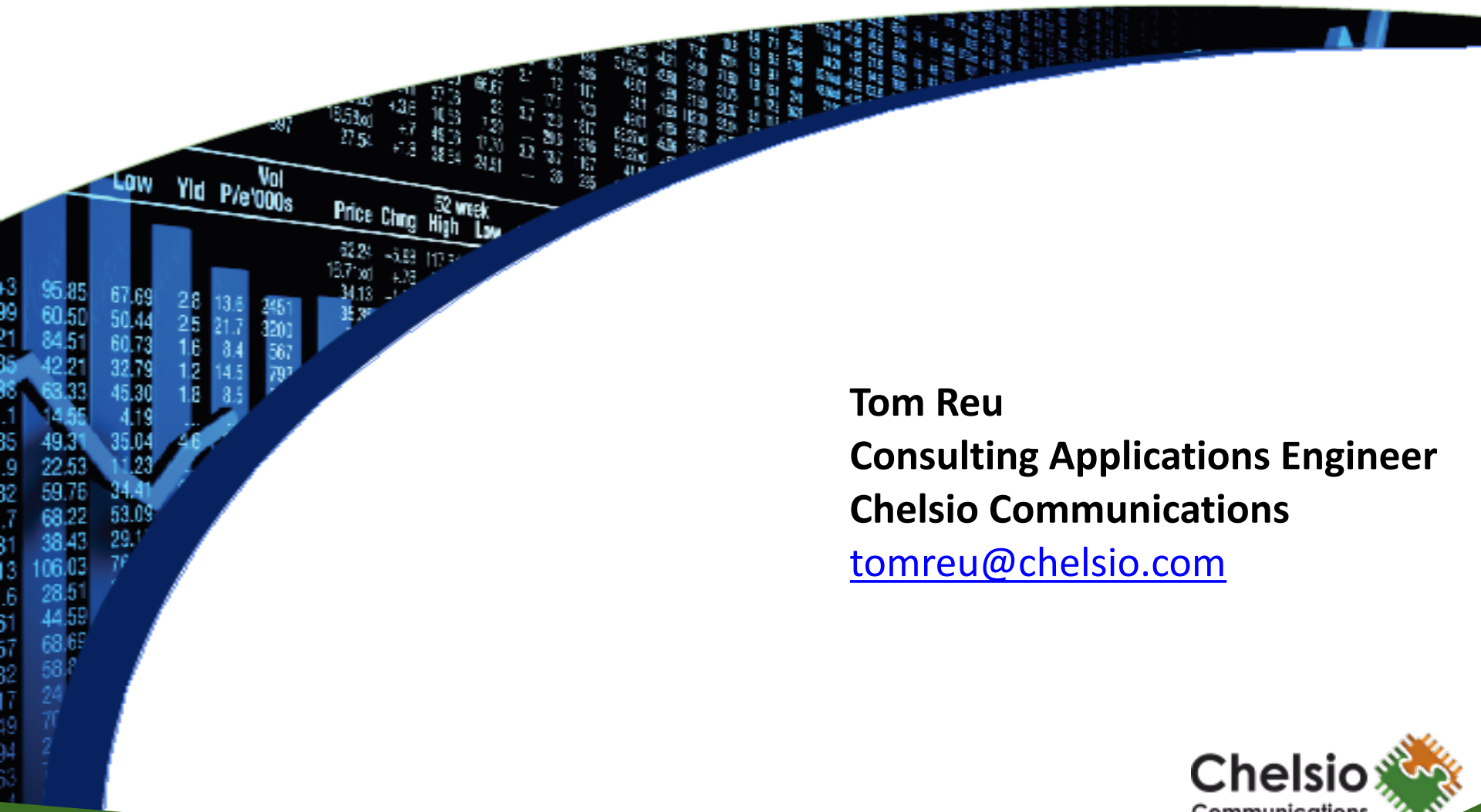


# High-Performance GPU Clustering: GPUDirect RDMA over 40GbE iWARP



**Tom Reu**  
**Consulting Applications Engineer**  
**Chelsio Communications**  
[tomreu@chelsio.com](mailto:tomreu@chelsio.com)

