# Computing System Trends and Performance Portability

Bálint Joó – OLCF, Oak Ridge National Laboratory

Lattice 2021 Virtual Meeting (MIT)

July 30, 2021

U.S. DEPARTMENT OF **ENERGY**

# Contents

- Machine Trends

- Node Programming Models

- Porting Stories

OAK RIDGE | LEADERSHIP
National Laboratory | COMPUTING FACILITY

Open slide master to edit

# Introduction – a brief history

- I recall at Lattice 2006 in Tuscon, first hearing about the idea of LQCD calculations on GPUs in a talk by Daniel Nogradi.
  - Programming was done using OpenGL.
  - I thought: "This looks fun, but it will never catch on!"

- I had to eat my words as I gave a talk about GPU accelerated computing in Squaw Valley at Lattice 2011 – 10 years ago.

- 2012-2018 was an era of 'friendly competition' between NVIDIA GPUs and Intel Knights
  - The Knights fought well, but were discontinued
  - Remaining fighting Knights are getting close to retirement

- OLCF Summit exceeded 1.88 ExaOps (in 32 and 16-bit precisions) on a Genomics Machine Learning Application in 2018, kicking off the Exascale era
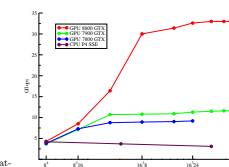


Lattice QCD as a video game

Győző I. Egri[a], Zoltán Fodor[abc], Christian Hoelbling[b],
Sándor D. Katz[ab], Dániel Nógrádi[b] and Kálmán K. Szabó[b]

[a] Institute for Theoretical Physics, Eötvös University, Budapest, Hungary
[b] Department of Physics, University of Wuppertal, Germany
[c] Department of Physics, University of California, San Diego, USA

**Abstract**

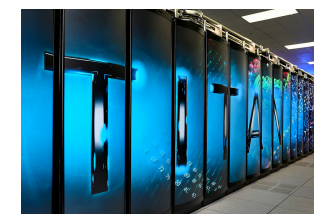The speed, bandwidth and cost characteristics of today's PC graphics cards make them an at-

*Egri et. al. "Lattice QCD as a video game" Comput. Phys. Commun. 177:631-639,2007 arXiv;hep-lat/0611022*



*JLab 9g cluster (2009)*
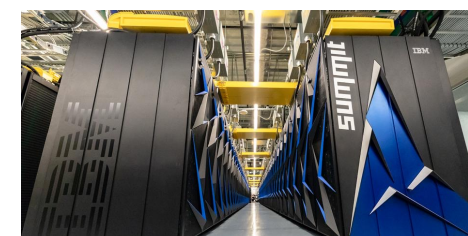


*JLab 10g cluster (2010)*



*OLCF Titan (2012)*



*JLab 12k/m cluster (2012)*



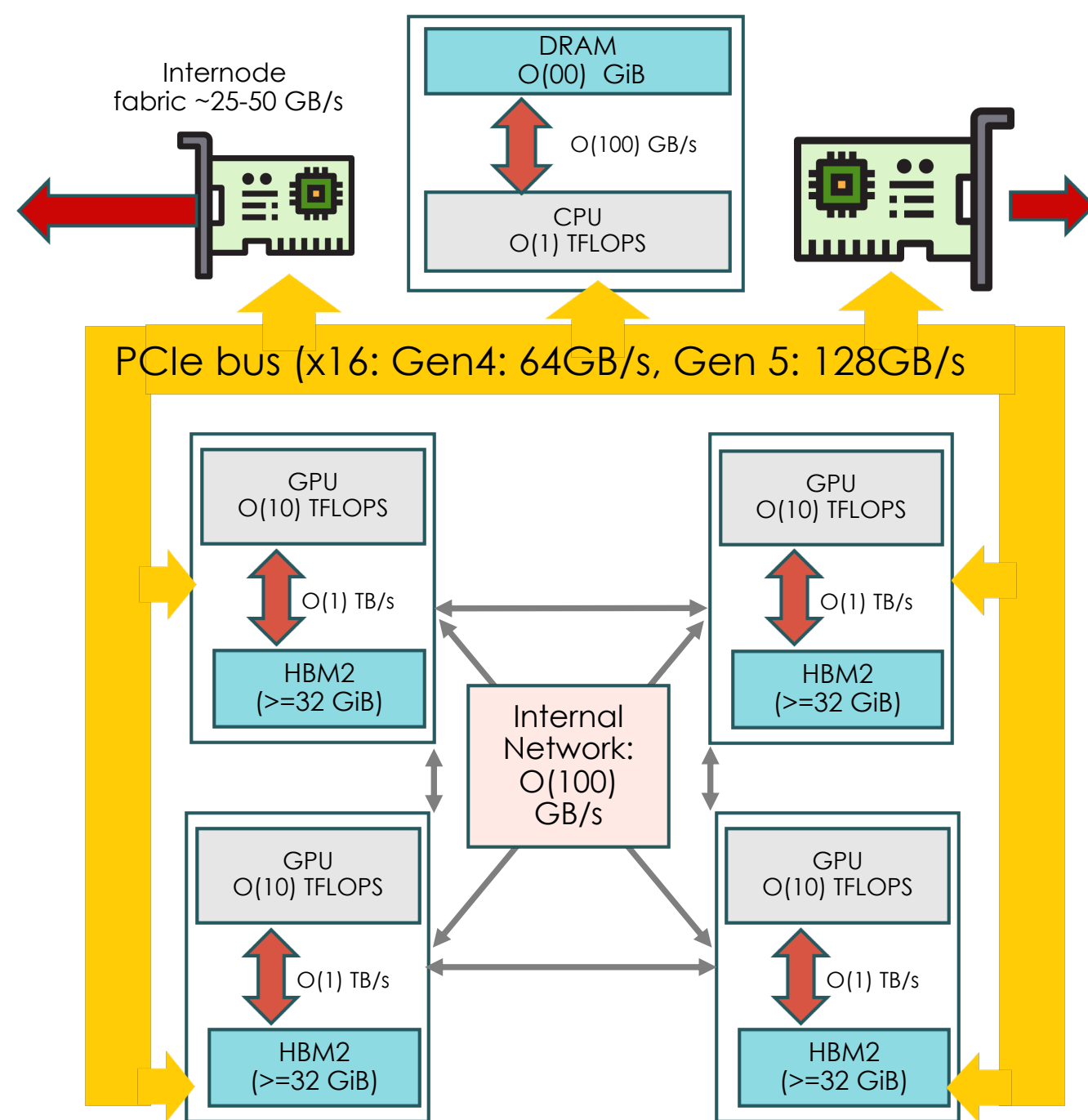*NERSC Cori KNL (2016)*

*JLab 16p and 18p KNL Cluster (2016-...)*



*OLCF Summit 2017*

**OAK RIDGE** National Laboratory | LEADERSHIP COMPUTING FACILITY

# Hardware

OAK RIDGE
National Laboratory | LEADERSHIP
COMPUTING
FACILITY

Open slide master to edit
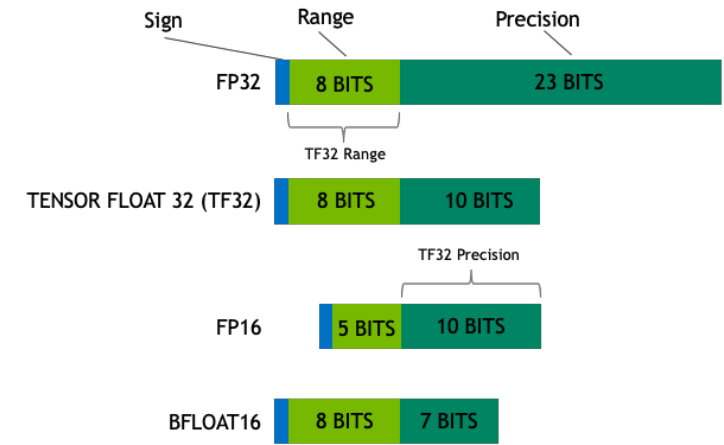
# Typical Pre-Exascale Node

