

Platform independent profiling of a QCD code

Marina Krstić Marinković*

CERN, TH Department, 1211 Geneve, Switzerland

School of Mathematics, Trinity College Dublin, Dublin 2, Ireland

E-mail: Marina.Marinkovic@cern.ch

Luka Stanisic

Inria Bordeaux Sud Ouest, Bordeaux, France

E-mail: Luka.Stanisic@inria.fr

The supercomputing platforms available for high performance computing based research evolve at a great rate. However, this rapid development of novel technologies requires constant adaptations and optimizations of the existing codes for each new machine architecture. In such context, minimizing time of efficiently porting the code on a new platform is of crucial importance. A possible solution for this common challenge is to use simulations of the application that can assist in detecting performance bottlenecks. Due to prohibitive costs of classical cycle-accurate simulators, coarse-grain simulations are more suitable for large parallel and distributed systems. We present a procedure of implementing the profiling for `openQCD` code [1] through simulation, which will enable the global reduction of the cost of profiling and optimizing this code commonly used in the lattice QCD community. Our approach is based on well-known `SimGrid` simulator [2], which allows for fast and accurate performance predictions of HPC codes. Additionally, accurate estimations of the program behavior on some future machines, not yet accessible to us, are anticipated.

34th annual International Symposium on Lattice Field Theory

24-30 July 2016

University of Southampton, UK

*Speaker.

1. Introduction

The numerical simulations in lattice QCD are large scale. Besides constant algorithmic development, they require massive parallelization and consume a substantial amount of resources of the leading supercomputing centers in Europe and worldwide. In order to efficiently exploit ever-changing HPC architectures, one continually needs to adapt and optimize the existing codes for new hardware. To develop a lattice QCD code that scales well on diverse platforms, besides the access to different machines, the developer needs to be able to perform large experimental campaigns in a fairly short period of time. It would therefore be advantageous to design a universal procedure that would allow for benchmarking of the native code executions on a wide range of supercomputing platforms. We propose a way to use an adaptable simulator for parallel systems, `SimGrid` [2] in combination with `openQCD` [1] code, in order to measure the scalability of the latter without accessing the target machine. The procedure presented here can in the same way be applied to other up-to-date QCD codes.

The aim of this setup is to speed up the benchmarking of the chosen QCD code on an unknown machine and to enable benchmarking of platforms to which the user does not yet have access. Moreover, many worlds' supercomputing centers, after awarding computer time, let the choice of the suitable machine within the given center to the users. Being able to estimate the optimal choice of the machine with the appropriate number of nodes in these situations in timely manner and without wasting computational resources is another possible application of this procedure.

In Section 2, we briefly describe the program package `openQCD` and perform a conventional native set of benchmarks for a pure gauge simulation code. `SimGrid` simulation framework is introduced in Section 3. We outline a procedure for simulating `openQCD` programs using `SimGrid` in Section 4. After reporting our preliminary results, we discuss next steps envisioned for this study in Sections 5 and 6, respectively. Finally, we conclude and address possible future applications of this work in Section 7.

2. `openQCD`: benchmarking and scaling

The simulation package `openQCD` contains programs for generating both pure gauge and dynamical QCD configurations, the latter with a choice of $O(a)$ improved Wilson fermion action. Although original motivation was to create a portable simulation tool for recently introduced open boundary conditions [3], the programs in this package allow for several different choices of boundary conditions¹: open, periodic, Schrödinger Functional (SF), and open-SF. The simulation program is based on the Hybrid Monte Carlo algorithm [6] and supports parallelization in 0,1,2,3 or 4 dimensions. All the programs in this package are highly optimized for machines with modern Intel or AMD processors, but will run correctly on any system that complies with the ISO C89 and the MPI 1.2 standards. Optimized version of the code for BlueGene/Q machines is available from [7]. The code is open source with GPL license [1].

¹Extending the code to include QED field dynamics and an additional choice of C* boundary conditions [4, 5] is in progress.

