# New Computational Trends in Lattice Gauge Theory

Peter Boyle (BNL)
pboyle@bnl.gov

- Lattice Gauge theory

- Electronic and hardware trends: what can we say about the future?

- Algorithmic response and trends

- Software response and trends

arXiv:2204.00039 "Lattice QCD and the Computational Frontier"
arXiv:2202.05838 "Applications of Machine Learning to Lattice Quantum Field Theory"
arXiv:1904.09725 "Status and Future Perspectives for Lattice Gauge Theory Calculations to the Exascale and Beyond"

# Lattice Gauge Theory

- **Numerical Euclidean path integral** in discretised 'femto-box' large $t \leftrightarrow$ hadronic ground state physics

- Wilson: preserves SU(N) gauge symmetry with $U_\mu(x) = e^{iagA_\mu(x)}$ connecting lattice sites

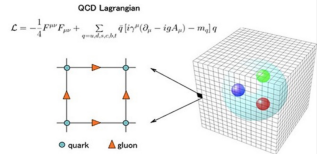- Partial derivatives replaced by finite difference

$$(\partial_\mu - igA_\mu(x))\psi(x) \to \frac{U_\mu(x)\psi(x+a\hat\mu) - U_\mu^\dagger(x-a)\psi(x-a\hat\mu)}{2a}.$$

- **Importance sample partition function:** $10^{10}$ degrees of freedom.
  Auxiliary momentum $(\pi)$ and pseudofermion $(\phi)$ integrals. Fermion determinant included stochastically via gaussian integral

$$Z = \int d\pi \int d\phi \int dU \quad e^{-\pi^2/2} e^{-S_G[U]} e^{-\phi^*(D^\dagger D)^{-1}\phi}.$$

- **Compute correlation functions**

$$\langle \mathscr{O} \rangle = \frac{1}{Z} \int_U e^{-S[U]} \mathscr{O}(U) dU \to \frac{1}{N} \sum_i \mathscr{O}(U_i)$$



QCD Lagrangian

$$\mathcal{L} = -\frac{1}{4}F^{\mu\nu}F_{\mu\nu} + \sum_{q=u,d,s,c,b,t} \bar{q}\left[i\gamma^\mu(\partial_\mu - igA_\mu) - m_q\right]q$$

2205.15373, 2204.00039, 2203.15810, 2203.10998
Scientific need for at least 10x current exascale computers in Lattice QCD
$128^3 \times 512 \times 16$ lattices for direct B simulation

  - g-2; hadronic vacuum polarisation and light by light

  - weak matrix elements and flavor physics

  - hadronic form factors, structure, PDFs

  - CP violation in B and Kaon sector

  - rare decay amplitudes

# Algorithms

**Phase-1: serially dependent gauge sampling**

- Metropolis-Hastings algorithms: Markov Chain Monte Carlo sampling arbitrary probability distributions

- Other key developers:
  M. Rosenbluth, A. Rosenbluth, E. Teller and A. Teller.
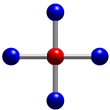
- Hybrid Monte Carlo (HMC)

**Phase-2: trivially parallelisable correlation function measurement**

- Dirac solvers:

  - Krylov space methods $\leftrightarrow$ polynomial approximation

  $$\mathscr{P}(\nabla^2) = c_n p^{2n} \sim \frac{1}{p^2}$$

  In a background gauge field eigenmodes are not plane waves: approximate Green's function as polynomial of discrete Dirac operator.
  - newer covariant multigrid methods

Stencil for a nearest neighbour finite difference in 2D



**Arianna Wright Rosenbluth** was an American physicist who contributed to the development of the Metropolis–Hastings algorithm. She wrote the first full implementation of the Markov chain Monte Carlo method.