- Schematically not much has changed from Summit
  - Internal network
    - NVLink; 50GB/s per link, A100: up to 12 links=>600 GB/s
    - Infinity Fabric: 92 GB/s link, MI100: up 3 links=> 276 GB/sec
  - External Network
    - Cray/HPE Slingshot (Slingshot **10**: 12.5 + 12.5 = 25 GB/s)
    - Infiniband: (NDR: 50 GB/sec)
  - GPU compute ~ ~10x CPU compute
  - GPU mem-bandwidth ~10x CPU bandwidth

- GPUs have fast low precision matrix-matrix multiply capability
  - NVIDIA Tensor Cores (V100, A100)
  - AMD CDNA (MI100)

- In these accelerated nodes GPUs provide majority of performance
  - I will **blissfully ignore CPU details from here on in**

- Most (not all) LQCD stencils are memory bandwidth bound
  - 1EFLOPS peak => 100 PB/sec peak bandwidth
  - F/B ratio ~0.5-1 for Wilson like fermions in 64 Bit

Internode fabric ~25-50 GB/s

DRAM O(00) GiB

O(100) GB/s

CPU O(1) TFLOPS

PCIe bus (x16: Gen4: 64GB/s, Gen 5: 128GB/s

GPU O(10) TFLOPS

O(1) TB/s

HBM2 (>=32 GiB)

GPU O(10) TFLOPS

O(1) TB/s

HBM2 (>=32 GiB)

Internal Network: O(100) GB/s

GPU O(10) TFLOPS

O(1) TB/s

HBM2 (>=32 GiB)

GPU O(10) TFLOPS

O(1) TB/s

HBM2 (>=32 GiB)

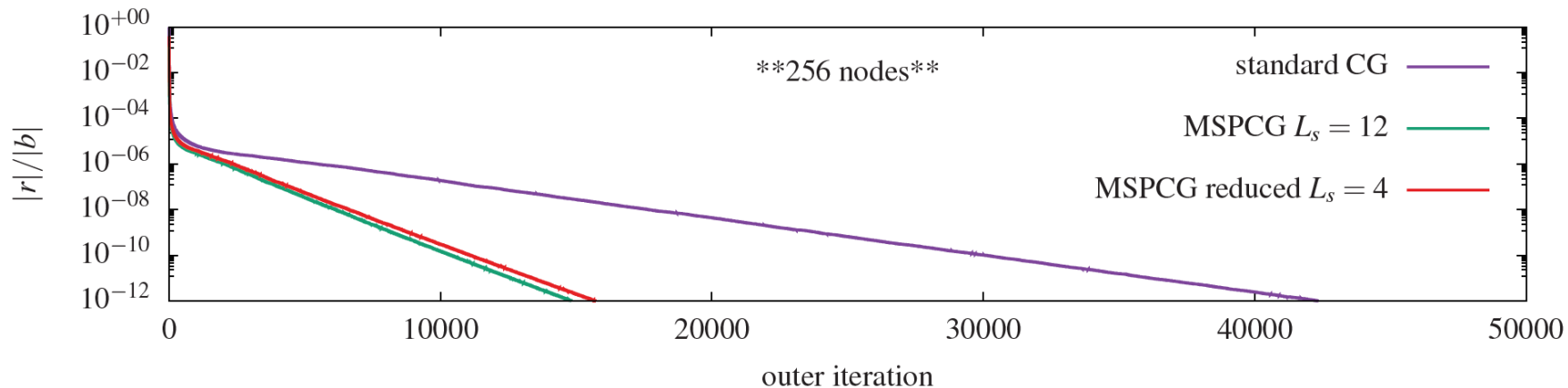OAK RIDGE National Laboratory | LEADERSHIP COMPUTING FACILITY

Open slide master to edit

# QUDA and tensor core

- Speed up compute intensive kernels in QUDA with tensor cores that accelerate MMA operations on NVIDIA GPUs [poster D9]
  - [WIP] Multi-BLAS
  - m5inv in MSPCG/additive Schwarz: **1.5x** time-to-solution speed up [arXiv:2104.05615, J. Tu's parallel 07/28 14:30]
  - Multi-grid setup: as high as **13x** speed up for `computeUV` on Volta

- [WIP] MMA support for all four precisions: double (fp64), single (fp32), half (16-bit fixed-point), quarter (8-bit fixed-point)
  - (hi + lo) * (hi + lo) = hi * hi + hi * lo + lo * hi + ~~lo * lo~~
  - Use 3xTFLOAT32 for single, and 3xBFLOAT16 for half



| Stage | V100 TFLOPS | V100 speed up | A100 TFLOPS | A100 speed up |
|-------|-------------|---------------|-------------|---------------|
| $UV$ | 22.7 | 13.0 | 48.4 | 12.1 |
| $V^\dagger UV$ | 6.5 | 1.5 | 11.4 | 2.0 |
| $X^{-1}Y$ | 42.0 | 9.2 | 70.6 | 11.6 |



577.6 seconds
458.3 seconds
381.5 seconds

OAK RIDGE National Laboratory | LEADERSHIP COMPUTING FACILITY

Slide Courtesy of Jiqun Tu, NVIDIA
Open slide master to edit

# Latency Or Bandwidth?  Both Please...



→ Message size for Face (DP) - Wilson Clover Fermions, $V=L^4$ sites

→ Message size for Face (DP) - DWF Fermions, $V=L^4$ sites, $L_s=16$

L=24 (648K)

L=32 (1.5M)

L=16 (192K)

L=12 (81K)

L=8 (24K)

L=2 (1.5K)

L=2 (3K)

Latency bound

L=16 (3M)

L=12 (1.3M)

L=8 (384K)

Bandwidth bound

½ bandwidth

**Unidirectional Network Bandwidth (MB/sec)**

25000
20000
15000
10000
5000
0

0  1  2  4  8  16  32  64  128  256  512  1024  2048  4096

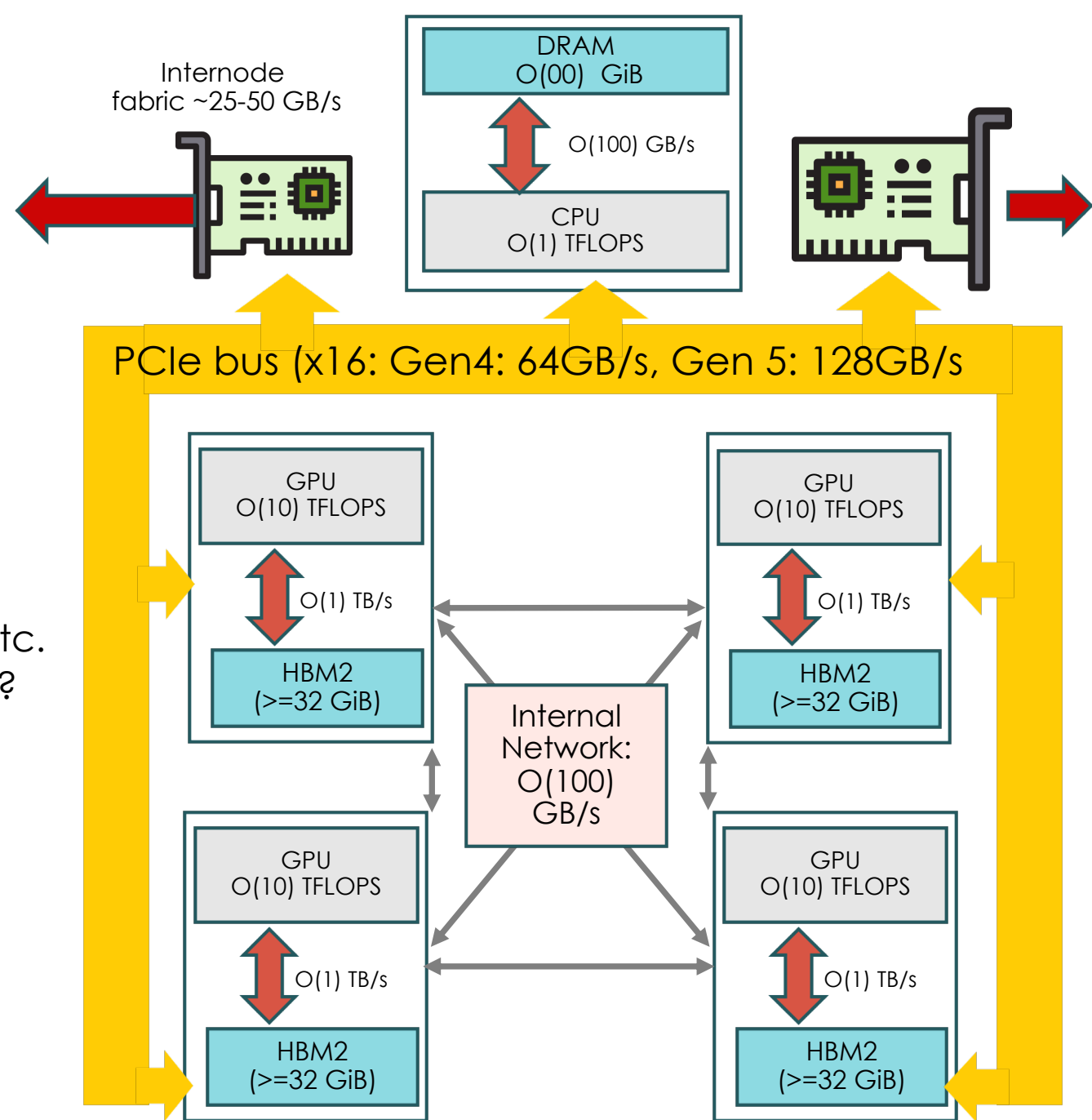→ Message Size (HP) – MG $V=2^4$ $N_\chi=2$ $N_{vec}=24$