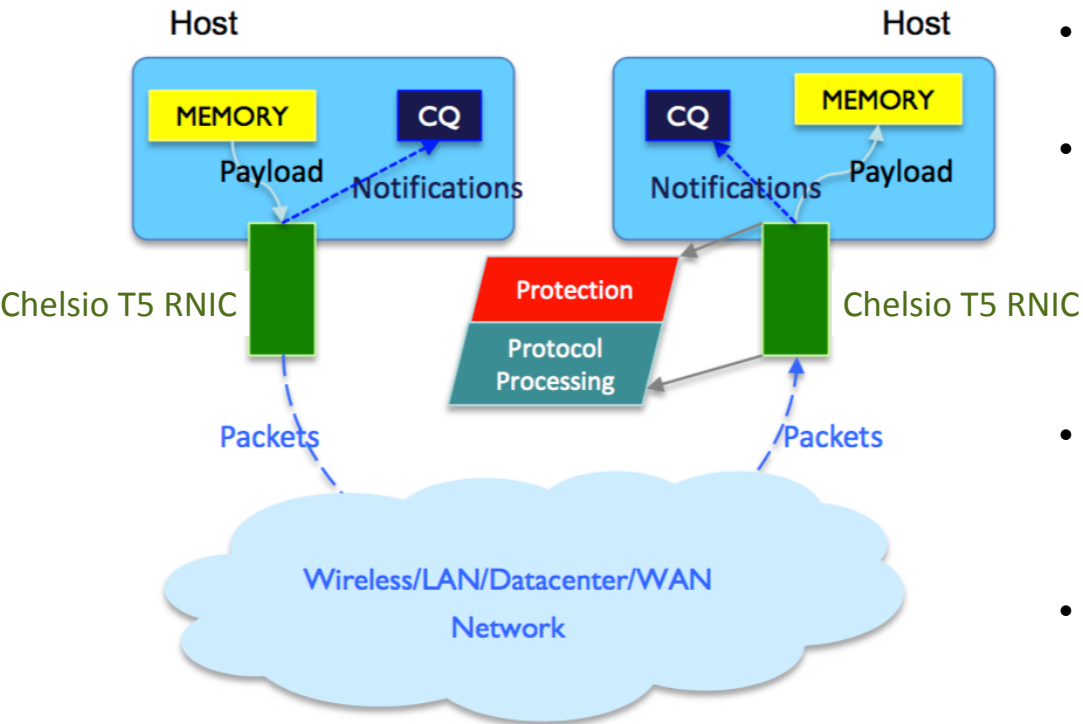
# Chelsio Corporate Snapshot

## Leader in High Speed Converged Ethernet Adapters

- Leading 10/40GbE adapter solution provider for servers and storage systems
  - ~800K ports shipped
- High performance protocol engine
  - 80MPPS
  - 1.5μsec
  - ~5M+ IOPs
- Feature rich solution
  - Media streaming hardware/software
  - WAN Optimization, Security, etc.
- Company Facts
  - Founded in 2000
  - 150 strong staff
- R&D Offices
  - USA – Sunnyvale
  - India – Bangalore
  - China - Shanghai



# RDMA Overview



- Direct memory-to-memory transfer
- All protocol processing handling by the NIC
  - Must be in hardware
- Protection handled by the NIC
  - User space access requires both local and remote enforcement
- Asynchronous communication model
  - Reduced host involvement
- Performance
  - Latency - polling
  - Throughput
- Efficiency
  - Zero copy
  - Kernel bypass (user space I/O)
  - CPU bypass

*Performance and efficiency in return  
for new communication paradigm*

# iWARP

## What is it?

---

- Provides the ability to do Remote Direct Memory Access over Ethernet using TCP/IP
- Uses Well-Known IB Verbs
- Inboxed in OFED since 2008
- Runs on top of TCP/IP
  - Chelsio implements iWARP/TCP/IP stack in silicon
  - Cut-through send
  - Cut-through receive
- Benefits
  - Engineered to use “typical” Ethernet
    - No need for technologies like DCB or QCN
  - Natively Routable
  - Multi-path support at Layer 3 (and Layer 2)
  - It runs on TCP/IP
    - Mature and Proven
    - Goes where TCP/IP goes (everywhere)

# iWARP

---

- iWARP updates and enhancements are done by the IETF STORM (Storage Maintenance) working group
- RFCs
  - RFC 5040 A Remote Direct Memory Access Protocol Specification
  - RFC 5041 Direct Data Placement over Reliable Transports
  - RFC 5044 Marker PDU Aligned Framing for TCP Specification
  - RFC 6580 IANA Registries for the RDDP Protocols
  - RFC 6581 Enhanced RDMA Connection Establishment
  - RFC 7306 Remote Direct Memory Access (RDMA) Protocol Extensions
- Support from several vendors, Chelsio, Intel, QLogic

# iWARP

## Increasing Interest in iWARP as of late

- Some Use Cases
  - High Performance Computing
  - SMB Direct
  - GPUDirect RDMA
  - NFS over RDMA
  - FreeBSD iWARP
  - Hadoop RDMA
  - Lustre RDMA
  - NVMe over RDMA fabrics

# iWARP

## Advantages over Other RDMA Transports

- It's Ethernet
  - Well Understood and Administered
  - Uses TCP/IP
    - Mature and Proven
    - Supports rack, cluster, datacenter, LAN/MAN/WAN and wireless
    - Compatible with SSL/TLS
  - Do not need to use any bolt-on technologies like
    - DCB
    - QCN
- Does not require a totally new network infrastructure
  - Reduces TCO and OpEx

# iWARP vs RoCE

iWARP	RoCE
Native TCP/IP over Ethernet, no different from NFS or HTTP	Difficult to install and configure - “needs a team of experts” - <b>Plug-and-Debug</b>
Works with ANY Ethernet switches	Requires DCB - <b>expensive</b> equipment upgrade
Works with ALL Ethernet equipment	Poor <b>interoperability</b> - may not work with switches from different vendors
No need for special QoS or configuration - <b>TRUE Plug-and-Play</b>	<b>Fixed</b> QoS configuration - DCB must be setup identically across all switches
No need for special configuration, preserves network <b>robustness</b>	Easy to break - switch configuration can cause performance <b>collapse</b>
TCP/IP allows reach to <b>Cloud</b> scale	Does not <b>scale</b> - requires PFC, limited to single subnet
No distance limitations. Ideal for <b>remote</b> communication and HA	Short <b>distance</b> - PFC range is limited to few hundred meters maximum
WAN <b>routable</b> , uses any IP infrastructure	RoCEv1 not <b>routable</b> . RoCE v2 requires lossless IP infrastructure and restricts router configuration
Standard for whole stack has been <b>stable</b> for a decade	ROCEv2 <b>incompatible</b> with v1. More fixes to missing reliability and scalability layers required and expected
<b>Transparent and open</b> IETF standards process	Incomplete specification and <b>opaque</b> process



