

# Porting of Lattice QCD simulation software to GPUs

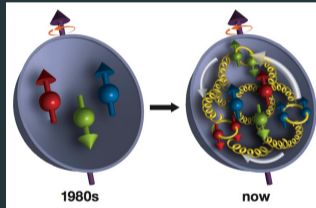
---

Nils Meyer, Stefan Solbrig, Tilo Wettig (University of Regensburg)

HPC Cafe, KONWIHR Workshop  
Feb. 14, 2023

# Overview of Lattice QCD

- QCD = Quantum Chromodynamics
  - Fundamental theory in the standard model of particle physics
  - Describes interaction of quarks and gluons and formation of bound states (e.g., proton)
  - Needed for interpretation of experiments at accelerators like CERN and in searches for “new physics”



Source: <https://www.energy.gov/science/np/articles/zooming-gluons-contribution-proton-spin>

- Lattice QCD = QCD on a 4-dimensional space-time lattice
  - Dominant problem: inversion of Dirac matrix, which is a very large sparse matrix whose dimension can be as large as  $O(10^9)$

# Main goals of the project

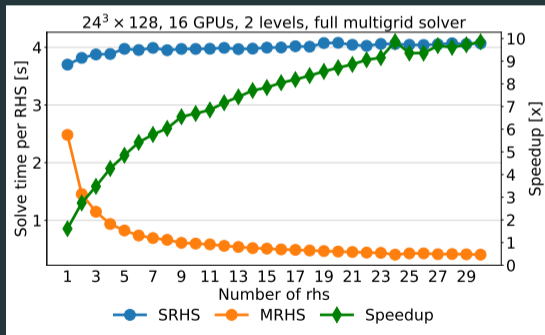
- “Grid” Lattice QCD framework already ported to GPUs (using CUDA, HIP and SYCL), but it does not meet all demands of U Regensburg particle physics group
- Here, focus on Wilson clover fermions and multigrid methods in Grid
- Work plan
  1. Optimization of global reductions (in progress)
    - In strong-scaling limit, global reductions become performance relevant
    - Root causes can be different in GPU systems compared to CPU systems
  2. Porting and optimization of multigrid methods (done)
    - Port our Grid multigrid code to GPUs
    - Algorithmic improvements (MRHS = multiple right-hand sides)
  3. Elimination of performance bottlenecks due to vendor libraries (in progress)
    - Performance analysis of dslash on available machines
    - Root cause of performance issues can be in vendor libraries (e.g., UVM issues in HIP)
    - Work with vendors to eliminate them or develop workarounds

# Porting and optimization of multigrid methods

- We already implemented multigrid methods in Grid, but only single right-hand side and not optimized for GPU
- Major overhaul and extension of existing codebase thanks to KONWIHR funding
  - Now supports multiple right-hand sides
  - Developed on (and optimized for) JUWELS Booster (NVIDIA A100)
  - By-product for free: Other architectures also supported (incl. Intel/AMD CPU, Arm NEONv8 and 512-bit SVE, AMD GPU), but performance not necessarily optimal
- JUWELS Booster results presented at annual Lattice QCD conference in 2022, proceedings paper submitted to PoS (LATTICE 2022), preprint available [arXiv:2211.13719]

# Porting and optimization of multigrid methods

- Performance
  - Solve time per RHS of FGMRES preconditioned by two-level multigrid method for CLS configuration U101 (volume =  $24^3 \times 128$ ) on 16 GPUs on JUWELS Booster
  - Speedup  $\sim 10x$  solving 30 RHS simultaneously compared to 30 sequential solves (SRHS)



# Access to production grade AMD machine

- LUMI: Fastest (and greenest) supercomputer in Europe, located in Finland
  - TOP500 #3 as of Nov. 2022
  - Rpeak = 428.70 PFlop/s
- LUMI-G: GPU partition <https://docs.lumi-supercomputer.eu/hardware/compute/lumig/>
  - 2560 nodes, each with
    - 4x AMD Instinct MI250X GPUs
    - 4x 200 Gb/s HPE slingshot interconnect
- We now have access to LUMI-G (via PRACE; until end of March 2023)
  - Joint project with U Bielefeld
  - Project plan (in progress)
    - Benchmark Grid and Bielefeld Lattice QCD codes
    - Test software stack and tune application parameters
    - Compare performance of Grid with JUWELS Booster