→ Message Size (SP) – MG $V=2^4$ $N_\chi=2$ $N_{vec}=24$

**Message Size (KiB)**

Data for Summit from SummitWorkshop-SummitNodePerformane-WJ talk by W. Joubert

OAK RIDGE National Laboratory | LEADERSHIP COMPUTING FACILITY

Open slide master to edit

# And so, nuances matter

- Hardware capability
  - Direct path from GPUs to NIC?
  - Can GPUs initiate communications?
  - Number of hops in internode networks
  - Amount of GPU memory & H2D2H transfers
  - Unified/Managed Memory Capability?

- Programming model & implementation
  - Kernel launch latency
  - GPU aware MPI, Interfaces for P2P, SHMEM etc.
  - Persistent kernels, synchronizing with atomics?
  - Access to Unified/Managed Memories

- Algorithmic Pushes
  - MG & Reduced precision -> latency
  - Multi-RHS/Split Grid algorithms -> B/W
  - DD algorithms -> device locality

Internode fabric ~25-50 GB/s

DRAM
O(00) GiB

O(100) GB/s

CPU
O(1) TFLOPS

PCIe bus (x16: Gen4: 64GB/s, Gen 5: 128GB/s

GPU
O(10) TFLOPS

O(1) TB/s

HBM2
(>=32 GiB)

GPU
O(10) TFLOPS

O(1) TB/s

HBM2
(>=32 GiB)

Internal
Network:
O(100)
GB/s

GPU
O(10) TFLOPS

O(1) TB/s

HBM2
(>=32 GiB)

GPU
O(10) TFLOPS

O(1) TB/s

HBM2
(>=32 GiB)

**OAK RIDGE**
National Laboratory | LEADERSHIP COMPUTING FACILITY
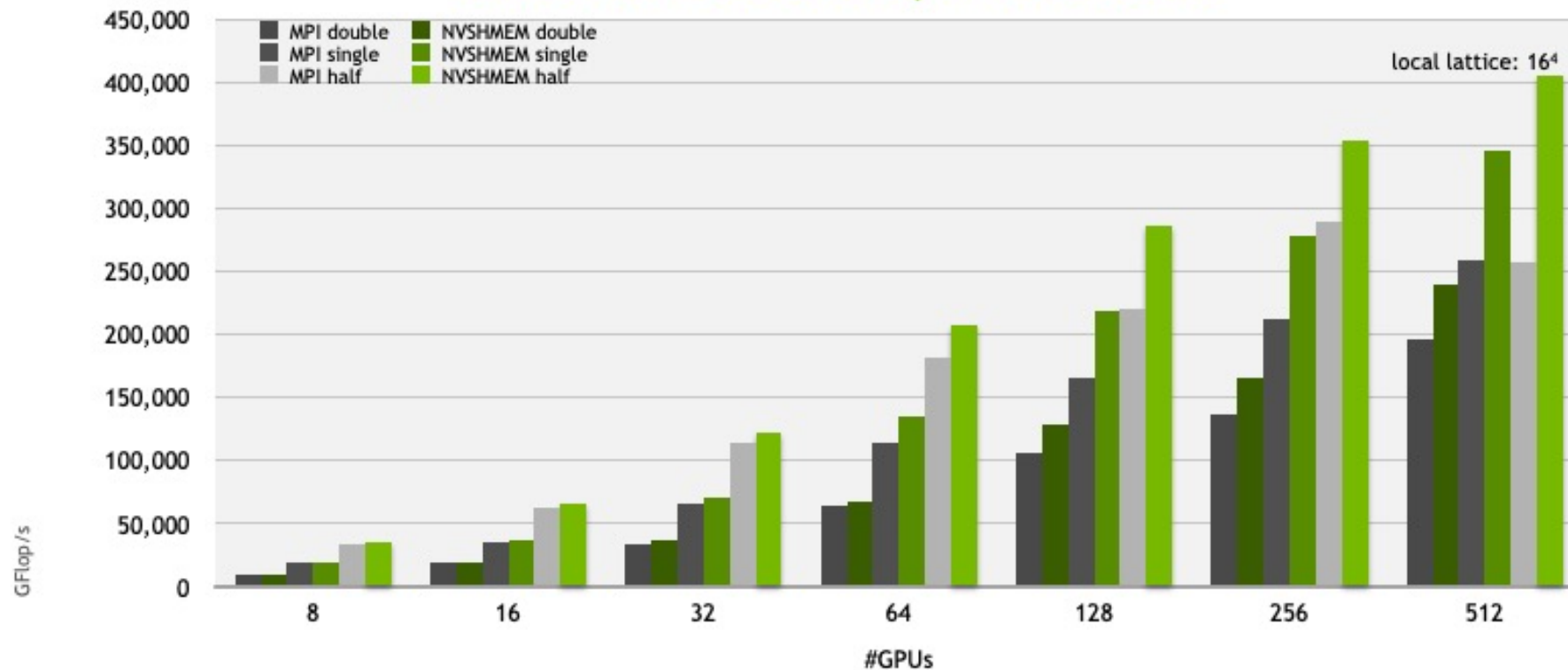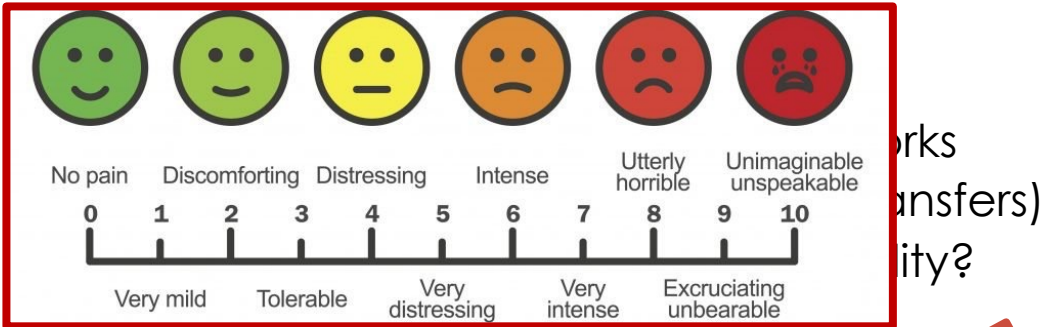
Open slide master to edit

# Latencies affect strong scaling



Figure Courtesy of M. Wagner, NVIDIA (Talk on Wednesday afternoon in Machines and Software session in this conference)

