

# LATTICE QCD ON MODERN GPU SYSTEMS

Mathias Wagner

### **STEPS IN AN LQCD CALCULATION**

 Generate an ensemble of gluon field configurations "gauge generation" Produced in sequence, with hundreds needed per ensemble Strong scaling required with 100-1000 TFLOPS sustained for several months 50-90% of the runtime is in the linear solver
 O(1) solve per linear system Target 16<sup>4</sup> per GPU

2. "Analyze" the configurations

Can be farmed out, assuming ~10 TFLOPS per job Task parallelism means that clusters reign supreme here 80-99% of the runtime is in the linear solver Many solves per system, e.g., O(10<sup>6</sup>) Target 24<sup>4</sup>-32<sup>4</sup> per GPU  $D_{ij}^{\alpha\beta}(x,y;U)\psi_{j}^{\beta}(y) = \eta_{i}^{\alpha}(x)$ or Ax = b

Simulation Cost ~ a<sup>-6</sup> V<sup>5/4</sup>

#### MAPPING THE DIRAC OPERATOR TO CUDA

Finite difference operator in LQCD is known as Dslash Assign a single space-time point to each thread

V = XYZT threads, e.g., V =  $24^4$  =>  $3.3 \times 10^6$  threads

Looping over direction each thread must

Load the neighboring spinor (24 numbers x8)

Load the color matrix connecting the sites (18 numbers x8)

Do the computation

Save the result (24 numbers)

Each thread has (Wilson Dslash) 0.92 naive arithmetic intensity



X[0]



### QUDA

- Effort started at Boston University in 2008, now in wide use as the GPU backend for BQCD, Chroma, CPS, MILC, TIFR, tmLQCD, etc.
- Provides:

Various solvers for all major fermionic discretizations, with multi-GPU support Additional performance-critical routines needed for gauge-field generation

- Maximize performance
  - Exploit physical symmetries to minimize memory traffic
  - Mixed-precision methods
  - Autotuning for high performance on all CUDA-capable architectures
  - Eigenvector and deflated solvers (Lanczos, EigCG, GMRES-DR)
  - Multigrid solvers for optimal convergence Multi-source solvers
  - Domain-decomposed (Schwarz) preconditioners for strong scaling
  - Strong-scaling improvements
- A research tool for how to reach the exascale

## **QUDA - LATTICE QCD ON GPUS**

#### http://lattice.github.com/quda, BSD license

lattice /	quda	רפעuests ז פעניג אין	• Actions	Projects 4	🗐 Wiki	● Watch ▼	49
<> Code	(!) Issues 136						
Releases	Tags						
La	atest release		0				
	♡ v1.0.0 •• 66729fd	mathiaswagner rele	ased this on J	an 10			
	Verified	/ersion 1.0.0 -	- 10 Janı	uary 2020			
	Compare 🔻	<ul> <li>Add support for C using either GCC</li> <li>6.x or clang &gt;= 6.</li> </ul>	CUDA 10.2: Q or clang com .x are highly r	UDA 1.0.0 is sup pilers. CUDA 10. ecommended.	ported on x and eithe	CUDA 7.5-10. er GCC >=	2

 Significant improvements to the CMake build system and removal of the legacy configure build.



## **QUDA CONTRIBUTORS**

10+ years - lots of contributors

Ron Babich (NVIDIA) Simone Bacchio (Cyprus) Michael Baldhauf (Regensburg) Kip Barros (LANL) Rich Brower (Boston University) Nuno Cardoso (NCSA) Kate Clark (NVIDIA) Michael Cheng (Boston University) Carleton DeTar (Utah University) Justin Foley (Utah -> NIH) Joel Giedt (Rensselaer Polytechnic Institute) Arjun Gambhir (William and Mary) Steve Gottlieb (Indiana University) Kyriakos Hadjiyiannakou (Cyprus)

Dean Howarth (LLNL) Bálint Joó (Jlab) Hyung-Jin Kim (BNL -> Samsung) Bartek Kostrzewa (Bonn) Claudio Rebbi (Boston University) Hauke Sandmeyer (Bielefeld) Guochun Shi (NCSA -> Google) Mario Schröck (INFN) Alexei Strelchenko (FNAL) Jigun Tu (Columbia -> NVIDIA) Alejandro Vaguero (Utah University) Mathias Wagner (NVIDIA) Evan Weinberg (NVIDIA) Frank Winter (Jlab)

## AUTOTUNING





PASCAL





MAXWELL





FERMI

## **RECOMPILE AND RUN**

#### Autotuning provides performance portability



#### QUDA'S AUTOTUNER ensuring optimal kernel performance

virtual C++ class "Tunable" that is derived for each kernel you want to autotune

By default Tunable classes will autotune 1-d block size, shared memory size, grid size Derived specializations do 2-d and 3-d block size tuning

Tuned parameters are stored in a std::map and dumped to disk for later reuse

Built in performance metrics and profiling

User just needs to

State resource requirements: shared memory per thread and/or per block, total number of threads Specify a tuneKey which gives each kernel a unique entry and break any degeneracy

## **GENERATIONAL COMPARISON**

F<sub>µv</sub> kernel - batched 3x3 multiplication



10 📀 nvidia

# **MIXED PRECISION**

#### LINEAR SOLVERS

QUDA supports a wide range of linear solvers CG, BiCGstab, GCR, Multi-shift solvers, etc.

