

UNIVERSITY NAME (IN BLOCK CAPITALS)

Thesis Title

by

Attila Döme Lehóczy

A thesis submitted in partial fulfillment for the
degree of Doctor of Philosophy

in the

Faculty Name

Department or School Name

2013. június 3.

Declaration of Authorship

I, AUTHOR NAME, declare that this thesis titled, THESIS TITLE' and the work presented in it are my own. I confirm that:

- This work was done wholly or mainly while in candidature for a research degree at this University.
- Where any part of this thesis has previously been submitted for a degree or any other qualification at this University or any other institution, this has been clearly stated.
- Where I have consulted the published work of others, this is always clearly attributed.
- Where I have quoted from the work of others, the source is always given. With the exception of such quotations, this thesis is entirely my own work.
- I have acknowledged all main sources of help.
- Where the thesis is based on work done by myself jointly with others, I have made clear exactly what was done by others and what I have contributed myself.

Signed:

Date:

„Write a funny quote here.“

If the quote is taken from someone, their name goes here

UNIVERSITY NAME (IN BLOCK CAPITALS)

Abstract

Faculty Name

Department or School Name

Doctor of Philosophy

by Attila Döme Lehóczy

The Thesis Abstract is written here (and usually kept to just this page). The page is kept centered vertically so can expand into the blank space above the title too...

Acknowledgements

The acknowledgements and the people to thank go here, don't forget to include your project advisor...

Tartalomjegyzék

| | |
|---|-------------|
| Declaration of Authorship | i |
| Abstract | iii |
| Acknowledgements | iv |
| List of Figures | vii |
| List of Tables | viii |
| Abbreviations | ix |
| Physical Constants | x |
| Symbols | xi |
| | |
| 1. Introduction | 1 |
| 2. Literature Review | 2 |
| 2.1. Introduction | 2 |
| 2.2. MPI | 2 |
| 2.2.1. OpenMPI | 3 |
| 2.2.2. MPICH | 3 |
| 2.3. Modelling and Simulation | 4 |
| 2.3.1. Advantages of Modeling | 4 |
| 2.3.2. Analytical and Simulation Models | 5 |
| 2.4. SimGrid | 5 |
| 2.5. SMPI | 6 |
| 2.6. STAR-MPI | 7 |
| 3. Problem Description | 9 |
| 3.1. Reproducible research | 9 |
| 3.2. Testing framework | 9 |
| 3.3. The testing process | 10 |
| 3.3.1. SG traces | 10 |

| | |
|---|-----------|
| 3.3.2. RL traces | 10 |
| 3.4. Instrumentation overhead | 11 |
| 3.5. Clock synchronization | 11 |
| 4. Implementation | 12 |
| 5. Evaluation Plan | 13 |
| 6. Results | 14 |
| 7. Conclusion | 15 |
| | |
| A. Appendix Title Here | 16 |
| | |
| Irodalomjegyzék | 17 |

Ábrák jegyzéke

Táblázatok jegyzéke

Abbreviations

LAH List Abbreviations **Here**

Physical Constants

Speed of Light $c = 2.997\,924\,58 \times 10^8 \text{ ms}^{-\text{s}}$ (exact)

Symbols

| | | |
|----------|-------------------|------------------------|
| a | distance | m |
| P | power | W (Js^{-1}) |
| ω | angular frequency | rads^{-1} |

For/Dedicated to/To my...

1. fejezet

Introduction

2. fejezet

Literature Review

simulation: system clocks dont differ tau and measurement: delays operations, communications smpi: no high-performance networking hardware -i just gigabit ethernet make sure to give openmpi/mpich parameters to run gigabit ethernet, not infiniband

2.1. Introduction

In the literature review, we discuss the area of research covered by the thesis, citing relevant papers and articles that serve as a base of ideas for this document.

First, we talk about the Message-Passing Interface (MPI), the parallel programming API used for the purposes of this thesis. We also talk about the different implementations of this API, notably OpenMPI and MPICH. We discuss the methods of modeling and simulation in general. Then, we talk about SimGrid, which is a simulation-based framework, and SMPI, the implementation of MPI that runs on top of SimGrid. After that, we also describe StarMPI, a set of MPI communication routines, presenting a technique that can improve the performance of MPI applications.

