# Lessons-learned developing performance portable QMCPACK

## IDEAS Seminar May 2023

Paul Kent (kentpr@ornl.gov)

Oak Ridge National Laboratory

# Outline

- Brief introduction to Quantum Monte Carlo & QMCPACK

- Performance portability goals

- Challenges of using GPUs

- Development approach

- Summary

Aim: illustrate for other developers & code owners what has been productive for us and our ongoing pain points.

**OAK RIDGE**
National Laboratory

# Acknowledgements

ECP QMCPACK team including

- Peter Doak (ORNL)

- William Godoy (ORNL)

- Ye Luo (ANL)

ECP SOLLVE project [OpenMP+LLVM]

OLCF, ALCF staff

AMD, Intel, NVIDIA, HPE engineers

**OAK RIDGE**
National Laboratory

# Quantum Monte Carlo

- The most accurate, general approach for solving Schrodinger's equation for "real" materials. [Foulkes RMP 2001]

- The few approximations in QMC can be tested, unlike standard methods. Nominally $N^3$. **Tradeoff: large computational cost.**

- Not exact, but very accurate today, can treat "strong" electron correlation, applicable to metals, insulators & molecules.

- For details and tutorials, see QMCPACK YouTube channel & https://github.com/QMCPACK/qmc_workshop_2021

$$\frac{\partial |\psi\rangle}{\partial \tau} = -\hat{H}|\psi\rangle$$

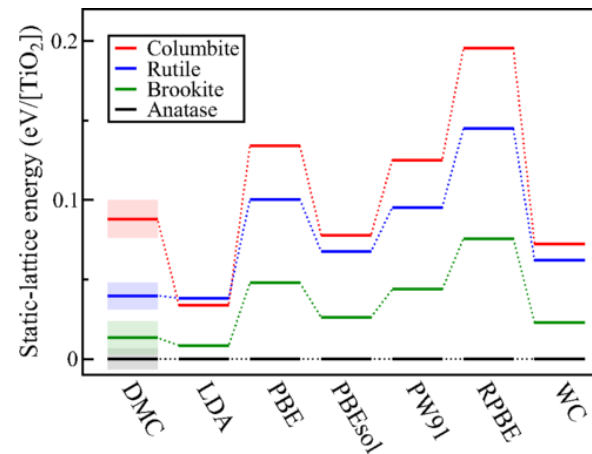$$|\psi(\delta\tau)\rangle = \sum_{i=0}^{\infty} c_i e^{-\epsilon_i \delta\tau}|\phi_i\rangle$$

Anatase
Didier Descouens CC-BY-SA 4.0

Rutile
Robert Lavinsky CC-BY-SA 3.0

Diffusion QMC for bulk $VO_2$

J. Trail et al. PRB **95** 121108 (2017)
Y. Luo et al. NJP **18** 113049 (2016)

# Performance Portability Goals

1. Run performantly on the full range of hardware, from laptops through to the #1 HPC machine and all 3 main vendor GPUs.

2. Use a single code path on all architectures, to the extent possible. Minimize maintenance burden, increase quality.

3. Retain ability to use specialized hardware & software, where merited.

**OAK RIDGE**
National Laboratory

# QMCPACK

- Open source, openly developed on GitHub, ~quarterly releases. Contributors credited on citation papers.

- C++17, HDF5, OpenMP+optional CUDA/HIP/SYCL+vendor dense linear algebra libraries. Highly vectorized, mixed precision supported.

- $O(2 \times 10^5)$ source lines.

- Science production using OpenMP target offload on NV GPUs with release versions of LLVM.

- New design has flexible dispatch, solves data movement and CPU fallback problem. Will always run unlike "legacy" GPU version.

- Code has undergone several major transitions: AoS to SoA CPU code for KNL, removal of "legacy" GPU version, ongoing removal of old CPU code paths.



**QMCPACK: an open source *ab initio* quantum Monte Carlo package for the electronic structure of atoms, molecules and solids**

OAK RIDGE
National Laboratory

# Recent QMC studies using QMCPACK

CrI$_3$ monolayers Staros JCP **156** 014707 (2022), H phase diagram Niu PRL **130** 076102 (2023), >10$^3$ molecules Huang JCTC **19** 1712 (2023). We aim to support and accelerate all of these calculations.

O(10$^3$) electrons          O(10$^2$) electrons          O(10$^{1-2}$) electrons

OAK RIDGE National Laboratory

# Challenge of exploiting GPUs

- Exascale-generation GPUs from NVIDIA, AMD, and Intel have >$10^4$ compute elements. Need >$10^6$ similar operations in flight for optimum performance.

- If we only have $10^{2-4}$ electrons, there naively will not be enough work.

- Generally, no single hot kernel. Kernels are both compute & memory bound.