OAK RIDGE National Laboratory | LEADERSHIP COMPUTING FACILITY

Open slide master to edit

# And so, Nuances matter

- Hardware capability



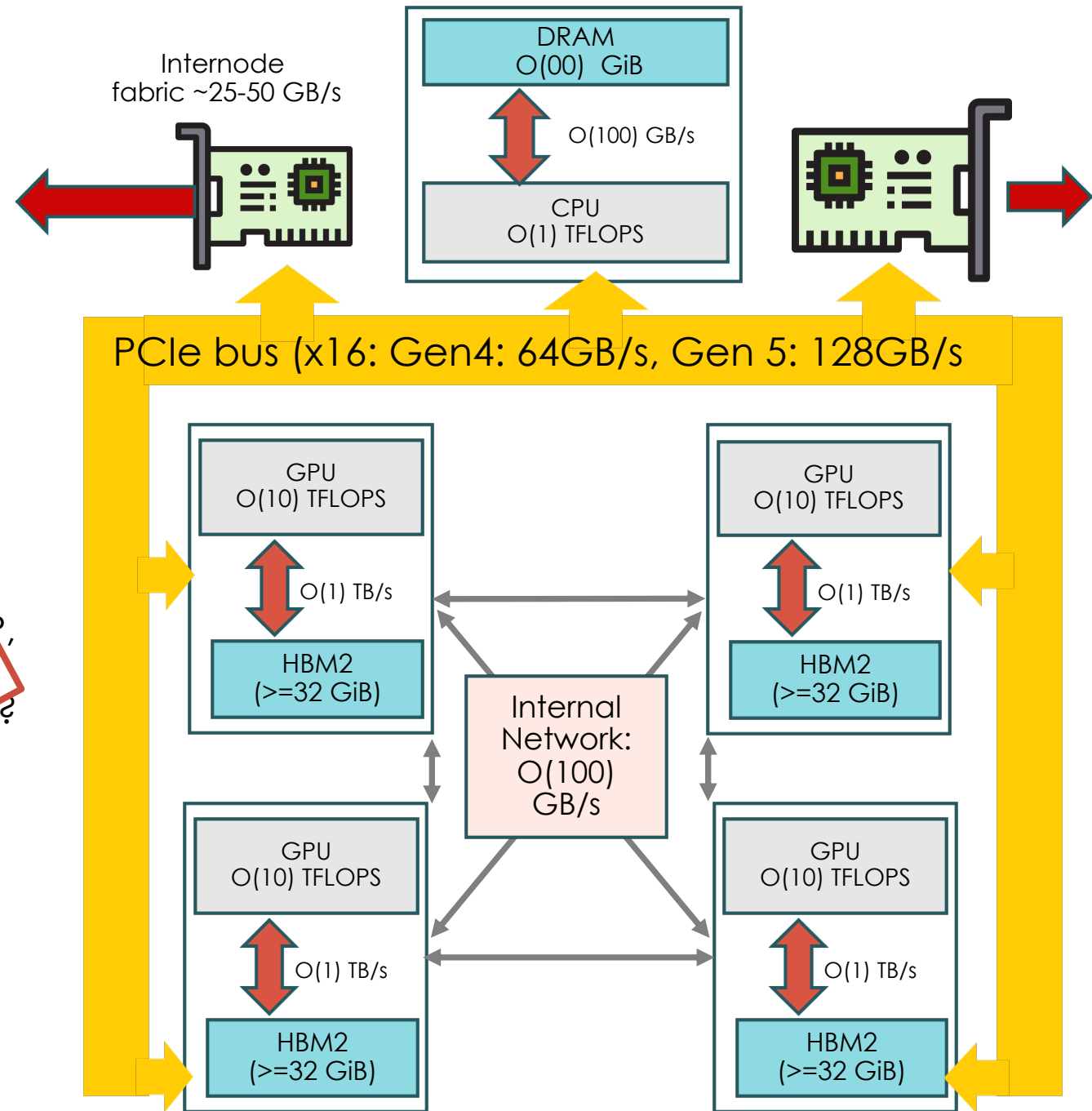| No pain | Discomforting | Distressing | Intense | Utterly horrible | Unimaginable unspeakable |
|---|---|---|---|---|---|
| 0 | 1 2 | 3 4 | 5 6 | 7 8 | 9 10 |
| | Very mild | Tolerable | Very distressing | Very intense | Excruciating unbearable |

works
...nsfers)
...ity?

- Programming model & implementation
  - Kernel launch latency
  - User access to Hardware c...   ...P, SHMEM etc)
  - Persistent kernels, s...
  - Access to Unifi...

- Algorithm...
  - M...   ...-> latency
  - ...gorithms -> B/W
  - ...device locality

**YOUR LEVEL OF PAIN MAY VARY !!!**

DRAM
O(00) GiB

O(100) GB/s

CPU
O(1) TFLOPS

Internode fabric ~25-50 GB/s

PCIe bus (x16: Gen4: 64GB/s, Gen 5: 128GB/s

GPU
O(10) TFLOPS

O(1) TB/s

HBM2
(>=32 GiB)

GPU
O(10) TFLOPS

O(1) TB/s

HBM2
(>=32 GiB)

Internal Network:
O(100) GB/s

GPU
O(10) TFLOPS

O(1) TB/s

HBM2
(>=32 GiB)

GPU
O(10) TFLOPS

O(1) TB/s

HBM2
(>=32 GiB)

**OAK RIDGE**
National Laboratory | LEADERSHIP COMPUTING FACILITY
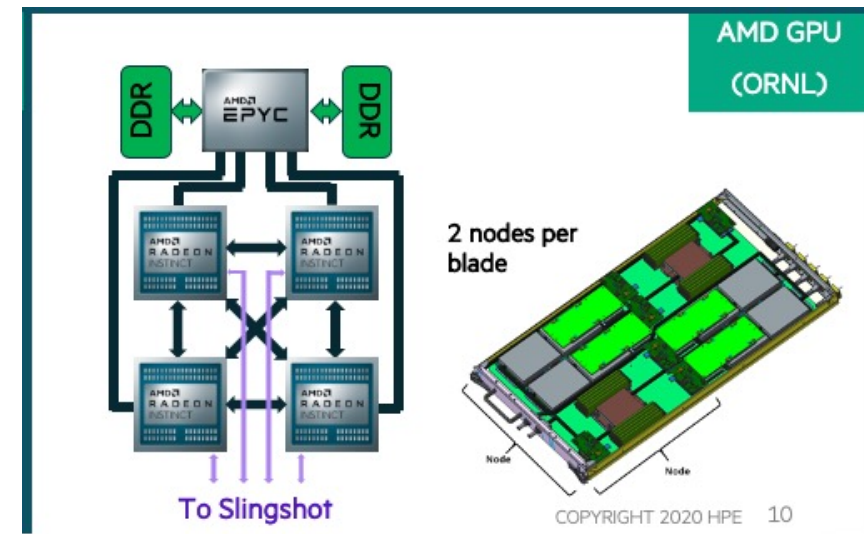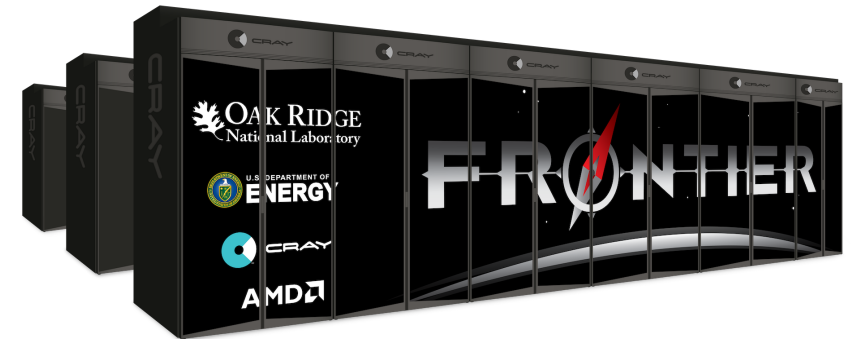
# US Exascale and Pre-Exascale Systems: NERSC Perlmutter

- Perlmutter is the New HPE/Cray "Shasta" system recently arrived at NERSC
  - #5 on Top 500 List: 89.8 PF (Rpeak), 64.5 PF (Rmax)

- Accelerators are NVIDIA A100 (Ampere) GPUs
  - Memory B/w: 1555.2 GB/sec
  - FLOPs:
    - 19.5 TF (FP32)   9.7 TF ( FP64 )
    - Tensor ops: 311.9 (FP16), 155.9 (TF32), 19.5 (FP64)

- NVLink 3
  - 4 GPUs on node have all-to-all connection
    - 100 GB/s/direction between a pair of GPUs
    - 600 GB/sec total into/out-of a GPU