2.2. MPI

Distributed computing is a very active and important subject of research in computer science, including fields such as cluster computing, grid computing, Cloud computing, or peer-to-peer computing. Communication between the different processes in a distributed application can be implemented in a number of ways. As communication is necessary in most cases, a standardized communication protocol can be a lot of help when developing a distributed program. The Message-Passing Interface (MPI) is a language-independent

message-passing library interface specification. It is not a language, but a standard - there exist multiple MPI implementations. Since its take-off, it has become a de facto standard for inter-process communication. The standard provides vendors a clear set of routines, that they can implement efficiently, or in a way that it suits the hardware they provide.^[1]

2.2.1. OpenMPI

OpenMPI is an MPI implementation with the goal of being able to achieve good performance on a wide range of different aspects of high-performance computing. To efficiently support multiple types of parallel machines, high performance “drivers” for all established interconnects are developed. These include TCP/IP, shared memory, Myrinet, Quadrics, and Infiniband. Features for checking data integrity are provided in order to account for network transmission errors. With the utilization of message fragmentation and striping over multiple (potentially heterogeneous) network devices, OpenMPI provides an increased bandwidth to applications, as well as the ability to handle the failure of network devices during runtime.^[2]

2.2.2. MPICH

MPICH was originally developed during the MPI standards process starting in 1992 to provide feedback to the MPI Forum on implementation and usability issues. This original implementation was based on the Chameleon portability system to provide a light-weight implementation layer (hence the name MPICH from MPI over CHameleon). Around August 2001, development begun on a new implementation called MPICH2.^[3] This implementation introduced improvements on collective communication operations by using multiple algorithms, choosing between them depending on certain variables - for example the message size.^[4] Another important result during the development of MPICH2 was the design of the Nemesis communication subsystem and the porting of MPICH2 on that system. The efficient implementation of shared-memory communication helped Nemesis MPICH2 achieve low latency and high bandwidth.^[5] Starting with November 2012, the project is renamed to MPICH, with version number 3.0.^[3]

2.3. Modelling and Simulation

In distributed computing, modelling means creating an abstraction of a real system by taking only the aspects of it that are relevant to the system's behavior into account. Once constructed, such a model becomes a tool with which we can investigate the behavior of the system.[6]

2.3.1. Advantages of Modeling

Modelling and simulation techniques have been used extensively in parallel computing and is an ongoing research topic, with new challenges continuously arising. There are various reasons for its importance.

Conducting experiments on real-world systems can be infeasible because experimenting would disrupt the service that is provided by the system. For example, in the case of a mail server, experiments or monitoring could cause delay, or maybe even data loss. Service disruption can sometimes be even dangerous, in addition to being an inconvenience: in the case of a nuclear reactor, delay or loss of data can prove fatal. Timeliness can be as important in such systems as correctness. However, performance analysis and monitoring might be crucial to draw conclusions about maintenance, for example. Another problem with direct experimentation is that the information we are looking for may not be available, or may be complicated to get. For example, in most operating systems, it is difficult to obtain the exact timing of instruction-level events.[6] Also, when conducting experiments on a real-world system, results are often non-reproducible, due to resource dynamics.[7] Another argument on the side of modelling is that it provides the ability of experimenting on different configurations. Investing in a large-scale computer cluster, or the setup of a distributed grid environment is an expensive and tedious process. Investors want to make sure that they get what they want: they impose performance constraints on the system. This means that they want to know how the system will behave before buying it and setting it up. To predict the behavior, experiments are needed to be conducted. We need to do these experiments on different setups, before finding out which one is the best in the current situation. Changing the hardware or software configuration parameters on a real-world system is very inconvenient - in most cases, it's not doable, because of time and money constraints. Thus, the solution is to simulate the desired system, and run the experiments there. This way, changing the configuration is simple and costless.[6] Another great benefit of simulation is that in a classroom setting, students can learn the principles of high-performance and distributed computing without actual access to a parallel platform.[8]

2.3.2. Analytical and Simulation Models

The accuracy of a model can vary: we can make an analytical, or qualitative model, in which all definite values are abstracted away - in this case, we get a representation of the system, which can be analysed mathematically to deduce its behavior. When using this method, no experiments can be conducted, we solely rely on theoretical analysis. In contrast, a simulation model is a stochastic model, which is an algorithmic abstraction of the real-world system that can be executed to reproduce the system's behavior. This model is also called a quantitative model, as we can get estimates of the modeled system's quantitative attributes, such as response time or throughput. In other words, we can use a simulation model to conduct performance analysis on a system.[6][9]

When wanting to get a prediction about how a specific system would perform, a theoretical model, in most cases, produces unreliable and unrealistic results - it's not feasible for such accurate predictions. Most research results are obtained via empirical evaluation of experiments.[7] For these reasons, we use the simulation model in this thesis. As we stated before, such a model can be executed, which is called simulation. During simulation, the model behaves like the real system would. Such a model contains more aspects of the real system, compared to the theoretical model, in order to accurately represent the system, while still avoiding unnecessary detail.[6] Creating and executing a simulation model is complicated, computationally expensive and poses a large number of challenges, thus, a good simulation framework can prove of much help when conducting experiments.