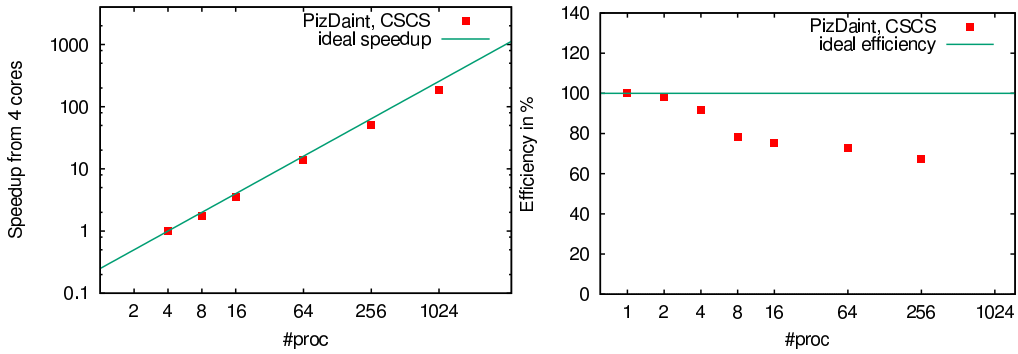


Figure 1: Strong and weak scaling of the `ym1` program on a typical Cray machine (PizDaint, CSCS). Strong scaling plot shown in the left panel is obtained by running `ym1` on a global lattice of a size 32^4 . Local lattice size used for the weak scaling plot shown in the right panel is 8^4 . Both plots are obtained with run parameters from the sample input file `ym1.in` distributed with `openQCD` [1].

We apply the first set of benchmarking on the program `ym1` that generates ensembles of $SU(3)$ gauge configurations. Exactly which theory is simulated depends on the parameters passed to the program. For the scaling tests shown in Figure 1, we chose tree-level improved Symanzik gauge action and open-SF boundary conditions. The scaling tests are performed at the PizDaint² machine of the CSCS supercomputing center.

3. SimGrid simulation framework

`SimGrid` is a versatile simulation framework for distributed systems that allows conducting studies in different fields, such as grid, cloud, peer-to-peer and volunteer computing. In the HPC context, more precisely regarding MPI applications, `SimGrid` provides `SMPI` module [8] which can be used in two different modes: offline and online. In the offline mode, a trace of a previous native execution on a supercomputer is replayed in the simulator. Quick and accurate predictions can be obtained, but such results are hard to faithfully extrapolate to different scenarios. The main reason behind this is that modern hardware and software stacks are so complex, that even a minor change to a platform parameter or an application parameter can cause substantial change to the overall executions. On the other hand, online mode, which relies only on the target platform description without any traces, executes directly the application. Therefore, simulations in this mode can mimic much better any modifications to the application or machine setup. Although such approach is much harder to implement and validate, its predictive capabilities are more advanced. There is a plethora of MPI simulators used in HPC community ([9], [10], and many more), but we have chosen `SimGrid` as its faithful online simulations with accurate modeling of network contention are needed for this type of study.

²As of December 2016, PizDaint machine has been upgraded. The benchmarks performed in this work correspond to the configuration of June 2016.