- HPE/Cray Slingshot-10 Interconnect

- System is expected to open to users soon

- Programming Models
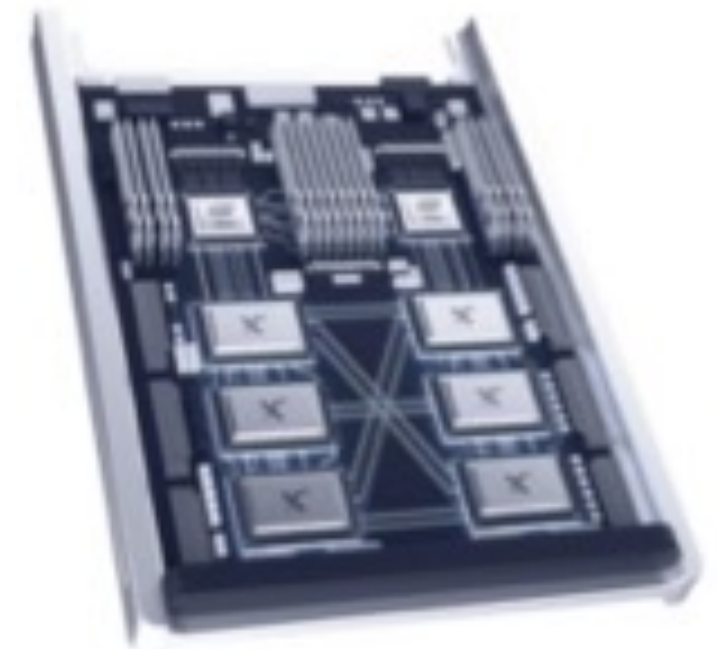  - MPI + CUDA, OpenMP-5 offload, NVSHMEM

- https://www.nersc.gov/systems/perlmutter/

OAK RIDGE
National Laboratory | LEADERSHIP COMPUTING FACILITY

Open slide master to edit

# US Exascale and Pre-Exascale Systems: OLCF Frontier

- Frontier will be the New HPE/Cray "Shasta" system at Oak Ridge Leadership Computing Facility
  - \>1.5 Exaflops peak performance

- Accelerators will be AMD Radeon Instinct GPUs

- AMD Infinity Fabric
  - 4 GPUs on node to have all-to-all connection

- HPE/Cray Slingshot
  - "Multiple Slingshot NICs providing 100GB/sec network bandwidth"
  - GPUs will have direct connection to Slingshot

- Programming Models
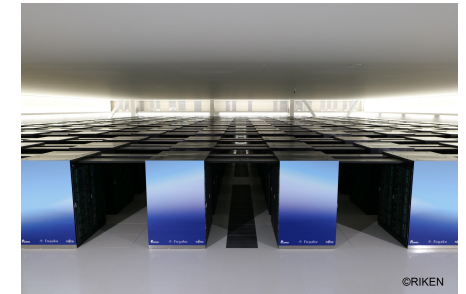  - MPI, AMD HIP, OpenMP-5-offload, ...

- https://www.olcf.ornl.gov/frontier

**OAK RIDGE** National Laboratory | LEADERSHIP COMPUTING FACILITY

Open slide master to edit

# US Exascale and Pre-Exascale Systems: ALCF Aurora

- Aurora will be the New HPE/Cray "Shasta" system at Argonne Leadership Computing Facility
  - >1 Exaflop peak (DP) performance

- Accelerators will be Intel Xe architecture based "Ponte Vecchio" GPUs

- Unified Memory Architecture
  - Across GPU and CPU

- Low Latency, High Bandwidth all-to-all connectivity within Node

- HPE/Cray Slingshot interconnect
  - 8 fabric endpoints per node

- Programming Models
  - MPI, Intel oneAPI DPC++ (based on SYCL), OpenMP-5-offload, ...

- https://alcf.anl.gov/aurora

**OAK RIDGE** National Laboratory | LEADERSHIP COMPUTING FACILITY

Open slide master to edit

# Other Noteworthy Systems for LQCD

- ## Supercomputer Fugaku, RIKEN CCS Japan
  - \- #1 Top 500 list: 537PF (Rpeak), 442 PF (Rmax)
  - – Please see talk by Yoshifumi Nakamura!

- ## Summit at OLCF, U.S.A
  - – #2 Top 500 list: 200PF (Rpeak), 148.6 PF (RMax)

- ## Sunway TaihuLight, National Supercomputing Center, Wuxi, China
  - – #4 on on Top 500 list: 125.4 PF(Rpeak) 93PF(Rmax)

- ## Tianhe-2A, National Super Computer Center, Guangzhou, China
  - – #7 on the Top 500 list; 100.6 PF (Rpeak) 61PF (Rmax)

- ## JUWELS Booster, Forschungszentrum Jülich, Germany
  - – #8 on Top500 List: 71PF (Rpeak), 44 PF (Rmax)
  - – 4 NVIDIA A100 GPUs, switched directly to HCAs
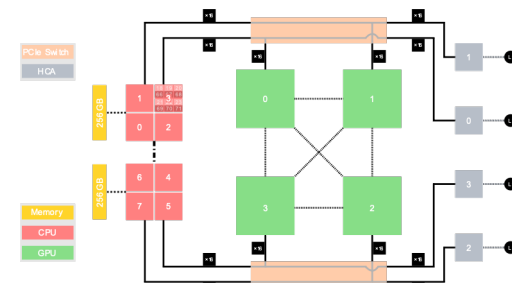  - – Infiniband Network in DragonFly+ configuration.



*Supercomputer Fugaku*



*OLCF Summit*



*TaihuLight*



*JUWELS Booster Node Architecture*

**OAK RIDGE** National Laboratory | LEADERSHIP COMPUTING FACILITY
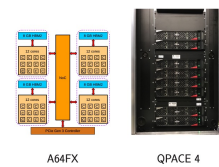
Open slide master to edit

# Custom QCD Systems?

- QPACE 4: 177/354 TFlop/s peak
  - Please see poster by the Regensburg group: N. Meyer, P.Georg, S. Solbrig & T. Wettig

- A64FX CPUs - Fugaku like
  - 48 cores/CPU 1.8 GHz
  - 512-bit Arm SVE SIMD: 2.76 TF/cpu
  - 32 GB HBM2 memory per CPU

- InfiniBand EDR interconnect

- Open-source software stack
  - no assistant cores like Fugaku
  - lots of kernel tuning to reduce O/S noise.

Open slide master to edit

# Software

OAK RIDGE
National Laboratory | LEADERSHIP COMPUTING FACILITY

Open slide master to edit

# Node Programming Models: Pick your standard

- Vendor options
    - NVIDIA CUDA: the CUDA you know, evolved..
        - Support for Tensor cores in cuBLAS, cuDNN and CUDA (wmma:: namespace)
        - ecosystem support: cuBLAS, cuFFT, cuRAND etc.
        - C++ Standard Library ( cuda::std:: )
        - Strong push to implement **ISO C++ standard** features in full.
    - **AMD HIP**
        - very similar to NVIDIA's CUDA core library ( replace 'cuda' with 'hip' )
        - **open source implementations for ROCm, also for CPUs** (HIP CPU)'
        - LLVM/Clang based compiler
        - ecosystem support: hipBLAS, hipFFT/rocFFT, hipRAND/rocRAND, etc
        - access to CDNA features through compiler
    - Intel oneAPI
        - based on **SYCL Standard from the Khronos Group, with OpenCL heritage**
        - Compiler in oneAPI SDK and in Intel's contributions to LLVM
        - ecosystem support: oneMKL, oneDNN, ...
        - oneAPI Tutorial at this conference by P. Steinbrecher

- Directive based approaches: **OpenMP-offload and OpenACC**
    - OpenMP: `#pragma omp teams distribute parallel for simd`
    - all vendors committed to supporting OpenMP-5.0 offload
    - OpenACC supported in NVIDIA HPC-SDK and GCC (work by Mentor Graphics)
    - OpenMP is likely the best option for Fortran  ( although nvfortran can offload DO CONCURRENT )

- **OpenCL** is available  through all vendors as well
    - but beware vendor specificity in the implementations...

OAK RIDGE National Laboratory | LEADERSHIP COMPUTING FACILITY

Open slide master to edit

# Performance Portability via Kokkos and RAJA

- Kokkos and RAJA provide portability to C++ programs via 'back ends'
  - Kokkos now has a Fortran interface project

- Both prog. models provide basic parallel patterns
  - forall, reduction, scan

- Kokkos provides 'View' datatype
  - multi-dimensional array, a precursor to std::mdspan
  - RAJA uses other packages (e.g. Umpire) for memory management

- Views and execution patterns are bound to memory and execution spaces via policies.
  - Policies have default values, which can be set in some architecture specific header.