- Few proven designs. QMC less mature than, e.g., quantum chemistry and classical molecular dynamics where multiple performant implementations are available.

NVIDIA A100 GPUs & AMD Milan CPUs

AMD GPUs and AMD CPUs

Intel Xe GPUs and Xeon CPUs

# Parallel Scalability

Despite needing communications every timestep, scalability is high due to high computational cost/step, careful MPI implementation.

See Kim et al. JPCM (2018) 10.1088/1361-648X/aab9c3



Strong scaling (fixed total samples) on Summit

Weak scaling (fixed samples per node) on Summit

OAK RIDGE
National Laboratory

# Key operations

Real space QMC uses both particle-based and dense linear algebra operations. Particle counts + matrix sizes can be small ($10^2$-$10^4$), requiring different choices to standard classical molecular dynamics or quantum chemistry techniques.

**Particle operations**
movement, interparticle distances, functions of position, minimum image when periodic

**Dense vector, matrix operations**
Spline and Gaussian basis set evaluation, determinant update, wavefunction optimization, BLAS1-3



Example Qs: Is it worth maintaining neighbor lists? Benefits, tradeoffs from sparsity?

**OAK RIDGE**
National Laboratory

# Miniqmc miniapp for design & performance experiments

- https://github.com/QMCPACK/miniqmc

- Order of magnitude smaller than QMCPACK

- Resulted in new design of QMCPACK with revised internal APIs and flexible runtime dispatch.

- Picked OpenMP target offload as default implementation route, supplemented if needed by vendor specific optimized code.

- Miniapp requires ongoing effort to maintain and update to keep synced with main application. => Unfortunately, miniqmc is currently out of date...



Profile validated vs main app ("proxy not imposter")



Easier profiling of miniapp aids in testing multithreaded offload strategy

OAK RIDGE
National Laboratory

# Algorithmic Challenge – How to map QMC to GPUs?

**Canonical QMC Algorithm**

do time step i   [ 1K-100K ]

<span style="color:red">do walker j  [ M walkers, ~1 per core, OpenMP ]</span>

<span style="color:green">do electron k [ N=$10^2$-$10^4$ ]</span>

propose new position $\underline{\mathbf{r}}_k$' [ O(~1-$N^2$) cost ]

evaluate $\boldsymbol{\Psi}(\underline{\mathbf{r}}_k$')  [ O(~N-$N^2$) cost ]

accept/reject using ~$|\boldsymbol{\Psi}'|^2/|\boldsymbol{\Psi}|^2$

if (accept) **update $\boldsymbol{\Psi}$** [ O($N^2$) cost ]

<span style="color:green">end k</span>

**evaluate Hamiltonian, Observables**

<span style="color:red">end j</span>

spawn/kill walkers, load balance

end i

- Works well on CPUs. Usually, 1 CPU thread per Monte Carlo walker.

- For GPUs, simply offloading the compute for each walker is not efficient since there is not enough numerical work.

**OAK RIDGE**
National Laboratory

# Previous GPU approach:
# Batching many independent walker moves

**Batched Metropolis QMC Algorithm, M walkers/node**

do time step i   [ 1K-100K ]

        do electron k [ $N=10^2-10^4$ ]
                propose M new positions $\{\underline{r}_k{'}\}_M$ [ $O(\sim M\text{-}MN^2)$ cost ]
                evaluate $\{\Psi(\underline{r}_k{'})\}_M$  [ $O(\sim MN\text{-}MN^2)$ cost ]
                accept/reject using $\sim|\Psi{'}|^2/|\Psi|^2$
                if (accept) **update** $\{\Psi\}_M$  [ $O(MN^2)$ cost ]
        end k
        **evaluate Hamiltonian, Observables for all M walkers**

        spawn/kill walkers, load balance
end i

- Batch (group) all operations over M walkers (Markov chains), now operations are $O(MxN)$ or $O(MxN^2)$. Choose M large enough to saturate GPU, typically 10-1000.

- CUDA version runs very efficiently, Esler et al. CISE **14** 40 (2012)

OAK RIDGE
National Laboratory

# New approach: multithreaded offload
# Batching smaller "crowds" of walkers

Time

**Single large crowd**

| EvalOrbitals | | EvalDeterminants |

GPU empty while CPU work performed

**Three smaller crowds, same total work**

EvalOrbitals — EvalDeterminants

EvalOrbitals — EvalDeterminants

EvalOrbitals — EvalDeterminants

Kernel running on GPU

Illustrative only. Not to scale

Launched but can't run. GPU busy.