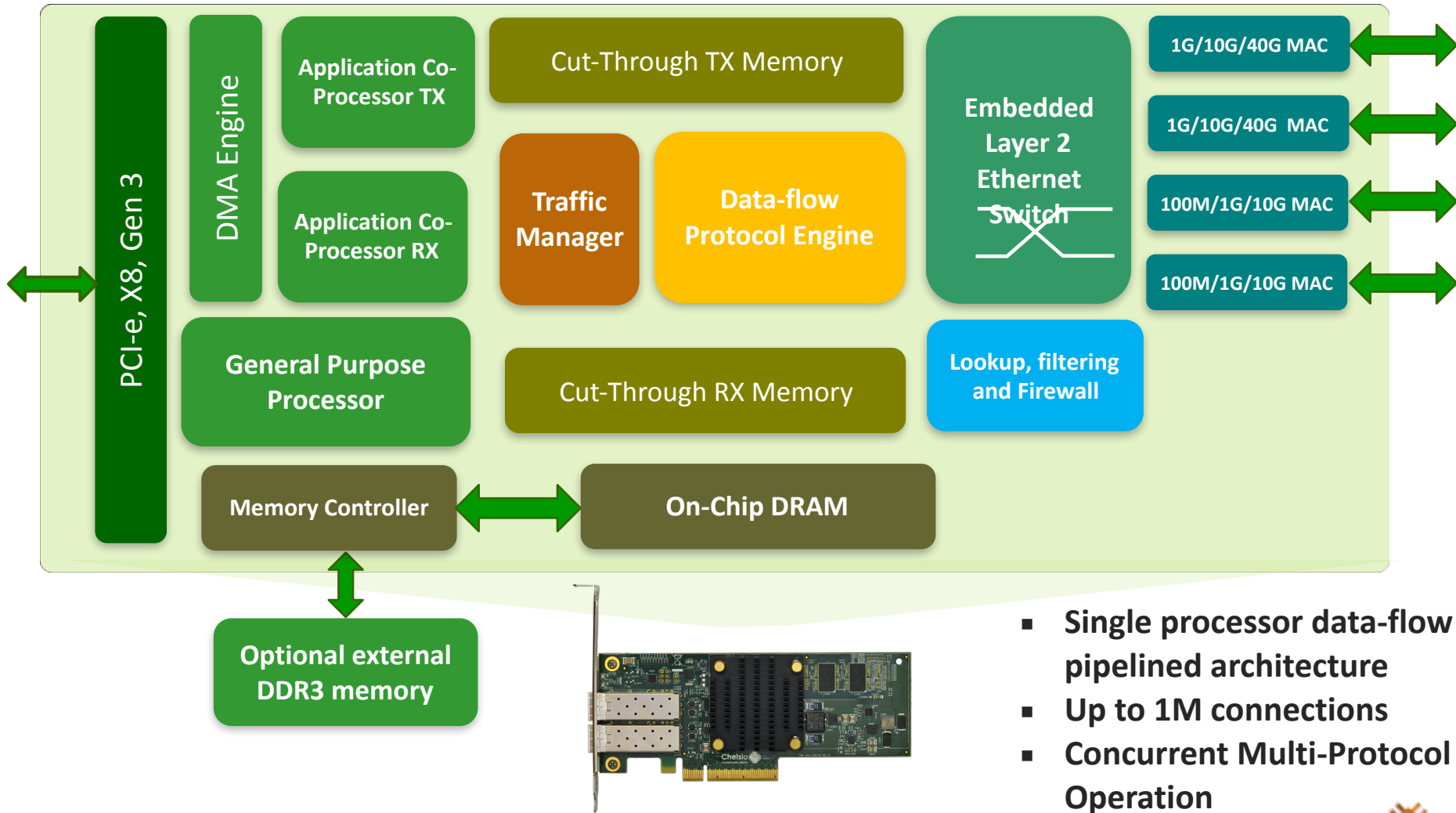
# Chelsio's T5

## Single ASIC does it all

- High Performance Purpose Built Protocol Processor
- Runs multiple protocols
  - TCP with Stateless Offload and Full Offload
  - UDP with Stateless Offload
  - iWARP
  - FCoE with Offload
  - iSCSI with Offload
- All of these protocols run on T5 with a SINGLE FIRMWARE IMAGE
  - No need to reinitialize the card for different uses
  - Future proof e.g. support for NVMf yet preserves today's investment in iSCSI