### A History of the Metropolis–Hastings Algorithm

David B. HITCHCOCK

Historical DOE/Manhattan project connection

# Electronic and hardware trends: what can we say now about the future?

| Location | System | Interconnect (GB/s) per node (X+R) | Floating point performance (GF/s) per node | Memory Bandwidth (GB/s) per node | Year | System peak (PF/s) | FP / Interconnect | FP / Memory | Memory / Interconnect |
|---|---|---|---|---|---|---|---|---|---|
| LLNL | BlueGene/L | 2.1 | 5.6 | 5.5 | 2004 | 0.58 | 2.7 | 1.0 | 2.6 |
| ANL | BlueGene/P | 5.1 | 13.6 | 13.6 | 2008 | 0.56 | 2.7 | 1.0 | 2.7 |
| ANL | BlueGene/Q | 40 | 205 | 42.6 | 2012 | 20 | 5.1 | 4.8 | 1.1 |
| ORNL | Titan | 9.6 | 1445 | 250 | 2012 | 27 | 150.5 | 5.8 | 26.0 |
| NERSC | Edison | 32 | 460 | 100 | 2013 | 2 | 14.4 | 4.6 | 3.1 |
| NERSC | Cori/KNL | 32 | 3050 | 450 | 2016 | 28 | 95.3 | 6.8 | 14.1 |
| ORNL | Summit | 50 | 42000 | 5400 | 2018 | 194 | 840.0 | 7.8 | 108.0 |
| RIKEN | Fugaku | 70 | 3072 | 1024 | 2021 | 488 | 43.9 | 3.0 | 14.6 |
| NERSC | Perlmutter/GPU | 200 | 38800 | 6220 | 2022 | 58 | 194.0 | 6.2 | 31.1 |
| ORNL | Frontier | 200 | 181200 | 12800 | 2022 | >1630 | 906.0 | 14.2 | 64.0 |

- All DOE Exascale computing is GPU accelerated
- Huge gains in floating point
  not matched by gains in memory (14x) and interconnect (300x)
- Machines increasingly better suited for dense matrices and machine learning
- Lots of diversity and difficulty: ECP and SciDAC essential
  - Systems with AMD, Intel, Nvidia GPUs
  - Systems with CPU cores
  - HIP, SYCL, CUDA and conventional programming
  - Host memory, GPU memory
- Why so complicated??
- What can we actually say about short and mid-term future??



Aurora
The Argonne Leadership Computing Facility's exascale system will be used to dramatically advance scientific discovery and innovation.

Forthcoming systems will increase floating point performance dramatically, but not interconnect.
- Lattice gauge theory algorithms for gauge field sampling *must* change to exploit.
- Lattice gauge theory correlation function calculations can run brilliantly

# Why? You canny change the laws of physics

- Wire delay: time to charge a $L \times \pi r^2$ rod of metal depends on resistance & capacitance (RC time constant)[†]

$$RC \sim 2\rho\varepsilon \frac{L^2}{r^2} / \log(r_0/r) \sim \text{const} \times \frac{L^2}{r^2},$$

| Aspect ratio dictates wire delay - does not change when rod is shrunk |
| --- |

- Long thin wires are slow, short broad wires are fast

- 2.5D = chips edge to edge; 3D = vertical stacking : lots of short broad wires !

- Partitioned memory of modern HPC/GPU servers is dictated by electronics

  - On package memory, and advanced 3D packaging are here to stay.
  - Pools of fast small memory + pools of large slow memory.
    Data must be physically close to computational elements.
  - Future may bring additional layers (e.g. non-volatile)
    Soon: fast memory, slow host memory, SLOW host memory



TSV stacked HMC Micron

3D chipstacking enables many more wires and less energy per bit
Logical limit is silicon cubes even for computational logic to reduce wire delay
Cooling is a practical issue; TSV's conduct.

†: $\rho \equiv$ metal resistivity, $\varepsilon \equiv$ dielectric permittivity. $r_0 \ll L$ length scale small enough that potential well approximated by long rod approximation

# Use transistors more effectively post-Moore?

Why GPU's?

- Simpler arithmetic units tailored to array processing increase floating point density
- Constraining: consecutive loop iterations should perform same operation on consecutive elements of data
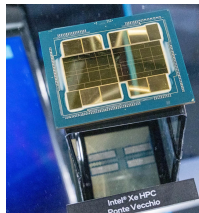- Less constrained than increasing SIMD vector lengths in CPUs
  - ORNL/Summit - Nvidia 6x V100 GPU + IBM Power CPU
  - NERSC/Perlmutter - Nvidia 4x A100 GPU + AMD x86 CPU
  - ORNL/Frontier - AMD 4x MI250-X GPU + AMD x86 CPU
  - ANL/Aurora (2023) - Intel 6x Ponte Vecchio GPU + Xeon
- As Moore's law approaches the atom barrier: better use of transistors could increase efficiency.
  - We have up to 100 billion of transistors in modern chips
  - 256MB of SRAM cache consumes 14Bn transistors but does no computation
- Near term: continued innovation
  - Nvidia Ampere → Hopper,
  - Intel Ponte Vecchio → Rialto Bridge (ISC 2022)
  - AMD → MI-300
- Nvidia, AMD, Intel competition is intense.
  Drives significant gains for Lattice gauge theory
- MPI messaging between nodes is a growing bottleneck that requires algorithm response