Condition number inversely proportional to mass Light (realistic) masses are highly singular Naive Krylov solvers suffer from critical slowing down at decreasing mass

Entire solver algorithm must run on GPUs Time-critical kernel is the stencil application Also require BLAS level-1 type operations while  $(|\mathbf{r}_k| \geq \varepsilon)$  {  $\beta_k = (\mathbf{r}_k, \mathbf{r}_k)/(\mathbf{r}_{k-1}, \mathbf{r}_{k-1})$   $\mathbf{p}_{k+1} = \mathbf{r}_k - \beta_k \mathbf{p}_k$   $\mathbf{q}_{k+1} = \mathbf{A} \mathbf{p}_{k+1}$   $\alpha = (\mathbf{r}_k, \mathbf{r}_k)/(\mathbf{p}_{k+1}, \mathbf{q}_{k+1})$   $\mathbf{r}_{k+1} = \mathbf{r}_k - \alpha \mathbf{q}_{k+1}$   $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha \mathbf{p}_{k+1}$  k = k+1} conjugate gradient



#### double-half

- Maintain solution vectors in high precision
  - Including the partial accumulator
- When true residual is injected, re-project the direction vector
- Use Polak-Ribière formula

$$eta_k := rac{\mathbf{z}_{k+1}^{\mathsf{T}}\left(\mathbf{r}_{k+1} - \mathbf{r}_k
ight)}{\mathbf{z}_k^{\mathsf{T}}\mathbf{r}_k}$$

double-half alt

 Residual replacement strategy of van der Worst and Ye



#### double-half

- Maintain solution vectors in high precision
  - Including the partial accumulator
- When true residual is injected, re-project the direction vector
- Use Polak-Ribière formula

$$eta_k := rac{\mathbf{z}_{k+1}^\mathsf{T}\left(\mathbf{r}_{k+1} - \mathbf{r}_k
ight)}{\mathbf{z}_k^\mathsf{T}\mathbf{r}_k}$$

#### double-half alt

 Residual replacement strategy of van der Worst and Ye



#### double-half

- Maintain solution vectors in high precision
  - Including the partial accumulator
- When true residual is injected, re-project the direction vector
- Use Polak-Ribière formula

$$eta_k := rac{\mathbf{z}_{k+1}^{\mathsf{T}} \left(\mathbf{r}_{k+1} - \mathbf{r}_k
ight)}{\mathbf{z}_k^{\mathsf{T}} \mathbf{r}_k}$$

#### double-half alt

• Residual replacement strategy of van der Worst and Ye

mixed precision:

apply Dslash in sloppy precision (single, half)

reliable residual replacement in high precision

ensures double-precision accuracy of final result

virtually identical iteration count



#### EIGENSOLVERS

Multiple workflows require repeated solution with different RHS with the same matrix

Multigrid not amenable to all linear operators

Eigenvector deflation is a robust alternative applicable to all operators Deflate out low modes from linear operator to accelerate the solver Cost of eigensolver is amortized if we solve enough RHS Aside: also use deflation to accelerate multigrid

Memory overheads can be limiting factor

May require storage of 1000s of vectors, ideally in fast memory

Configuration provided by HotQCD collaboration (Mukherjee *et al*)

### **MIXED-PRECISION DEFLATION**

V=48<sup>3</sup>x12, HISQ operator, physical light quarks, tol 10<sup>-10</sup>, 2xV100

📕 CG double 📕 CG double-single 📕 CG double-half 📕 CG double-quarter



### **DEFLATION STABILIZES LOW PRECISION**

V=48<sup>3</sup>x12, HISQ operator, physical light quarks, tol 10<sup>-10</sup>

16000

double-single-single



### **DEFLATION STABILIZES LOW PRECISION**

V=48<sup>3</sup>x12, HISQ operator, physical light quarks, tol 10<sup>-10</sup>

16000

double-single-single

double-half-single







Configuration provided by HotQCD collaboration (Mukherjee et al)

### **MIXED-PRECISION DEFLATION**

#### V=48<sup>3</sup>x12, HISQ operator, physical light quarks, tol 10<sup>-10</sup>, 2xV100



# SCALING

#### **STRONG SCALING**

Multiple meanings

Same problem size, more nodes, more GPUs

Same problem, next generation GPUs

Multigrid - strong scaling within the same run (not discussed here)

To tame strong scaling we have to understand the limiters Bandwidth limiters Latency limiters

#### **MULTI-GPU BUILDING BLOCKS**



Halo packing Kernel

Interior Kernel

Halo communication

Halo update Kernel

#### **BENCHMARKING TESTBED**

#### **NVIDIA Prometheus Cluster**

DGX-1 nodes

8x V100 GPUs connected through NVLink4x EDR for inter-node communicationOptimal placement of GPUs and NIC