# T5 ASIC Architecture

## High Performance Purpose Built Protocol Processor

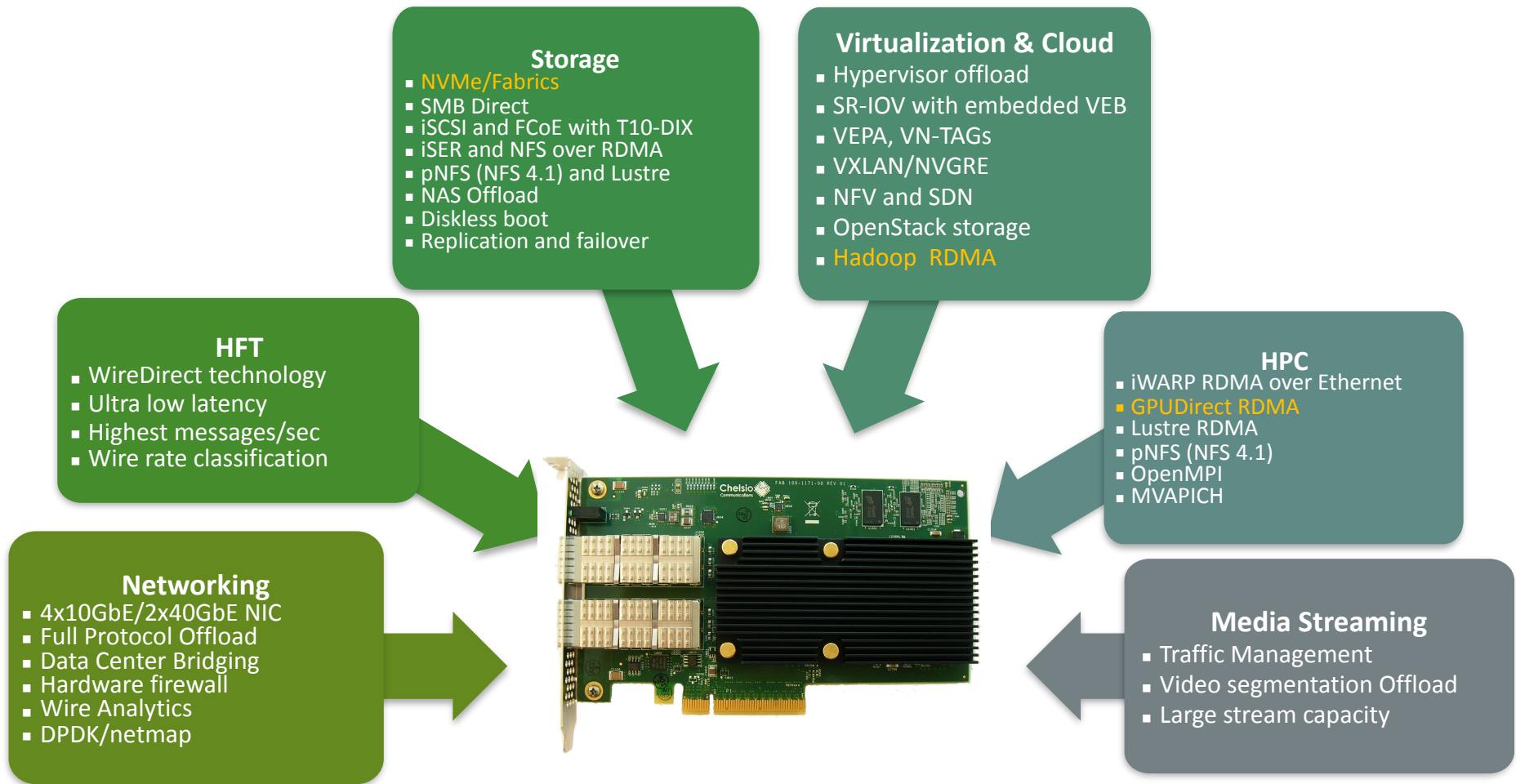


Single connection at 40Gb. Low Latency.

- Single processor data-flow pipelined architecture
- Up to 1M connections
- Concurrent Multi-Protocol Operation

# Leading Unified Wire™ Architecture

## Converged Network Architecture with all-in-one Adapter and Software



**Single Qualification – Single SKU**  
**Concurrent Multi-Protocol Operation**

# GPUDirect RDMA

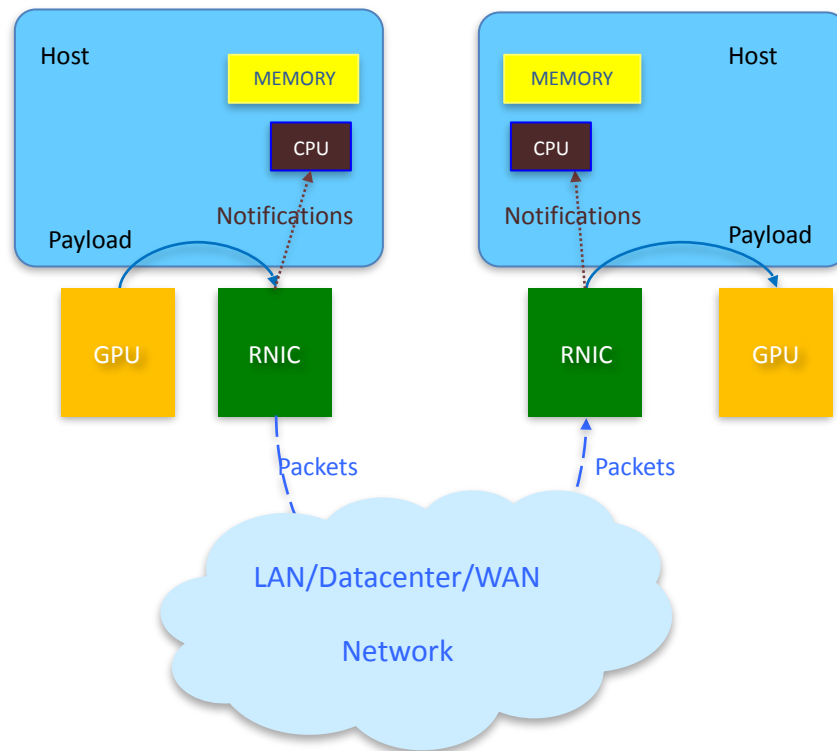
---

- Introduced by NVIDIA with the Kepler Class GPUs. Available today on Tesla and Quadro GPUs as well.
- Enables Multiple GPUs, 3rd party network adapters, SSDs and other devices to read and write CUDA host and device memory
- Avoids unnecessary system memory copies and associated CPU overhead by copying data directly to and from pinned GPU memory
- One hardware limitation
  - The GPU and the Network device **MUST** share the same upstream PCIe root complex
- Available with Infiniband, RoCE, and now iWARP

# GPUDirect RDMA

## T5 iWARP RDMA over Ethernet certified with NVIDIA GPUDirect

- Read/write GPU memory directly from network adapter
  - Peer-to-peer PCIe communication
  - Bypass host CPU
  - Bypass host memory
- Zero copy
- Ultra low latency
- Very high performance
- Scalable GPU pooling
  - Any Ethernet networks



# Modules required for GPUDirect RMDA with iWARP

---

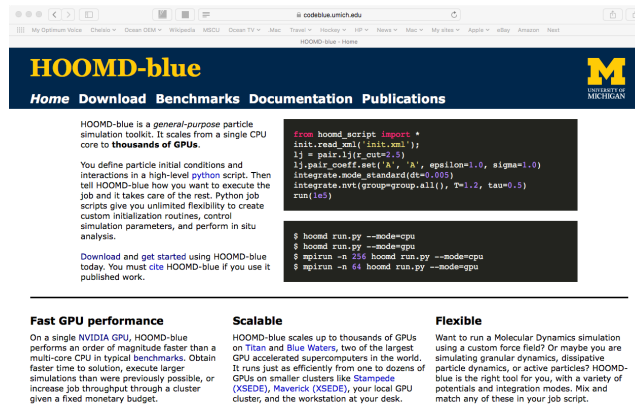
- Chelsio Modules
  - cxgb4 - Chelsio adapter driver
  - iw\_cxgb4 - Chelsio iWARP driver
  - rdma\_ucm - RDMA User Space Connection Manager
- NVIDIA Modules
  - nvidia - NVIDIA driver
  - nvidia\_uvm - NVIDIA Unified Memory
  - nv\_peer\_mem - NVIDIA Peer Memory



# Case Studies

# HOOMD-blue

- General Purpose Particle simulation toolkit
- Stands for: Highly Optimized Object-oriented Many-particle Dynamics - Blue Edition
- Running on GPUDirect RDMA - WITH NO CHANGES TO THE CODE - AT ALL!
- More Info: [www.codeblue.umich.edu/hoomd-blue](http://www.codeblue.umich.edu/hoomd-blue)



**HOOMD-blue**  
Home Download Benchmarks Documentation Publications

HOOMD-blue is a general-purpose particle simulation toolkit. It scales from a single CPU core to thousands of GPUs.