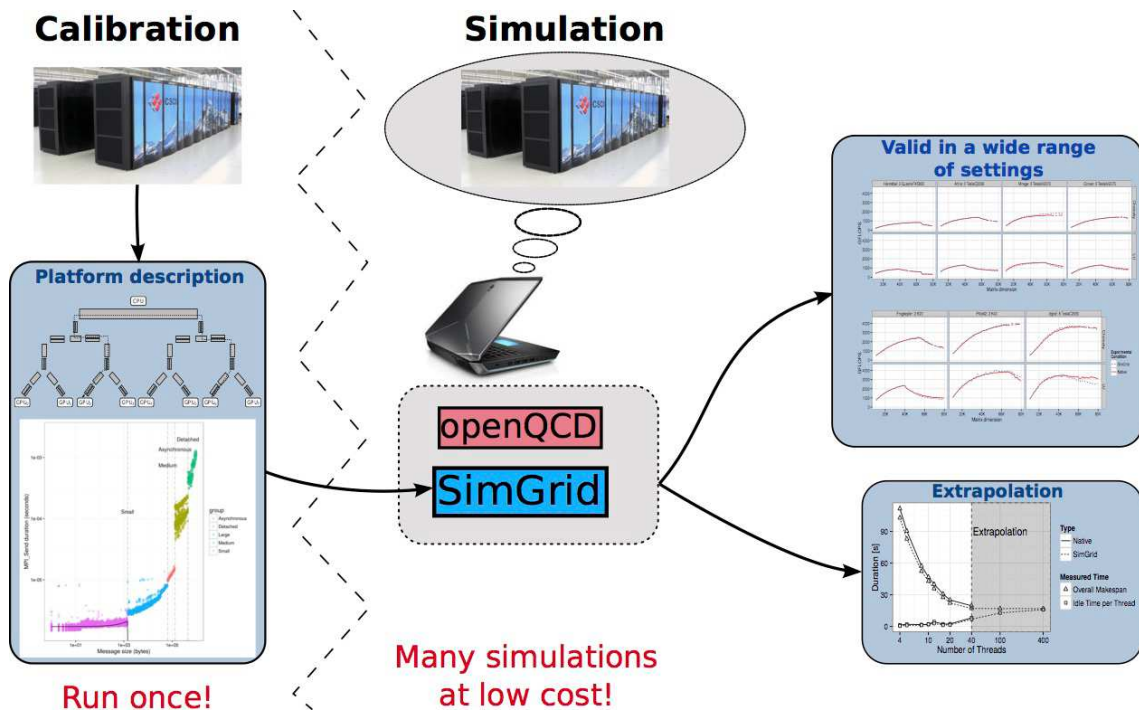


Figure 2: Workflow for coupling the `openQCD` package and `SimGrid`. Once platform characteristics have been benchmarked and network models have been instantiated, many simulations can be run at low cost to provide faithful predictions for a wide range of scenarios.

4. Methodology for profiling `openQCD` with a `SimGrid` simulator

Workflow for conducting simulations of `openQCD` is presented in Figure 2. In the initial phase, target machine needs to be carefully calibrated. One needs to reserve only few nodes of a possibly large supercomputer. These nodes are first benchmarked to capture their characteristics, typically number of processing units and their corresponding speed. Then, a different set of benchmarks containing a wide range of MPI calls is performed to measure the network characteristics and its behavior in the case of contention. The obtained values will be used by `SimGrid` when instantiating its communication models. These coefficients together with the previously captured platform details are saved in a separate XML file.

Once this platform description file is generated, no more access to the supercomputer is needed. One needs to simply compile `openQCD` application with the `SimGrid` SMPI, instead of classical MPI library. This new program can then be run on a single core of a commodity laptop and provide faithful performance predictions. `openQCD` application “thinks” that it is executed on a large scale machine using many nodes, while it is actually run inside a single `SimGrid` process.

Since almost no actual computations and no data transfers are performed, simulation is typically much faster than the native execution and additionally consumes much less memory. This allows for running many profilings at low cost and fully assessing the influence of an application input parameter (e.g. lattice size, bare gauge coupling, etc. in the case of a lattice QCD code). Since

platform description file can be easily manually edited, one may also explore “what-if” scenarios by increasing the number of nodes, number of processing unit, network latency and bandwidths, etc. This can provide invaluable information to the `openQCD` developers and users who have a challenging task of correctly dimensioning the problem size they want to study, as well as determine the optimal machine for such workload.

5. Initial evaluation of `ym1`

We used `SimGrid` to simulate pure Yang-Mills simulation program (`ym1`) of `openQCD`. We compared the time needed to run four trajectories in the generation of pure gauge configurations using Symanzik gauge action and open-SF boundary conditions, on a 8^4 local lattice. The first tests were run on 4 and 8 processors (2 Dodeca-core Haswell Intel Xeon 2,5 GHz 128GB RAM machine) on PlaFRIM computing cluster, corresponding to the global lattice sizes of $16^2 \times 8^2$ and $16^3 \times 8$. Actual execution times of the `ym1` program were compared to the execution times predicted by `SimGrid`. The outcome is shown in Figure 3. Native `openQCD` code demonstrates very good performance, as for these weak scaling experiments the timings between two executions are very similar. Although `SimGrid` provides an accurate estimation in the case of 4 processors, for 8 nodes prediction is largely overestimating the overall execution duration.