2.4. SimGrid

For the aforementioned reasons, simulation techniques have been historically widely utilised in several areas of computer science, e.g. microprocessor design, network protocol design. Due to this, a lot of effort went into developing the technology and as a result, widely used and reliable simulation frameworks have been developed in these areas. However, there hasn't been a well-developed standard simulation tool for what we talk about in this thesis: execution of distributed applications on distributed computing platforms. Rather, there has only been a number of in-house developed, highly specialized tools to satisfy the need of the community. SimGrid is a more generic simulation framework that is being developed to be one of the acknowledged and widespread tools for simulation in large-scale distributed computing.[7]

SimGrid's key features include:[7]

- A scalable and extensible simulation engine that implements several validated simulation models, and that makes it possible to simulate arbitrary network topologies, dynamic computational and network resource availabilities, as well as resource failures;
- High-level user interfaces for researchers (who are not necessarily computer science experts, but rather experts on their own field of research) to quickly assemble simulation prototypes in either C or Java;
- APIs for distributed computing developers to create distributed applications that can run seamlessly in either "simulation mode" or "real-world mode", in order to be able to test it on the simulated environment before actually deploying it.

SimGrid is a very active project, both in terms of research and in terms of development. It is a favored tool by researchers, which is proven by the increasing number of papers written where the research was conducted using SimGrid as a scientific instrument. In terms of development, the developer team envisions a number of directions for future work: addition of a model for disk resources; extension of scalability to improve usability in the P2P domain; ability to dispatch simulated nodes over several physical machines.[7] Another important field of research for the SimGrid team is the implementation of the API that has already been mentioned: the Message-Passing Interface (MPI).

2.5. SMPI

As stated before, MPI is one of the most widely used APIs for communication between nodes in distributed computing. SMPI is a framework for simulating on a single node the execution of parallel applications implemented using the MPI standard. It is part of the SimGrid project and as such, it is built on the SimGrid simulation kernel, benefiting from its fast, scalable and validated network models. SMPI also extends the existing model with other techniques, such as a validated piece-wise linear model for data transfer times between cluster nodes. SMPI simulations also account for network contention - this is done relying on analytical models, rather than packet-level simulation, thus helping the simulation to be fast and scalable.[8]

Three of the main challenges for simulating an MPI application are:

- Accuracy: The prediction of the real-world execution time (the "simulated time") needs to be as accurate as possible, so that reasonable conclusions can be drawn from the experiments.

- Scalability: We want to be able to simulate large-scale applications within a reasonable timescale.
- Speed: It would be advantageous if the simulation time (the actual time of running the simulation) would be as low as possible, compared to the simulated time (the predicted execution time of the real-world application).

Extensive testing was conducted in [8] to verify these qualities of the framework. In these tests, the OpenMPI and MPICH implementations were used to serve as verification benchmarks. The results showed that SMPI predicted the execution time of OpenMPI and MPICH applications for point-to-point, one-to-many and many-to-many applications with an average error value of under 10% in each cases. Using techniques to reduce the memory footprint, SMPI tests were successfully conducted on a scale of up to 448 processors. As a result of the RAM-folding, which was used to reduce the memory footprint of the application, the predicted execution times were underestimates with an average error value of 18.5%, which is higher than in the previous experiments. We have to note here, though, that certain tests weren't successful without the RAM-folding techniques, due to an out-of-memory error. This shows, that although it poses difficulties, reducing the memory usage is vital in SMPI. Also, relying on the CPU sampling approach to reduce the number of computations in an application, SMPI also manages to reduce simulation time by a great margin.

As SMPI is an actively developed project alongside SimGrid, there are a number of research directions. One major development to the project would be a testing framework that would aim to lessen the burdens of testing as much as possible. The goal is to provide a unified method to set up experiments across different environments and to do it with as little necessary adjustments on user part as possible. Another direction is related to the optimization of the framework's performance, with the utilisation of ideas from the STAR-MPI project.