Intel Ponte Vecchio
multi-die package



Figure 4. Grace Hopper Superchip

Nvidia Grace (ARM CPU)
+ H100 integrated device

# Memory innovation for CPUs

GPU's use aggressive memory systems. e.g. 2.5D HBM memory or GDDR

- Seeing memory system innovation spreading beyond GPUs
- New multi-die packaging options: provide many short wires (Fujitsu/ARM, AMD x86, Intel x86, Nvidia Grace)
- Multicore CPU's are acquiring novel memories - what is the right scheme?
- Integration of CPU and GPU
  - Nvidia Grace/Hopper
  - AMD APU's: LLNL/El Capitan AMD MI300 combines GPU and CPU
- Dataflow / reprogrammable hardware (FPGA/spatial acceleration) may further increase floating point density
- Machine learning specific acceleration



Fujitsu/ARM Fugaku A64FX



AMD vcache 3D chipstack



Intel Sapphire Rapids + HBMe (Vision 2022)

# Programming models

**Programming interfaces:** syntactical details can be wrapped and hidden

- No fundamental reasons for "native" HIP / SYCL / CUDA 'mayhem'
- Expect standardisation of interfaces
  - OpenMP target for acceleration
  - Parallel STL /C++17

**Memory models:** fundamental semantic differences cannot be hidden

- HIP, SYCL and CUDA introduce virtual memory page based migration
- Even if physically partitioned, single virtual address space simplifies
- Performance penalty if pages are only 4KB $\Rightarrow$ need bigger pages.
- DOE/ASCR is in a position to dictate more flexibility and good performance for virtual memory
- Expect will eventually include: non-volatile, dram, in-package and on-chip cache.

# Software impact

- Dealing with all this brings productivity challenges to scientists

- Field is heavily reliant on QUDA and Grid for performance on GPUs.

- Abstract the APIs in a thin layer, capture loop bodies in "device lambda functions" through macros
  - Prudent to expect to use "native" compiler tools for each system: e.g. HIP, SYCL, CUDA
  - ECP project effort to port QUDA to HIP and SYCL almost complete
  - Prudent to expect to add OpenMP and Parallel C++ in near future. These may simplify and replace.

- Explicit data motion vs. unified memory choice vs. software caching

- Continuous effort and support to target new programming models and optimise on emerging architectures is required.

- ECP & SciDAC are enormously important. Humanpower AND early access to prepare for forthcoming hardware.

# Algorithmic impact

Finer simulations must cover a growing range of length scales. Current algorithms display criticla slowing down. USQCD is actively researching the following multiscale algorithm directions:

- Domain wall and staggered fermion multigrid

  - Multigrid has been transformational for the Wilson fermion discretisation
  - Accelerates observable calculations with minimal memory footprint
  - Similar acceleration is required for domain wall and staggered fermions

- Gauge invariant fourier acceleration to address critical slowing down as continuum limit is taken

  - Riemannian Manifold HMC
  - field transformation HMC

- Domain decomposition rational Hybrid Monte Carlo (next slide)

# Why domain decomposition

https://arxiv.org/pdf/2203.17119.pdf

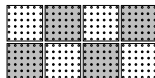Divide space into red and black hypercubes: Dirac operator may be written as

$$D = \begin{pmatrix} D_{\Omega} & D_{\partial} \\ D_{\bar{\partial}} & D_{\bar{\Omega}} \end{pmatrix}.$$

This can be LDU factorised: the factorised determinant is:

$$\det D = \det D_{\Omega} \det D_{\bar{\Omega}} \det \left\{ 1 - D_{\Omega}^{-1} D_{\partial} D_{\bar{\Omega}}^{-1} D_{\bar{\partial}} \right\},$$

The first two factors are *local* and last is non-local.

- Updating these on different timesteps in HMC ⇒ cost can be predominantly local.

- Avoids communication limits on future computers

- Size domains to computing nodes - tunes algorithm to islands of high performance created by multi-GPU computers.

- "fix-up" term can be bounded: if only evolve gauge links separated from boundary

- Generalised to odd flavours (2203.17119)

# Summary

- Scientific need for at least 10x current exaflop computers in lattice QCD
- Modern electronics industry is raising bandwidths by strapping chips together intelligently
- Optimal solution is unclear: recommended for elegant solutions to data placement and motion
  - e.g. improved virtual memory paging as the best method for hierarchical memory
- Current proliferation of APIs is unsustainable and standards should take over
- Using internal abstraction layers helps cope, needs continued software development support
- Massive opportunities in computing for lattice gauge theory, but also signficant challenges
- New algorithms are required to handle multiscale physics
  AND to match physics locality to supercomputer locality
- Continuous support (ECP, SciDAC) is essential to realise the potential
- CompF2 made recommendations on training, career paths, software support and lab-university joint positions.