You define particle initial conditions and interactions in a high-level python script. Then tell HOOMD-blue how you want to execute the job and it takes care of the rest. Python job scripts give you unlimited flexibility to create custom initialization routines, control simulation parameters, and perform in situ analysis.

Download and get started using HOOMD-blue today. You must cite HOOMD-blue if you use it published work.

```
from hoomd_script import *
init_read_xml('main.xml')
lj = pair_lj(r_out=2.5)
lj.pair_coeff.set('A', 'S', epsilon=1.0, sigma=1.0)
integrate.mode_standard(d=0.05)
integrate.nvt(group=group.all(), T=1.2, tau=0.5)
run(1e5)
```

```
$ hoomd run.py --mode=cpu
$ hoomd run.py --mode=gpu
$ mpirun -n 256 hoomd run.py --mode=cpu
$ mpirun -n 64 hoomd run.py --mode=gpu
```

**Fast GPU performance**  
On a single NVIDIA GPU, HOOMD-blue performs an order of magnitude faster than a multi-core CPU in typical benchmarks. Obtain faster time to solution, execute larger simulations than were previously possible, or increase job throughput through a cluster given a fixed monetary budget.

**Scalable**  
HOOMD-blue scales up to thousands of GPUs on Titan and Blue Waters, two of the largest GPU accelerated supercomputers in the world. It runs just as efficiently from one to dozens of GPUs on smaller clusters like Stampede (KSEDE), Maverick (KSEDE), your local GPU cluster, and the workstation at your desk.

**Flexible**  
Want to run a Molecular Dynamics simulation using a custom force field? Or maybe you are simulating granular dynamics, dissipative particle dynamics, or active particles? HOOMD-blue is the right tool for you, with a variety of potentials and integration modes. Mix and match any of these in your job script.

## Molecular Dynamics capabilities

### Integrators

Apply any number of integrators to separate particle groups. HOOMD-blue has integrators built in for many different thermodynamic ensembles and energy minimization. Many of them support integration of orientational degrees of freedom.

- NVE, NVT, NPH, NPT
- Langevin dynamics
- Brownian dynamics

### Pair potentials

- CGCM
- DPD
- Lennard-Jones
- Gaussian
- Mie
- Molierie
- Morse
- Yukawa
- ZBL
- User-defined (table)

### Bond potentials

- Harmonic
- FENE
- User-defined (table)

### Angle potentials

- Harmonic
- CGCM
- User-defined (table)



## Test Configuration

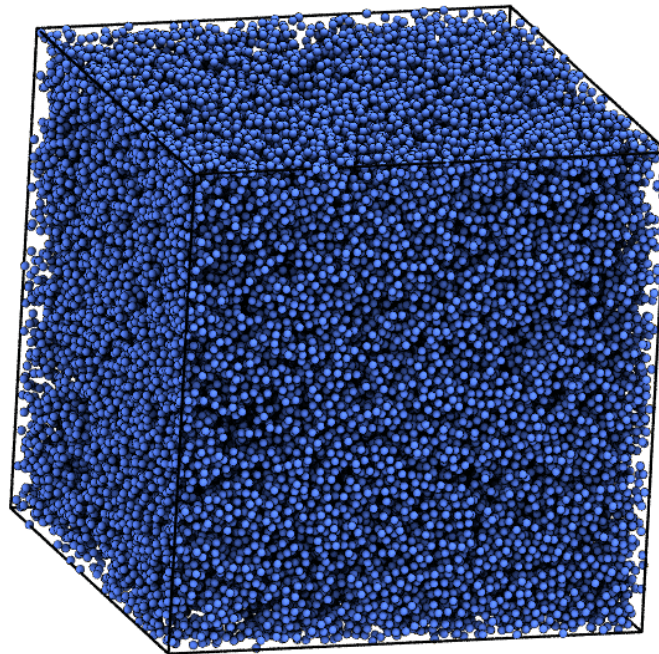
- 4 Nodes
- Intel E5-1660 v2 @ 3.7 Ghz
- 64 GB RAM
- Chelsio T580-CR 40Gb Adapter
- NVIDIA Tesla K80 (2 GPUs per card)
- RHEL 6.5
- OpenMPI 1.10.0
- OFED 3.18
- CUDA Toolkit 6.5
- HOOMD-blue v1.3.1-9
- Chelsio-GDR-1.0.0.0
- Command Line:

```
$MPI_HOME/bin/mpirun --allow-run-as-root -mca btl_openib_want_cuda_gdr 1  
-np X -hostfile /root/hosts -mca btl openib,sm,self -mca  
btl_openib_if_include cxgb4_0:1 --mca btl_openib_cuda_rdma_limit 65538 -  
mca btl_openib_receive_queues P,131072,64 -x CUDA_VISIBLE_DEVICES=0,1 /  
root/hoomd-install/bin/hoomd ./bmark.py --mode=gpu|cpu
```

# HOOMD-blue

## Lennard-Jones Liquid 64K Particles Benchmark

- Classic benchmark for general purpose MD simulations.
- Representative of the performance HOOMD-blue achieves for straight pair potential simulations

