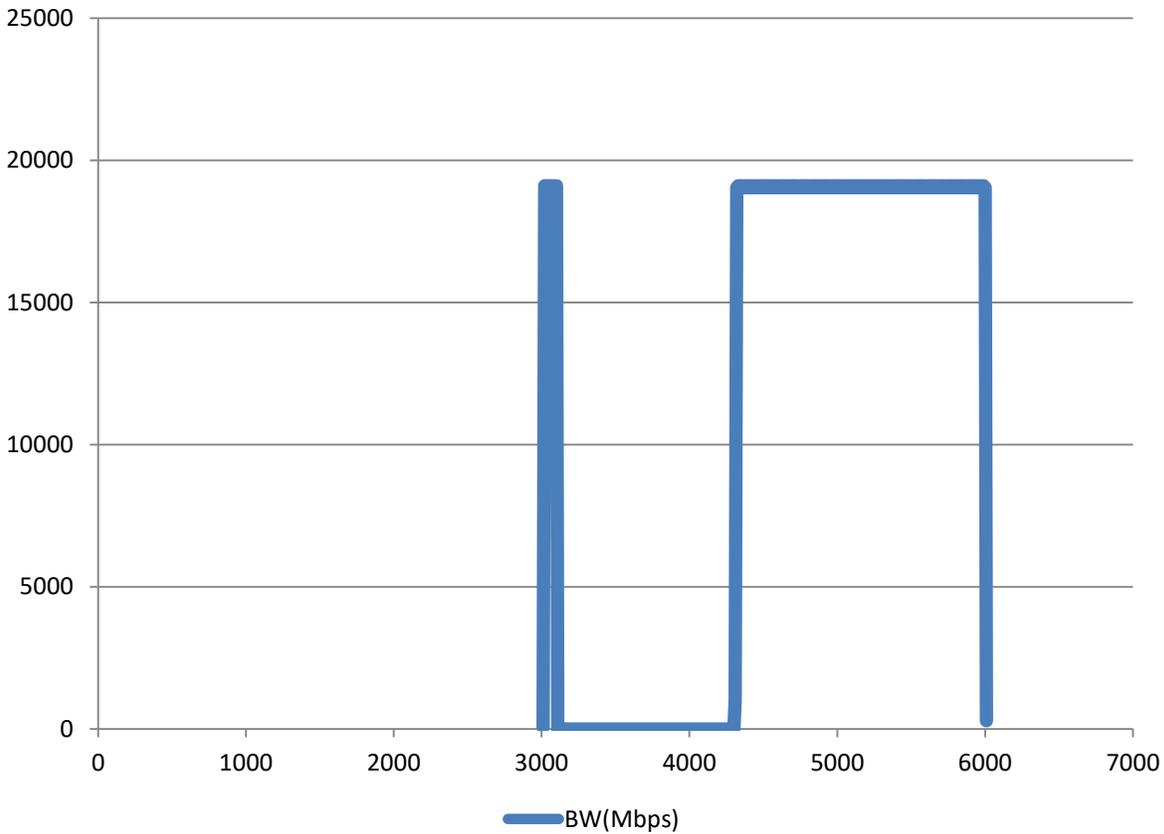
- **New protocol designed for AWS's unique datacenter network**
- **Implemented as part of our 3rd generation Nitro chip**
- **EFA exposes SRD as a reliable datagram interface**
- **Inspired by Infiniband Reliable Datagram, without the drawbacks**
 - No limit on the number of outstanding messages per context
- **Out-of-order delivery – no head-of-line blocking**
 - Messages are independent in many cases, application/middleware can restore ordering only if/when needed
 - Same motivation as weak/relaxed memory ordering
- **Packet spraying over multiple ECMP paths**
 - No hot-spots
 - Fast and transparent recovery from network failures
- **Congestion control designed for large-scale cloud**
 - Prevent packet drops
 - Minimize latency jitter