Kokkos & RAJA Abstractions → CUDA Back-End, OpenMP Back-end, HIP Back-end, SYCL DPC++ Back-end, OpenMP target Back-End

```
Kokkos::View<float[1000]> y("y"), x("x");
// Fill x and y somehow – not shown

float alpha=0.5;
Kokkos::parallel_for(1000,
            KOKKOS_LAMBDA(int i) {
        y(i) += alpha * x(i);
});
```

**OAK RIDGE** National Laboratory | LEADERSHIP COMPUTING FACILITY

Open slide master to edit

# Portability of the programming Models

- Directive based approaches: should target all main systems
  - require strong compiler support, and vendors to agree on interpretation of the standard

- Intel Data Parallel C++: main target: Intel GPUs/CPUs/FPGAs etc.
  - Aurora Proxies: Linux laptops with Intel Graphics (Gen 9 or later), Intel Dev Cloud, JLSE  Early access systems
  - NVIDIA GPUs: Intel LLVM compiler has CUDA back end from Codeplay
  - AMD GPUS:  hipSYCL Project, a pilot AMD Back End for Intel LLVM is being developed by Codeplay

- HIP: main target AMD GPUs
  - Frontier Proxies: NVIDIA GPUs systems (Summit), Early Access Systems (Spock)
  - NVIDIA GPUs: Native support up to a certain CUDA level (substitutions: 'hip→'cuda')
  - Intel GPUs: ECP HIP-LZ (HIP on Level Zero) project
  - nascent HIP-CPU project to support CPUs (as a development aid)

- CUDA: main target NVIDIA GPUs
  - Perlmutter Proxies: Summit, JUWELS-Booster, Desktop GPU systems,
  - AMD: HIP includes 'hipify' conversion tool for 'developer assisted' conversion to HIP.
  - Intel GPUs: Intel provides 'developer assisted' conversion to Data Parallel C++

- Kokkos & RAJA:  pick most performant back end (most likely the native one)



Comparing SYCL Options on V100 against QUDA, Kokkos+CUDA

OAK RIDGE National Laboratory | LEADERSHIP COMPUTING FACILITY

# Code porting/development efforts: QUDA

- QUDA is a software library for QCD components
  - Advanced Optimized solvers (e.g. Aggregation based MG preconditioner)
  - Key optimized routines: GaugeForce (MILC), Gauge Fixing,
  - Multiple precision use since Day-1
  - Built in Autotuning
    - adjusting grid, block and shared memory sizes,
    - choosing communications policy e.g. P2P, GDR, MPI

- See "Towards QUDA 2.0" talk by K. Clark (Wed, 1pm) for refactoring details
  - plus several other QUDA related talks at this conference (M. Wagner, E. Weinberg, J. Tu, ...)

- QUDA grew a 'performance portability' subgroup of developers:
  - K. Clark (NVIDIA) QUDA Lead developer, C++, std::par, pSTL, NVC++
  - B. Joo (OLCF) HIP Backend
  - D. Howarth (LLNL) HIP Backend
  - J. Osborn (ALCF) DPC++ Backend
  - X-Y. Jin (ALCF) OpenMP-offload Backend
  - A. Strelchenko (FNAL) – Intel CUDA-Conversion tool exploration
  - D. McDougall, C. Robeck (AMD) – HIP consultation
  - P. Steinbrecher (Intel) – Intel compiler consultation (DPC++ and OpenMP)
  - weekly meetings + notes of varying level of depth in our Slack #portability channel

OAK RIDGE
National Laboratory | LEADERSHIP COMPUTING FACILITY

ECP
EXASCALE COMPUTING PROJECT

# Refactoring for Perf Portability (led by K. Clark)

- Launch Abstractions
  - Kernel1D, Kernel2D, Kernel3D, Reduction, TransformReduce,
  - most portability frameworks do something like this: Kokkos, Raja, and also Grid

- API Elements
  - memory allocation/copy ( quda_malloc() .... )
  - Streams (fixed number and use index as 'stream')
  - Shared memory: sharedMemory manager class
  - Constant memory – for large list of arguments (e.g. small matrices)
  - API functions error check internally (architecture dependent status codes are isolated/contained)

- Device constraints & Macro Reduction/Localization
  - constexpr functions
  - SFINAE, and template-template classes.

- Libraries
  - CUB/hipCUB,     => namespace QudaCub aliases
  - cu/hipRand, cu/hipFFT => wrapped into 'shim headers'
  - certain math functions ( sincos and rsqrt and half prec implementations ) implemented on host/device
  - straightforward for HIP, since many libraries corresponding to CUDA exist, may be tricker for e.g. Intel

- Miscellany
  - vector types (e.g. float4, int4) – HIP returned proxy objects / accessors to some structure elements
  - IPC for Peer2Peer seems easy to port between CUDA and HIP

**CONSTEXPR**



**ALL THE THINGS!**

CppCon17 talk by Ben Deane and Jason Turner

**CONSTEXPR**



Also "Don't constexpr all the things" – David Sankel, CppNow 2021

(actually this talk is about Circle)

OAK RIDGE National Laboratory | LEADERSHIP COMPUTING FACILITY

Open slide master to edit

ECP EXASCALE COMPUTING PROJECT

# Spock: OLCF's Frontier Early Access System (EAS)

- 3 Cabinets w. 12 nodes/cabinet => 36 nodes, 144 GPUs total

- Each node has
  - 1x 64 core AMD EPYC 7662 "Rome" CPU (4 NUMA domains) + 256 GB DDR4 memory (205 GB/s)
  - 4x AMD Radeon Instinct MI100 GPUs (gfx908), 32 GB HBM2 @ (1.2 TB/sec)
    - GPU-GPU: All-to-all Infinity Fabric Interconnect,   Host-GPU: PCie Gen4: 32+32 GB/sec
  - 2 x NVMe SSDs (3.2 TB each): 6800 MB/sec read, 4200 MB.sec write

- Slingshot 10 Interconnect: 12.5 + 12.5 GB/sec

- Spock Training  Workshop materials: https://www.olcf.ornl.gov/spock-training

- Documentation: https://docs.olcf.ornl.gov/systems/spock_quick_start_guide.html

*Original Spock Compute Node*

# QUDA Porting Status

- (All of) the library builds with ROCm 4.2
  - feature/hip-compile-fixes branch
  - P2P & MPI Comms built. Currently undergoing testing and debugging

- All 'Ctest' pass on single AMD GPU

- Intel GPU porting in progress
  - James Osborn is still developing and testing the SYCL port and Xioyong-Jin is working on OpenMP-offload
  - focus is on making sure the refactorings which we discussed allow efficient porting to both programming models (access to shared memories, reductions, etc)

- *A100 data courtesy of E. Weinberg, NVIDIA. Data from **A100 80 GB** part*
  - *~**2.0TB/s** peak mem. B/W*
  - *Perlmutter has 1.56 TB/sec, 40GB parts*

- *V100 data from Summit*
  - *~**900** GB/sec peak mem B/W*

- *MI100 data from Spock*
  - *~**1.2TB**/sec peak mem B/W*

QUDA Dslash Test, V=$32^4$ sites, single device, wilson (scalar build)

AMD DATA PRELIMINARY

OAK RIDGE National Laboratory | LEADERSHIP COMPUTING FACILITY

Open slide master to edit

# Progress with QDP-JIT

- Frank Winter ported QDP-JIT to AMD GPUs: Presented this conference Wed 7/28

- QDP-JIT can run a full Chroma HMC trajectory on AMD GPU
    - 800+ GFLOPS from Native QDP-JIT Dslash implementation, 800+ GB/sec on SU(3)xSU(3) kernels
    - MI100 has ~1200 GB/sec maximum HBM bandwidth

- AMDGPU backend for LLVM generates ISA rather than an intermediate form like PTX
    - Needs additional linking step

- Intel GPUs next...