Balanced GPU / IB configuration

understand limiters on a single node first



### **MULTI-GPU PROFILE**

#### overlapping comms and compute



#### **REDUCING LOCAL PROBLEM SIZE**

single GPU performance



#### **REDUCING LOCAL PROBLEM SIZE**



#### **STRONG SCALING PROFILE**

overlapping comms and compute



DGX-1,1x2x2x2 partitionin<mark>g</mark>

### **STRONG SCALING PROFILE**

#### Latencies ate my scaling



### **REDUCING API OVERHEADS**

Packing kernel writes to remote GPU using CUDA IPC



#### NVSHMEM

Implementation of OpenSHMEM1, a Partitioned Global Address Space (PGAS) library

**NVSHMEM** features

Symmetric memory allocations in device memory Communication API calls on CPU (standard and stream-ordered) Kernel-side communication (API and LD/ST) between GPUs NVLink and PCIe support (intranode) InfiniBand support (internode) Interoperability with MPI and OpenSHMEM libraries

#### DSLASH NVSHMEM IMPLEMENTATION First exploration

Keep general structure of packing, interior and exterior Dslash

Use nvshmem\_ptr for intra-node remote writes (fine-grained) Packing buffer is located on remote device Use nvshmem\_putmem\_nbi to write to remote GPU over IB (1 RDMA transfer)

Need to make sure writes are visible: nvshmem\_barrier\_all\_on\_stream
or barrier kernel that only waits for writes from neighbors

### **NVSHMEM DSLASH**

#### first exploration



## **NVSHMEM + FUSING KERNELS**

#### no extra packing and barrier kernels needed



## **ÜBER KERNEL**



DGX-1,1x2x2x2 partitioning

N

#### LATENCY REDUCTIONS



### **MULTI-NODE STRONG SCALING**

DGX SuperPOD (DGX2 nodes: 16xV100 (32GB), 8xEDR IB)



#### **MULTI-NODE STRONG SCALING**

DGX SuperPOD (DGX2 nodes: 16xV100 (32GB), 8xEDR IB)



# **APPLICATION SCALING**

### MILC NERSC BENCHMARK OVERVIEW

- MILC NERSC Benchmark comes in 4 lattice sizes
  - small 18<sup>3</sup>x36, medium **36<sup>3</sup>x72**, large 72<sup>3</sup>x144, x-large 144<sup>3</sup>x288
- Benchmark runs the RHMC algorithm
  - Dominated by the multi-shift CG sparse linear solver (stencil operator)
  - Also have auxiliary "Force" and "Link" computations
- Since 2012 MILC has built-in QUDA support
  - Enabled through a Makefile option
  - All time-critical functions off loaded to QUDA



#### MILC HMC SCALING ON DGX-2

NERSC MEDIUM BENCHMARK 36<sup>3</sup>x72



### **MILC SOLVER SCALING ON DGX-2**

NERSC MEDIUM BENCHMARK 36<sup>3</sup>x72



# MULTIPLICATIVE SPEEDUP

#### CHROMA HMC MULTIGRID

HMC typically dominated by solving the Dirac equation, but

Few solves per linear system

Can be bound by heavy solves (c.f. Hasenbusch mass preconditioning)

#### Multigrid setup must run at speed of light

Reuse and evolve multigrid setup where possible

Use the same null space for all masses (setup run on lightest mass)

Evolve null space vectors as the gauge field evolves (Lüscher 2007)

Update null space when the preconditioner degrades too much on lightest mass

#### CHROMA HMC-MG ON SUMMIT



From Titan running 2016 code to Summit running 2019 code we see >82x speedup in HMC throughput

Multiplicative speedup coming from machine and algorithm innovation Highly optimized multigrid for gauge field evolution gradient integrator

Chroma - Force

Bálint Joó

DPP-JI7

30ram Yoon -Frank Winter

### NODE PERFORMANCE OVER TIME

Multiplicative speedup through software and hardware



Time to solution is measured time to solution for solving the Wilson operator against a random source on a 24x24x24x64 lattice,  $\beta$ =5.5,  $M_{\pi}$ = 416 MeV. One node is defined to be 3 GPUs

# SUMMARIZING

#### QUDA ROADMAP On to QUDA 2.0

Multi-rhs block solvers for all stencils Contraction framework Improved strong scaling through NVSHMEM Beyond just regular QCD

Longer term: Investigate how well QUDA runs on C++17 pSTL

Post feature requests here: https://github.com/lattice/quda/issues



## **QUDA - LATTICE QCD ON GPUS**

Widely used for Lattice QCD applications on GPUs

State of the art solvers using mixed precision

Multigrid Deflation Block-Krylov solver

All components for gauge field evolution

Portable high-performance kernels through auto-tuning and careful optimization Tuned Multi-GPU scaling

GPU centric communication with NVSHMEM takes CPU limitations out Multiplicative speedup from hardware and software: more science