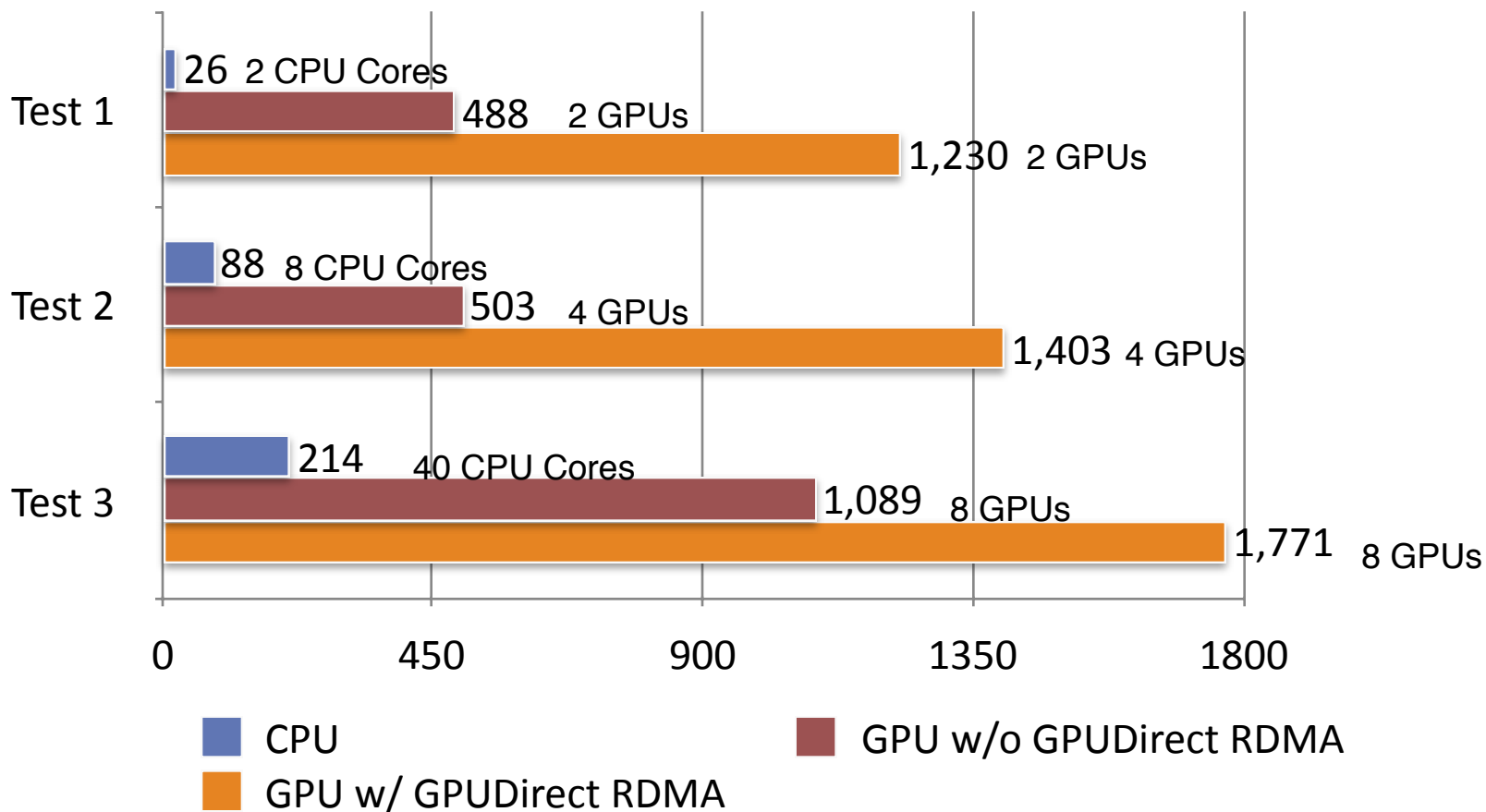


# HOOMD-blue



## Lennard-Jones Liquid 64K Particles Benchmark Results

Average Timesteps per Second



Longer is Better

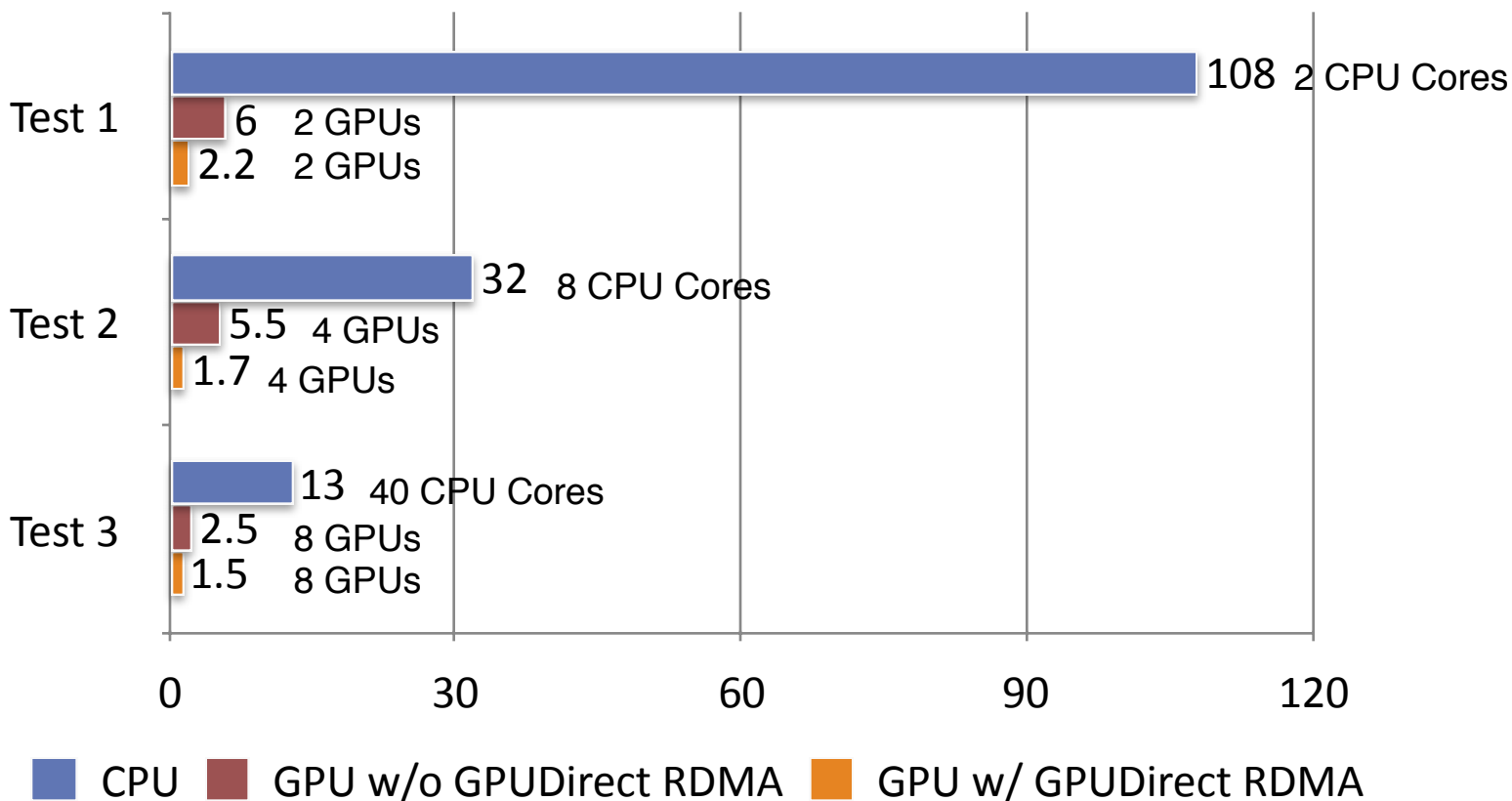


# HOOMD-blue



## Lennard-Jones Liquid 64K Particles Benchmark Results

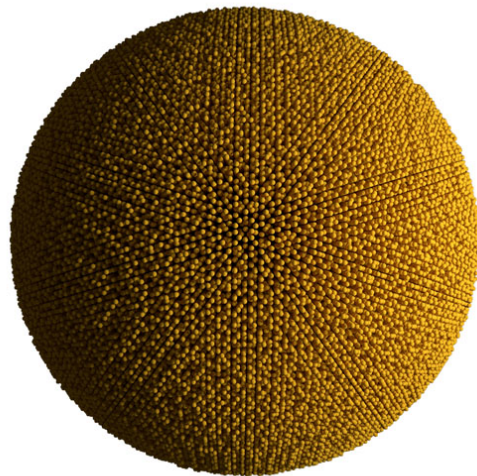
Hours to complete 10e6 steps