C++ Source code

```
LatticeFermion x;
gaussian(x);
```

Compile time Expression Template Magic

**Executable**

Code generator for gaussian() → IR for gaussian()

IR for gaussian() → LLVM CodeGen

LLVM CodeGen → PTX / AMD ISA

PTX / AMD ISA → LLVM Linker

LLVM Linker → load/link and cache

load/link and cache → dispatch()



DSlash on AMD gfx908 (ROCm 4.1)
f32 / f32 tuned / f64 / f64 tuned
Gflops vs Lattice size, $L^4$



SU3 x SU3 on AMD gfx908 (ROCm 4.1)
f32 / f32 tuned / f64 / f64 tuned
GB/s vs Lattice size, $L^4$

*Preliminary MI100 results from the ECP All Hands Meeting poster by F. Winter. Courtesy of Frank Winter.*

OAK RIDGE National Laboratory | LEADERSHIP COMPUTING FACILITY

Open slide master to edit

# Julia For LQCD

- Interesting work from Akio Tomiya -- poster at this conference.

- Julia uses JIT compilation via LLVM.

- "Very high level" with good C/C++/Python interoperability – ML friendly

- Support for Wilson-Clover & Staggered Fermions, Stout Smearing, RHMC, Heat-Bath, Self Learning Monte Carlo,

- Similar performance to Fortran code

- Related:
  - work by A. Strelchenko (2019 USQCD All Hands Meeting) for algorithimic epxloration

## LatticeQCD.jl : Lattice QCD code with Julia

Akio Tomiya (International professional university of technology in Osaka, Assistant Professor) akio@yukawa.kyoto-u.ac.jp
Yuki Nagai (Japan Atomic Energy Agency, Senior Scientist)

### 1. Introduction (1)

- Lattice QCD = Multi-dimensional integral over SU(3)

$$S[U,\psi,\bar{\psi}] = \sum_n \left[ -\frac{1}{g^2} \mathrm{Re\ tr}\ U_{\mu\nu} + \bar{\psi}(\slashed{D}+m)\psi \right]$$

$$\langle \mathcal{O} \rangle = \frac{1}{Z} \int \mathcal{D}U \mathcal{D}\bar{\psi}\mathcal{D}\psi e^{-S[U,\psi,\bar{\psi}]} \mathcal{O}[U,\psi,\bar{\psi}]$$

$$\mathcal{D}U = \prod_{n\in(\mathbb{Z}/L)^4} \prod_{\mu=1}^{4} dU_\mu(n) \quad >1000\ \text{dimension.}$$

- This integral gives non-perturbative information of QCD

- Monte-Carlo is used to calculate (Numerical error independent to the dimensionality!)

- C++/Fortran have been used for simulations since it costs a lot! Supercomputers are needed for large scale calculations

- We make an open source code for lattice QCD with **Julia language**!

### 2. Julia? (2)

- Programing language for science (Ref. 1) since 2012. Free, open

- Fast as C/Fortran (Fig1), productive as Python (Fig2)

Fig1.
N=2000, product of 2 NxN random matrices (Not BLAS)
Fortran (gfortran -O3)
Julia 1.6.1

Fig2.
```
for i in [1,2,3]
    println("Hello world ", i)
end
```

- 33.9k star on Github. NASA uses Julia [2]. Runnable on Supercomputers

- Easy start: Binary available for Win, macOS and Linux, run everywhere

- Good package control system unlike python environment.

- Just-in-compiling, dynamic type. We can use Python/Fortran/C libraries. **Machine learning friendly!**

### 3. Why we make? (3)

- To examine capability of Julia

- Ease of install/compiling

- Machine learning friendly lattice QCD code is needed

- Educational purpose/ Ease of modification

### 4 .USAGE   Only 4 steps! See our Github webpage in details

1. Download Julia binary from the official webpage
2. Add lattice QCD using built-in package control system
3. (optional) Make parameter file with the wizard (type `run_wizard()`)
4. Execute! (type `run_LQCD("my_parameters.jl")`)

**We also provide Google Colab notebook [4]!**

### 5. Features (4)

- General gauge action (plaq+rect+chair +…) for SU(N) is supported

- Dynamical clover-Wilson (Nf=2), staggered fermions (Nf = 2-8). Both can be run with/without stout.

- (R)HMC, Heatbath (for quenched), self-learning Monte-Carlo, etc are supported

- Measurements: Plaquette, Polyakov loop, Chiral condensate, Pion correlator, topological charge

- Gradient flow with general gauge action

- ILDG I/O support

- Work on Google colab/ batch job / REPL (Julia prompt)

- (parallelization is in progress)

- Parameter wizard

### 6. Benchmark (5)

Machine:
m1 mac mini
- Julia 1.6.1 + Rosetta2
- gfortran11 (with/wo O1)

L=4^3 x Lt
Lt = 4, 8, 10, 12, 16, 20
kappa = 0.141139
beta = 5.5
Nmd = 10, ε=0.1
CG eps = 10-8
(Default of Lattice tool kit [3])

- We compare with Lattice Tool kit (Fortran) , (same algorithm)

- Ls=4, Lt=4-20, beta = 5.5, kappa = 0.141139, full HMC

- Performance is good so far on single thread/core

### 7. Summary (6)

- Julia is fast as Fortran/C, productive as Python (easy to write)

- LatticeQCD.jl works well, fast as a fortran code

- Future work: Overlap, domain-wall, parallelization

- (Modified version of) this code used for arXiv 2010.11900 and arXiv 2103.11965. Talk in session 27th, 13:00-, Algorithms

### Reference

1. Julia: https://julialang.org
2. LTK: https://nio-mon.riise.hiroshima-u.ac.jp/LTK/
3. NASA: https://modelingguru.nasa.gov/docs/DOC-2783
4. Google Colab https://bit.ly/3yytQiG

LATTICE2021@MIT (Virtual)    Start Lattice calculation in 5 min  https://github.com/akio-tomiya/LatticeQCD.jl    Abstract: https://indico.cern.ch/event/1006302/contributions/4378500/

OAK RIDGE National Laboratory | LEADERSHIP COMPUTING FACILITY

Open slide master to edit

# Grid

- Grid is an LQCD framework headed up by P. Boyle and team
  - several Grid related talks/posters at this conference
  - Poster by P. Boyle, A. Yamaguchi, GPT Python bindings tutorial by C. Lehner, etc....

- C++ expression template based LQCD software system

- Portability to several "back-ends"
  - OpenMP
  - HIP
  - CUDA
  - SYCL

- coalescedRead()/Write() – abstraction

- accelerator_for() – abstraction (macro?)
  - wraps actual kernel launch via HIP/SYCL/CUDA

- for more details please contact the Authors

- excellent performance, including near wire-speed saturation of networks reported

OAK RIDGE National Laboratory | LEADERSHIP COMPUTING FACILITY



## Grid: OneCode and FourAPIs
### www.github.com/paboyle/Grid

**Azusa Yamaguchi** *(University of Edinburgh)*
**Peter Boyle** (*Brookhaven National Laboratory*)

**Background:**
Grid is a C++11 high level library for lattice Gauge theory[1,2]
It aims to be performance portable across all modern architectures
A number of LQCD software efforts make use of Grid for actions, algorithms (solvers, multigrid, HMC, contraction primitives) [3].
Portability in across the Exascale roadmap requires support for AMD (HIP), Intel (SYCL) and Nvidia (QUDA) GPUs in addition to vectorising multicore architectures. Grid is a single source and targets all of these APIs portably.

Frontier AMD CPU, AMD GPU; HIP  Perlmutter AMD CPU, Nvidia GPU; CUDA

Aurora Intel CPU, Intel GPU; SYCL  Summit IBM CPU, Nvidia GPU; CUDA

+ CPU computing is also *not* going away

*Covariant* programming:

**SIMD and SIMT differ semantically in whether local variables are vectors or scalars**

*Naively* poses a barrier to writing single source kernels which vectorise on CPUs and read coalesce on GPUs.

C++ automatic type inference lets you avoid naming the types (vector or scalar) so you can deduce the type according to the architecture following simply programming idioms.

Combined with *accelerator_for* abstraction capturing and offloading loop bodies in *device lambda functions* we can write high performance kernel code that runs on all four API's.