TCP VS INFINIBAND VS SRD

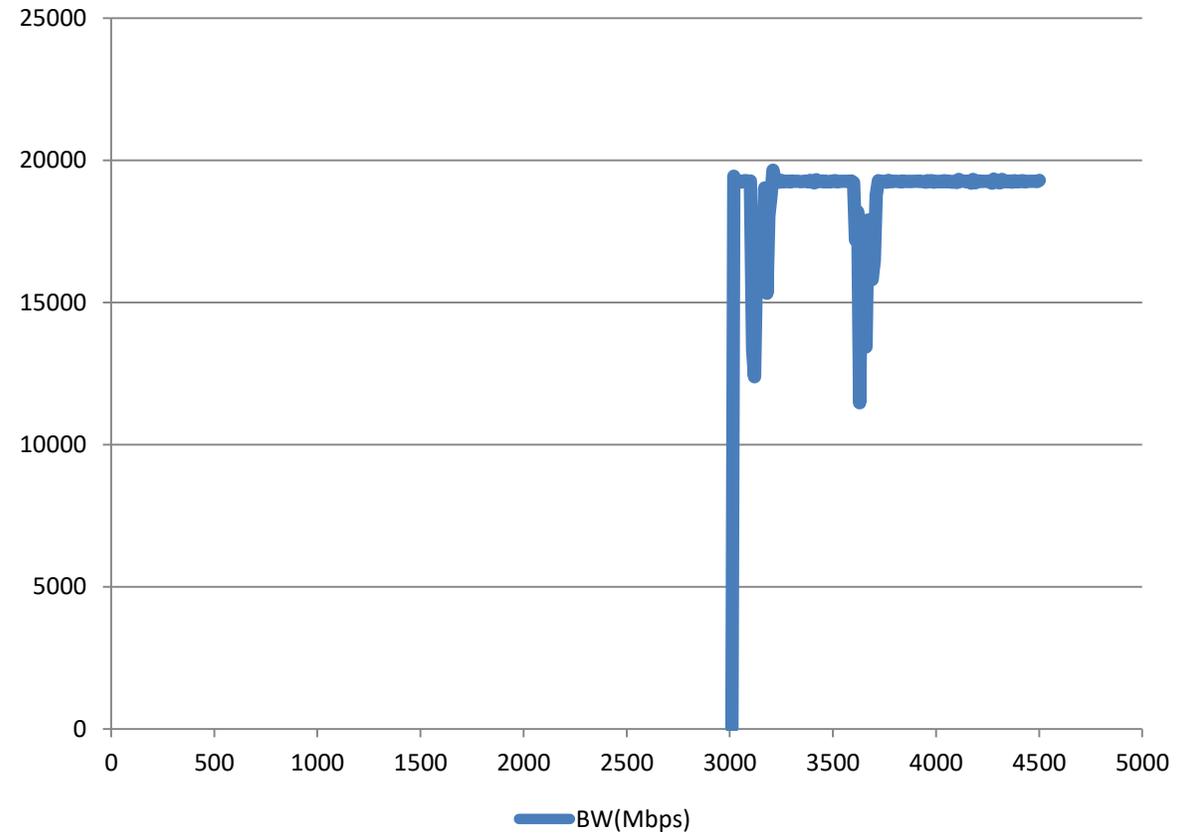
TCP	Infiniband	SRD
Stream	Messages	Messages
In-order	In-order	Out-of-order
Single path	Single (ish) path	ECMP spraying with load balancing
High limit on retransmit timeout (>50ms)	Static user-configured timeout (log scale)	Dynamically estimated timeout (usec resolution)
Loss-based congestion control	Semi-static rate limiting (limited set of supported rates)	Dynamic rate limiting
Inefficient software stack	Transport offload with scaling limitations	Scalable transport offload (same number of QPs regardless cluster size)

SRD LINK FAILURE HANDLING

TCP



SRD



EFA KERNEL MODULE AND RDMA-CORE

- RDMA subsystem in the Linux kernel
- Unreliable datagrams (UD)
- Scalable Reliable Datagram - driver QP type
- RC (and kernel ULPs) not currently supported
- Libibverbs provider for rdma-core
- Driver submitted to linux-rdma@ for upstreaming
 - <https://patchwork.kernel.org/cover/10852679/>