To understand the source of such discrepancy, we intend to introduce timestamps into `openQCD` code. Then, by comparing the traces between native and `SimGrid` executions, we expect to detect and later correct the error in our models, which will allow for achieving the good estimations for large scale executions. Currently, we suspect that the issues we encountered are caused by bogus network drivers and problems with MPI libraries on PlaFRIM cluster at the time of benchmarking, which lead to inaccurate generation of communication models. However, these assumptions need to be empirically verified.

6. Future prospects

For practical applications in lattice QCD, one would need to go beyond testing the program for pure gauge configurations generation `ym1`. On the example of `openQCD`, this would mean to get the good matching of the native and simulated execution times of the `qcd1` program. Providing faithful estimates for the dynamical gauge configuration generation program requires accurate predictions of the Dirac matrix inversions. However, since `SimGrid` has been already successfully used for the simulation of the solvers for sparse linear systems [11], we expect a smooth transition to more complex simulation programs such as `qcd1`.

Presented study was mostly exploratory, as we wanted to evaluate the feasibility of the simulation approach on `openQCD` application and identify the biggest challenges. Once our methodology is extensively validated on a set of small scale experiments and it becomes more robust, we plan to generalize our approach in order to use it for the performance predictions of the large scale

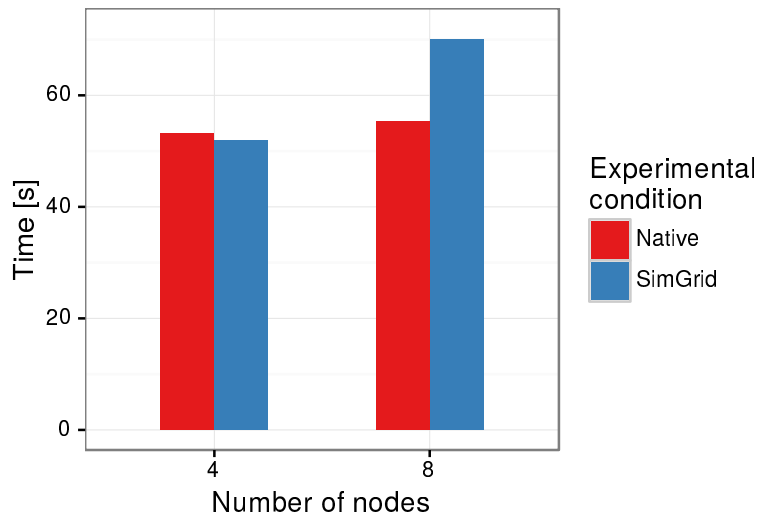


Figure 3: Native (red) and simulated (blue) times on 4 and 8 processors at 2 Dodeca-core Haswell Intel Xeon E5-2680 machine. We simulated four HMC trajectories of pure Yang-Mills simulation program that is a part of `openQCD`. Run parameters are taken from the sample input file `ym1.in` distributed with `openQCD` [1].

arbitrary machines. To accomplish this challenging task, it is mandatory to obtain accurate calibrations of other supercomputers. This is a common prerequisite shared by many researchers in `SimGrid` community, and thus there is an ongoing work by `SimGrid` developers in automatizing the calibration process and collecting platform description files for many supercomputers, such as: Stampede and Blue Waters in the USA, MareNostrum in Spain, CSCS in Switzerland, etc.

7. Summary and Outlook

We have proposed a procedure of implementing an intermediate profiling for a QCD code using `SimGrid` simulator, which is particularly convenient for simulating parallel applications based on MPI programming paradigm. We tested our methodology on the `ym1` program of `openQCD`, although, the same procedure can be easily applied to other lattice QCD codes. Our initial prototype works and provides encouraging results, but still needs some improvement in order to reach the desired accuracy, which would allow for doing experiments on a much larger scale. We expect these to come from the refinement in communication modeling that we are currently pursuing.