Shorter is Better

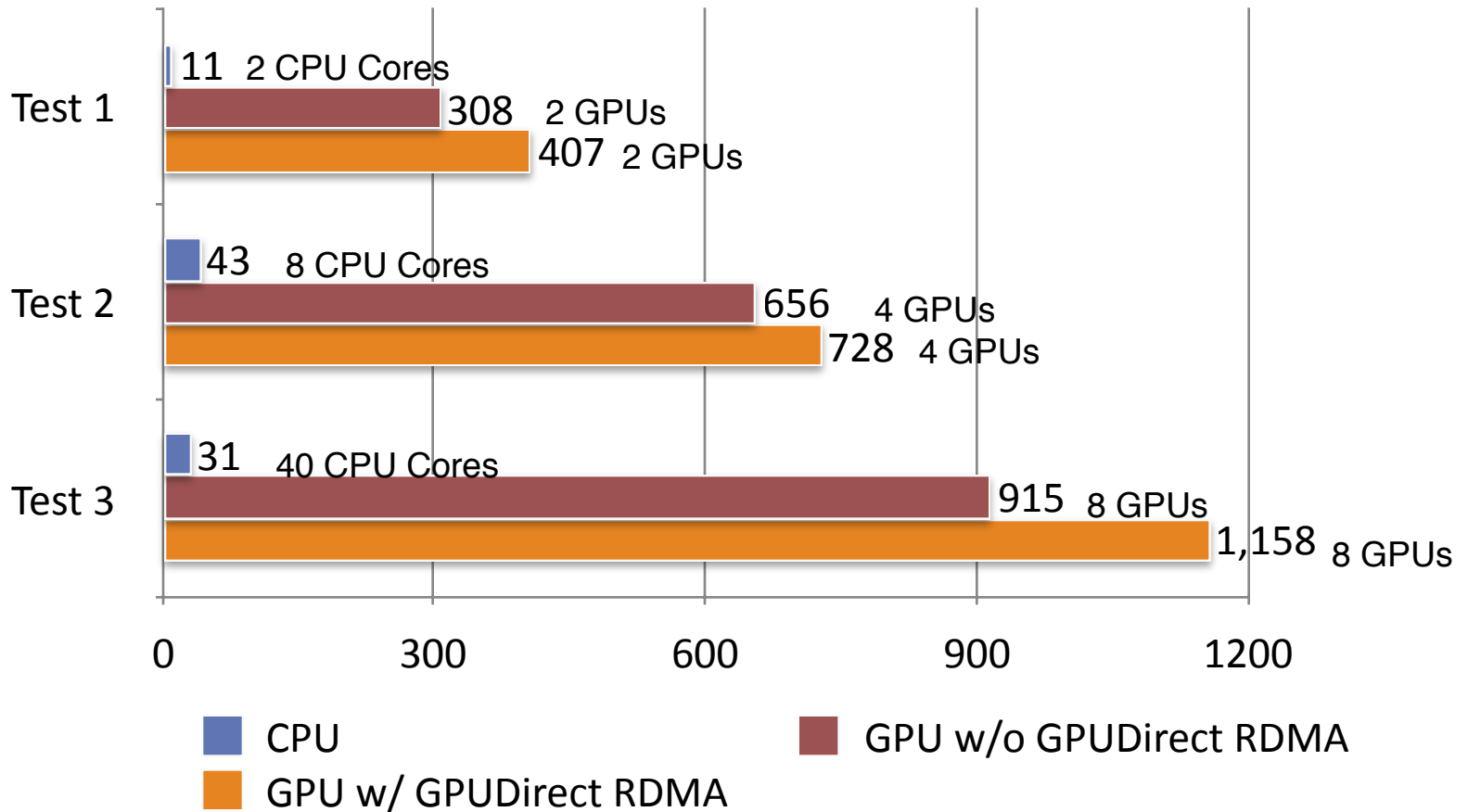


- runs a system of particles with an oscillatory pair potential that forms a icosahedral quasicrystal
- This model is used in the research article: Engel M, et. al. (2015) Computational self-assembly of a one-component icosahedral quasicrystal, Nature materials 14(January), p. 109-116.