EFA LIBFABRIC PROVIDER

```
provider: efa
fabric: EFA-fe80::82d:33ff:feb5:d1ac
domain: efa_0-rdm
version: 3.0
type: FI_EP_RDM
protocol: FI_PROTO_EFA
provider: efa
fabric: EFA-fe80::82d:33ff:feb5:d1ac
domain: efa_0-dgrm
version: 3.0
type: FI_EP_DGRAM
protocol: FI_PROTO_EFA
```

RDM

- Reliable, unordered datagrams
- ~8 KiB max message size
- Send/receive interface, with no tag matching
- Native multi-pathing; no “flow limit”

DGRAM

- Unreliable, unordered datagrams
- ~8 KiB max message size
- Send/receive interface
- Subject to same “flow limit” as TCP/IP and UDP/IP over ENA

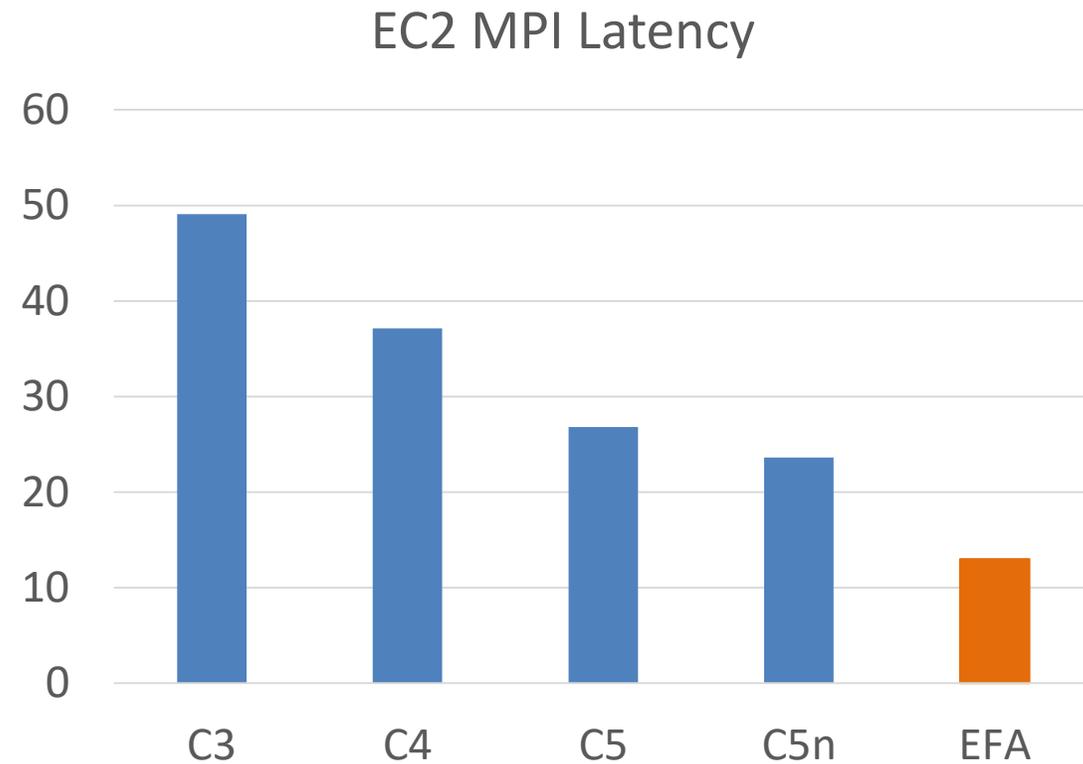
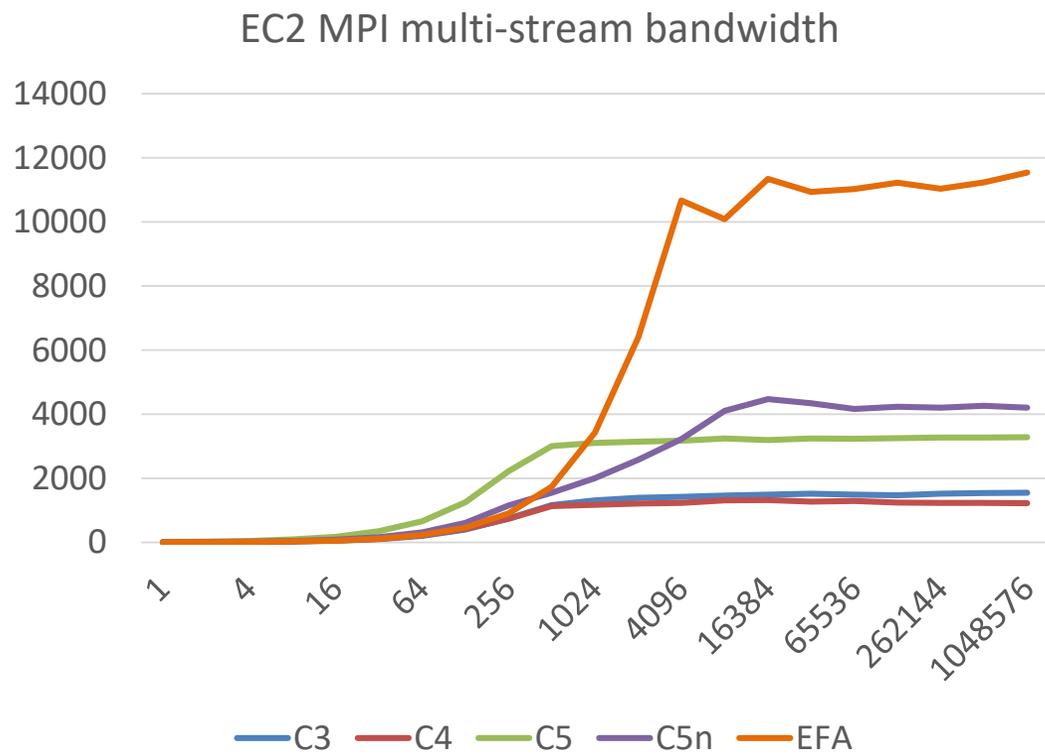
Memory registration cache and userfaultfd monitors