**FourAPIs and OneCode**

Portability 101 – abstract the interfaces    Performance Portability 102 – abstract the layout

Similar ideas to RAJA and Kokkos – use device lambda capture ; lean internal interface to offload - HIP and OpenMP similar

**HIP and OpenMP targets**

### Wilson dslash kernel (sketch)

**Same optimised kernel transform**s covariantly between GPU and CPU

Return type of coalescedRead dictates whether SOA or scalar structs are processed in each logical thread

*Thread interleaving happens naturally on GPU with memory resident data stored in SOA.*

*CPU processes loop as SOA data with good vectorisation.*

### Nvidia performance:

Excellent performance on ATOS sequana A100 x 4 nodes with 4x HDR infiniband *(e.g. Juelich Booster + Edinburgh Tursa systems)* – communication is key
6TF/s per node in multi-node operation.

90 % of wirespeed delivered to application
Aggregate MPI bidirectional bandwidth per node (GB/s) on 16 nodes
Wirespeed = 200 Gb/s



| Nodes | GPUs | Measured Perf / GPU | Measured TF/s | Ideal GPU scaling | Ideal Node scaling |
|---|---|---|---|---|---|
| 1 | 1 | 3.85 | 3.85 | 3.85 | |
| 1 | 4 | 3.075 | 12.3 | 15.4 | 12.3 |
| 16 | 64 | 1.671875 | 107 | 246.4 | 196.8 |
| 64 | 256 | 1.484375 | 380 | 985.6 | 787.2 |

Nsight compute indicates on A100-40
- 82% of L2 cache saturation,
- 76% of HBM saturation
- 36% FMA pipeline usage

**Intel / Aurora**
- Saturates memory bandwidth on both Iris XE max (DG1) and Arctic Sound
- AMD/Frontier – code ports and run, but performance is a work in progress.
- Fugaku port by Nils Meyer / Regensburg

| Device | Fp32 Dw GF/s | Memory BW GB/s |
|---|---|---|
| DG1 | 170-201 GF/s | 58 GB/s |
| V100 | 1750 GF/s | 850GB/s |

### Conclusions:
Supercomputing companies are not helping scientific productivity with a proliferation of programming models.
**There is hidden commonality as they are all based on vector computing**
The semantic differences between CPU and GPU can be abstracted and high performance portable source written
Other codes may wish to adopt these techniques

1. Grid: A next generation data parallel C++ QCD library : arXiv:1512.03487
2. Performance Portability Strategies for Grid C++ Expression Templates : arXiv:1710.09409
3. Hadrons: https://github.com/aportelli/Hadrons
   GPT: https://github.com/lehner/gpt
   CPS: https://github.com/RBC-UKQCD/CPS
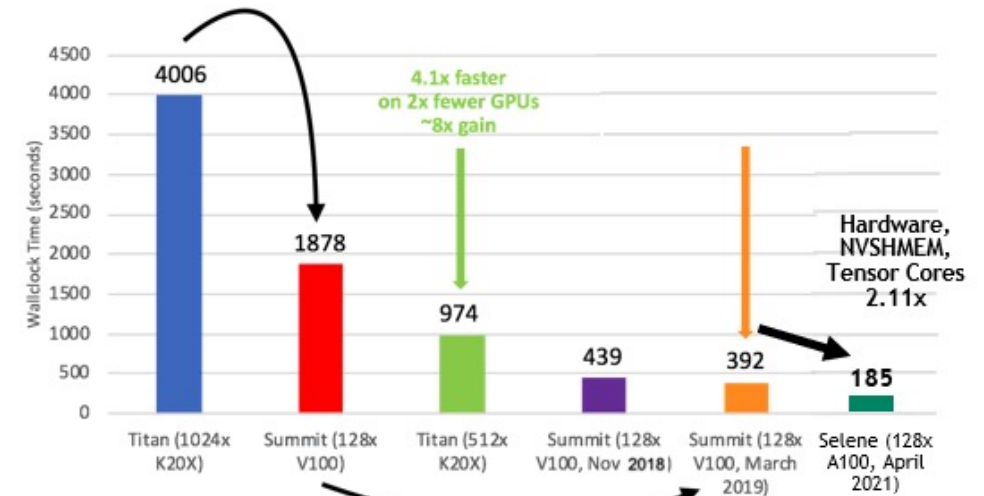   MILC: http://www.physics.utah.edu/~detar/milc/

# Summary

- Very exciting time as code porting work is coming to fruition
  - Perlmutter is here, we are looking forward to hopping on...
  - Frontier will be starting to arrive very soon
  - Slightly stressful: tooling is still maturing in some cases, bugs to iron out

- Some interesting new machines on the horizon
  - Aurora at ALCF
  - Large A64FX based system by Atos with GENCI and CEA in France
  - A system based on the Grace CPU from NVIDIA at Los Alamos

- With so many architectures it is important to have a good performance portability strategy

- Algorithms also advancing – sadly not enough time to cover in this talk
  - Multigrid algorithms advances (E. Weinberg, P. Boyle, ...)
  - Advances in Domain Wall algorithms: Multi-Grid and DDHMC (talk by Peter Boyle, Chulwoo Jung and others)
  - Algorithms to combat critical slowing down and adding machine learning to accelerate Monte-Carlo methods. (talks by Sam Foreman & Xiao-Yong Jin)
  - and many more...

**OAK RIDGE** | LEADERSHIP
National Laboratory | COMPUTING
FACILITY

Open slide master to edit

# Departing Plot

- In 2018 we measured a baseline of our HMC gauge generation algorithm on 1024 nodes of Titan at OLCF
  - No multigrid yet, but GPU accelerated – K20X GPUs

- On this figure we can track improvement of algorithms, software and algorithm tuning

- At the current time what took 4006 sec. on 1024 Titan nodes, takes only 185 sec. on the NVIDIA Selene System (#6 on Top500 system)
  - 21.6x speedup on 8x fewer devices
  - 173x "integrated improvement"

- Looking forward to adding some Frontier data on this plot soon

- Where will new advances take us next?



Hardware: 2.13x wall-time on 8x fewer GPUs = 17x

4.1x faster on 2x fewer GPUs ~8x gain

Hardware, NVSHMEM, Tensor Cores 2.11x

4006 — Titan (1024x K20X)
1878 — Summit (128x V100)
974 — Titan (512x K20X)
439 — Summit (128x V100, Nov 2018)
392 — Summit (128x V100, March 2019)
185 — Selene (128x A100, April 2021)

Algorithms, Software and Tuning: 4.79x

Chroma w/ QDP-JIT and QUDA, ECP FOM data, $V=64^3 \times 128$ sites, $m_\pi \sim 172$ MeV, (QDP-JIT by F. Winter, Jefferson Lab)

Original figure credit Balint Joo ~2 years ago, new numbers from M. Wagner

OAK RIDGE National Laboratory | LEADERSHIP COMPUTING FACILITY

Open slide master to edit

# Contribution and Collaboration Acknowledgements

- I would like to express my thanks to everyone in the QUDA Portability working group: K. Clark, D. Howarth, X-Y Jin, D. McDougall, J. Osborn, C. Robeck, P. Steinbrecher, A. Strelchenko

- Collaborators in the previous Kokkos and SYCL work (C. Trott, D. Sunderland, D. Ibanez, N. Liber, J. DeSlippe, T. Kurth, J. Kim, P. Steinbrecher, D. Doerfler)

- K. Clark, E. Weinberg, J, Tu, M. Wagner, NVIDIA for long time collaboration and specifically providing figures and helpful comments

- Corbin Robeck, Damon McDougall, Nick Curtis, Rene van Oostrum  from AMD at the Frontier Center Of Excellence, and Leopold Ginberg from the AMD Business Unit for detailed collaboration on porting and optimizing on AMD GPUs.

- F. Winter for use ot his QDP-JIT plot

- P. Boyle for their Grid Poster and details of Domain Wall algorithms,

- N. Meyer for information about QPACE 4, and the QPACE 4 Poster,

- Sam Foreman, and Xiao-Yong Jin for the l2hmc and ML Transformation based HMC algorithms

- Photographs of computers courtesy of their institutions and the internet.

- Use of computer systems for data generated here, including Summit and Spock at Oak Ridge Leadership Computing Facility (OLCF), the Tulip System operated by HPE/Cray,   Inte Gen-9 workstation at Jefferson Lab

- Philip Roth for suggesting the Hospital Pain Chart on Slide 8.

**OAK RIDGE** | LEADERSHIP COMPUTING FACILITY
National Laboratory

Open slide master to edit

# Funding Acknowledgements

**OAK RIDGE** National Laboratory | LEADERSHIP COMPUTING FACILITY

Open slide master to edit