## Quasicrystal results

Average Timesteps per Second



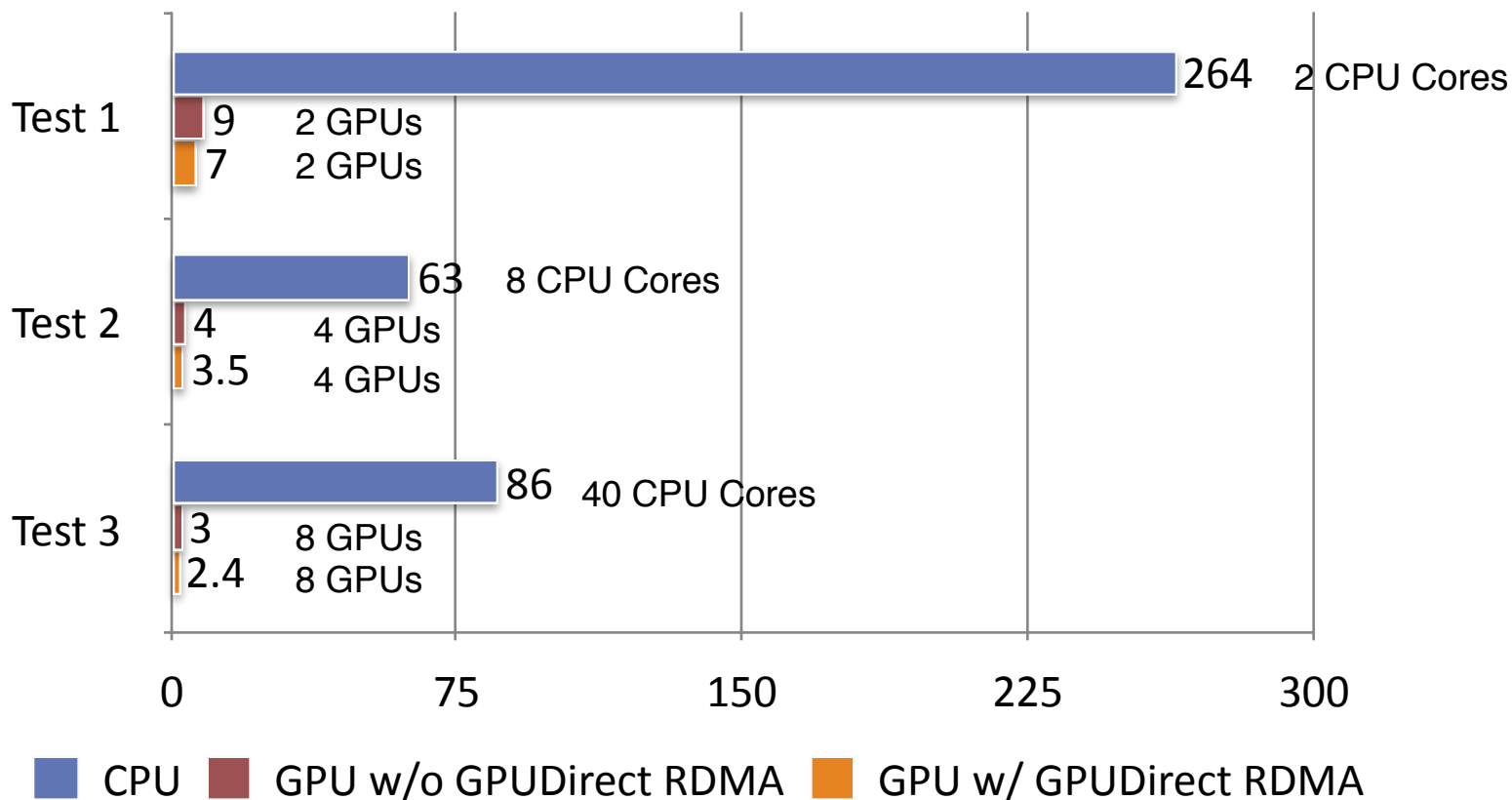
Longer is Better

# HOOMD-blue

## Quasicrystal results



Hours to complete 10e6 steps



Shorter is Better



# Caffe

## Deep Learning Framework

---

- Open source Deep Learning software from Berkeley Vision and Learning Center
- Updated to include CUDA support to utilize GPUs
- Standard version does NOT include MPI support
- MPI implementations
  - mpi-caffe
    - Used to train a large network across a cluster of machines
    - model-parallel distributed approach.
  - caffe-parallel
    - Faster framework for deep learning.
    - data-parallel via MPI, splits the training data across nodes



# Summary

## GPUDirect RDMA over 40GbE iWARP

- iWARP provides RDMA Capabilities to a Ethernet network
- iWARP uses tried and true TCP/IP as its underlying transport mechanism
- Using iWARP does not require a whole new network infrastructure and the management requirements that come along with it
- iWARP can be used with existing software running on GPUDirect RDMA which NO CHANGES required to the code
- Applications that use GPUDirect RDMA will see huge performance improvements
- Chelsio provides 10/40Gb iWARP TODAY with 25/50/100 Gb on the horizon

# More information

## GPUDirect RDMA over 40GbE iWARP

- Visit our website, [www.chelsio.com](http://www.chelsio.com), for more White Papers, Benchmarks, etc.
- GPUDirect RDMA White Paper: <http://www.chelsio.com/wp-content/uploads/resources/T5-40Gb-Linux-GPUDirect.pdf>
- Webinar : <https://www.brighttalk.com/webcast/13671/189427>
- Beta code for GPUDirect RDMA is available TODAY from our download site at [service.chelsio.com](http://service.chelsio.com)
- Sales questions - [sales@chelsio.com](mailto:sales@chelsio.com)
- Support questions - [support@chelsio.com](mailto:support@chelsio.com)

# Questions?



# Thank You