Use multiple smaller batches ("crowds") launched from different host threads, not a single large batch.
- Trades some kernel efficiency for more asynchronous work and potentially greater throughput. Highly dependent on problem, hw+sw stack.
- Can recover original GPU algorithm with 1 crowd/thread.
- Highly beneficial if there is any significant CPU work remaining

**OAK RIDGE**
National Laboratory

# Results: 128 atoms NiO / 1536 electrons

# Results: 128 atoms NiO / 1536 electrons

# Results: 128 atoms NiO / 1536 electrons



>2x Increase in throughput

OAK RIDGE
National Laboratory

# Revisiting key algorithms

- We have replaced older algorithms with more compute bound/less memory bound algorithms.

- **New delayed update algorithm counterintuitively increases the operation count for higher performance.**

- Matrix multiply rich but extra work per step. Uses Sherman-Morrison-Woodbury formula to obtain wavefunction ratios during a delay period $n$, then update inverse. Avoid recalculating intermediates. Improves on our earlier algorithm (McDaniel 2017).

- ~2x faster on GPUs, ~10x faster on CPUs.

Bulk NiO supercells. Key: Atoms(Electrons per spin)

OAK RIDGE
National Laboratory

# Development approach

OAK RIDGE
National Laboratory

# General approach

- Focus on making the best use of our time

- Pragmatically adopt "best practices". Refine based on actual data from code review, CI, tools experience.

- Keep barrier for new developers, open-source contributions low.

- Limit required dependencies. Define a support policy for compilers, libraries etc.


E.g. Transitioned documentation to use sphinx & readthedocs. Minimal barrier to doc edits, plus full CI on changes.

**OAK RIDGE**
National Laboratory

# Ecosystem challenges: version control of dependencies



Minimal QMCPACK dependencies excluding compilers, many optional python dependencies
(spack graph output)

- Even "low dependency" apps have many dependencies. These are all undergoing development… changing GPU support, python module changes, HDF5 API updates etc. can all lead to breakage.

- For nightly testing, "version control" achived via spack package manager (https://spack.io/) to cover a sparse matrix of older/newer software. However, users use ~any combination of versions…

OAK RIDGE
National Laboratory

# Testing and Continuous Integration

- Large set of unit, deterministic, & stochastic integration tests built up. Test subset run in CI, plus more extensive sets nightly and weekly for many different compiler, CPU, GPU, library combinations.

- Helps us make large changes to the code, onboard new developers, engage with vendors, contributors.

- Pull requests undergo review, testing, coverage reporting, sanitizer tests. Procedures for reviews developed, e.g., merger can not be at same institution as PR.

- CI uses GitHub actions, plus our own hardware, needed to test GPUs not available in cloud. Aim for O(1h) turnaround. Window for input to vendors for fixes in their next release is small.

- OpenMP target offload "GPU" code partly tested via offload to host CPUs with LLVM. Huge time, $ savings.



CI reporting on PRs



Initial failure rate of PRs



QMCPACK cdash

OAK RIDGE
National Laboratory

# Many edge case bugs found & resolved

For QMCPACK problems, testing regimen minimizes chance of recurrence. Statistics let us (re)focus testing effort on most critical areas. For external problems, comprehensive testing **allows us to give prompt feedback to the relevant developers.**

Examples

- QMCPACK:
  - State machines associated with efficient Monte Carlo.
  - Handling of optional features via legacy #ifdefs requiring separate builds.

- Wider software stack:
  - Compilers, particularly OpenMP offload. Problems with complex reductions, multithreading... Latest releases of LLVM in production on NV.
  - Libraries. E.g. Numerically incorrect results from GPU dense linear algebra libraries, threading problem in CPU OpenBLAS (fixed promptly).
  - Many transient packaging and compatibility issues associated with specific versions of libraries, compilers, tools. Important to have latest versions promptly in spack, in addition to older versions.

**OAK RIDGE**
National Laboratory

# Status

As of LLVM 15.0 (released 9/2022), performance portable version sufficiently close in performance to limited feature legacy GPU implementation for science production, but with full capabilities available.

# Ongoing challenges, open questions

- Further maturation of the ecosystem is necessary. Having Frontier, Aurora, Polaris (etc.) in production will help.

- Helpful to have "stable" and "leading edge" machines available for CI. It is not practical for every app team to run their own CI. Lack of access slows development velocity.

- Automated testing at facilities would help catch issues with vendor provided software, MPI, their unique environment etc.

- Can we obtain ~full performance with only OpenMP using newer/revised versions of the standard, or will some CUDA/HIP/SYCL still be required? What is needed in future C++ standards?

**OAK RIDGE**
National Laboratory

# Conclusions

- Performance portable QMC is a challenge!

- Performance portable QMCPACK is in science production.

- Next challenges: taming memory usage, science features.

- Modern development practices, particularly testing, has improved code quality, enabled large changes, and increased our efficiency overall. These practices require dedicated resourcing and staff.

Questions, comments? kentpr@ornl.gov

OAK RIDGE
National Laboratory