RXR UTILITY PROVIDER

```
provider: efa;ofi_rxr  
fabric: EFA-fe80::82d:33ff:feb5:d1ac  
domain: efa_0-rdm  
version: 1.0  
type: FI_EP_RDM  
protocol: FI_PROTO_RXR
```

- Generic RDM-over-RDM (RxR) utility provider
- Layers over EFA's FI_EP_RDM (SRD protocol)
- Rx packet reordering
- Tag-matching support (fi_t*)
- Segmentation and reassembly for >MTU messages
- Also supports FI_MULTI_RECV, FI_SOURCE, FI_SOURCE_ERR, FI_DIRECTED_RECV

HPC NETWORK PERFORMANCE WITH EFA



USING EFA

Supported platforms

- C5n.18xlarge, P3dn.24xlarge

EFA Kernel module

- <https://github.com/amzn/amzn-drivers>

Libfabric and rdma-core network stack

- AWS-custom version for first half 2019 until we upstream

MPI Implementation or NCCL

- Open MPI 3.1.3 or later or NCCL 2.3.8 or later
- Intel MPI and MPICH in development

1 EFA ENI per instance

See <https://aws.amazon.com/hpc/> for more details

THE ROAD AHEAD

- EFA currently in customer preview, will be Generally Available shortly
- Continue working with the linux-rdma and libfabric communities to upstream
 - Kernel module review: <https://patchwork.kernel.org/cover/10852679/>
 - rdma-core userspace provider review: <https://github.com/linux-rdma/rdma-core/pull/475>
 - Libfabric providers targeting v1.8 release in Summer
- Kernel ULP: We believe we can emulate RC. Looking for feedback to prioritize against other future enhancements.

THE ROAD AHEAD

- Intel MPI will work with EFA in Q2'19
- Constantly iterate on improving performance. Current expectations:
 - Less than 15 μ s $\frac{1}{2}$ RTT in placement group (osu_latency benchmark)
 - 70 Gbps single endpoint MPI bandwidth
 - 100 Gbps system bandwidth
- Extend providers' capabilities support any libfabric-enabled middleware



15th ANNUAL WORKSHOP 2019

THANK YOU

Raghu Raja

Amazon Web Services

(craghun@amazon.com)