In the future, with more robust and automatized solution, we envisage application of the proposed approach on multiple use cases. First, this can help users who are preparing the computer time demands for a machine they yet do not have access to. Additionally, it may be useful for choosing the most suitable machine within already awarded computer time allocation, in a timely manner. Furthermore, this procedure is not only applicable on currently existing machines. It can even be applied on machines that are still in the design phase, once the features such as machine topology, processor speed and communication characteristics become known. Nevertheless, such results would have to be interpreted carefully, as unknown systems could contain some yet undiscovered (and thus not modeled) phenomena.

Finally, source code and raw data generated during this study are publicly available³ to anyone interested in learning more details about this work.

Acknowledgments

This work was supported by a grant from the Swiss National Supercomputing Centre (CSCS) under project ID s642. Some experiments presented in this paper were carried out using the PlaFRIM experimental testbed, being developed under the Inria PlaFRIM development action with support from LABRI and IMB and other entities: Conseil Regional d'Aquitaine, FeDER, Universite de Bordeaux and CNRS (see <https://plafrim.bordeaux.inria.fr/>).

References

- [1] <http://luscher.web.cern.ch/luscher/openQCD/>.
- [2] Henri Casanova, Arnaud Giersch, Arnaud Legrand, Martin Quinson, and Frédéric Suter. Versatile, scalable, and accurate simulation of distributed applications and platforms. *Journal of Parallel and Distributed Computing*, 74(10), June 2014.
- [3] Martin Luscher and Stefan Schaefer. Lattice QCD without topology barriers. *JHEP*, 07:036, 2011.
- [4] U. J. Wiese. C periodic and G periodic QCD at finite temperature. *Nucl. Phys.*, B375:45–66, 1992.
- [5] Biagio Lucini, Agostino Patella, Alberto Ramos, and Nazario Tantalo. Charged hadrons in local finite-volume QED+QCD with C* boundary conditions. *JHEP*, 2015.
- [6] S. Duane, A. D. Kennedy, B. J. Pendleton, and D. Roweth. Hybrid Monte Carlo. *Phys. Lett.*, B195:216–222, 1987.
- [7] BlueGene/Q optimized openQCD.
http://hpc.desy.de/simlab/codes/openqcd_bgopt/.
- [8] Paul Bedaride, Augustin Degomme, Stéphane Genaud, Arnaud Legrand, George Markomanolis, Martin Quinson, Mark Stillwell, Lee, Frédéric Suter, and Brice Videau. Toward Better Simulation of MPI Applications on Ethernet/TCP Networks. In *Intl. Workshop on Perf. Modeling, Benchmarking and Simul. of HPC Systems*, November 2013.
- [9] A. F. Rodrigues, K. S. Hemmert, B. W. Barrett, C. Kersey, R. Oldfield, M. Weston, R. Risen, J. Cook, P. Rosenfeld, E. CooperBalls, and B. Jacob. The Structural Simulation Toolkit. *SIGMETRICS Perform. Eval. Rev.*, 38(4):37–42, March 2011.
- [10] Gengbin Zheng, Gunavardhan Kakulapati, and Laxmikant Kalé. BigSim: A Parallel Simulator for Performance Prediction of Extremely Large Parallel Machines. In *Proc. of the 18th International Parallel and Distributed Processing Symposium (IPDPS)*, April 2004.
- [11] Luka Stanisic, Emmanuel Agullo, Alfredo Buttari, Abdou Guermouche, Arnaud Legrand, Florent Lopez, and Brice Videau. Fast and Accurate Simulation of Multithreaded Sparse Linear Algebra Solvers. In *The 21st IEEE International Conference on Parallel and Distributed Systems*, Melbourne, Australia, December 2015.

³<http://gitlab.inria.fr/stanisic/openqcd-simgrid/>