2.6. STAR-MPI

Self-Tuned Adaptive Routines for MPI Collective Operations (STAR-MPI) is a set of MPI collective communication routines that are capable of dynamically adapting to system architecture and application workload. The main idea lays in a technique called "delayed finalization of MPI collective communication routines" (DF). For each operation, STAR-MPI maintains a set of communication algorithms. The aim is to postpone the decision of which algorithm to use until after the platform and/or the application are known. This technique bears the potential of platform-specific or application-specific

optimization of an MPI application.[\[10\]](#)

A development idea for SMPI is to apply the same technique there, thus, in a sense, to "implement" STAR-MPI in SMPI. A set of potentially choosable algorithms could be implemented alongside a set of selector mechanisms. With the right mechanisms, the performance of SMPI could greatly improve.

3. fejezet

Problem Description

3.1. Reproducible research

New scientific ideas, developments and results are only useful when they are documented and published. It is vital that results are announced, so others can be aware of the latest developments on their field of research. This helps in creating a linked data cloud, used by scientists to incorporate various output of other research into their own, using previous results as "stepping stones" to achieve something new.[11] But simply publishing results is not enough in order for others to make use of them. Besides announcing the achievements, the other goal of scientific publications is to convince the readers that the results it presents are correct. Besides theoretical reasoning, papers in experimental science should provide a documented methodology describing how the author has gotten to those results.[12] The methodology has to be detailed and precise enough so other researchers can repeat the same steps, thus reproducing the same results. This is vital in order to provide the possibility to verify those results and to fully understand them. Reproducing the results also makes for a starting point for further development, as the described methods used for reproduction can be extended to achieve something more or something different in the same area of research, or repurposed to gain useful results in a completely different area. This subject is relevant to SMPI and one of the main goals of the testing and validation framework is to make developments in the area.

3.2. Testing framework

SMPI is an actively developed project and as such, a lot of tests are run and a lot of measurements are taken. Previous papers ([8] [13]) have shown, amongst other results, how accurate the performance predictions SMPI makes are and how the time of the

simulation can be lowered, while getting very little differences in simulated time (which means the predicted performance). These results are obtained through extensive testing and as development continues, more and more test data is needed for verification purposes. In order to get conclusive results, both real-life (RL) tests and simulation tests using SimGrid (SG) are needed to be obtained and then the results need to be visualized, analyzed and compared. All this needs to be done in as many different kinds of environments as possible. Currently, this is a tedious task that needs a lot of "tinkering", meaning that there is no unified methodology for setting up experiments and obtaining results on different kinds of distributed environments and for different kinds of MPI implementations - one has to find his/her own way to make it work. Documentation only exists for specific systems, but obviously that doesn't always help with problems arising in an other environment. This statement is reimbursed by [13], where we can read that previous experimental data have been collected "by hand", suggesting that it's a tedious process.

3.3. The testing process

3.3.1. SG traces

In order to conduct tests on the simulator, the first step is to create the simulated environment. This is done by creating a platform file that can be fine-tuned to model the desired system. #maybe detail how a platform file looks like?#

3.3.2. RL traces

As talked about in detail in [14], conducting RL tests involves multiple steps. RL test data collection is done by collecting traces of MPI benchmarks. Currently, the favored tool in trace collection in the project is Tuning and Analysis Utilities (TAU)[15], which is a well-established tracing and profiling tool. Thus, on the system where we are running the benchmarks, TAU has to be deployed and configured, alongside other software that TAU depends on. One is PAPI, an interface which provides us with the possibility to get access to low-level hardware counters (to trace the number of instructions at processor level). We also need the Program Database Toolkit (PDT), which provides the ability of automatic performance instrumentation. In order for TAU to collect the traces we need, these toolkits have to be deployed and correctly linked with TAU. TAU has its own compiler scripts for MPI programs for both Fortran, C and C++. After compiling a benchmark using one of these scripts, they will generate TAU trace files upon their execution. One trace file (.trc) and one event file (.edf) is generated for each MPI

process.

We want to visualize the traces so we can compare the RL and SG traces more easily. Paje is a visualization tool that can be used for this purpose. It has its own trace format that it can comprehend, thus the TAU traces have to be converted to that. Also, we need to merge the traces into one file that we can give to Paje after the conversion. There is a TAU script that is able to do this, creating one trace and one event file. We now only have the task of converting the TAU trace to a Paje trace file. Another MPI tracing library, Akypuera provides this possibility, having its own tau-to-paje conversion script. Once done, we finally have a trace file that Paje can read and display to us.

In the future, it is very likely that more and more tests will have to be run in order to verify old and newly implemented features, as well as various experiments will be conducted using SMPI to test it against various MPI platforms. The main reason of importance of developing a systematic testing and validation framework is that it could make the previously discussed tedious testing process a lot smoother and faster. One of the main features of this framework needs to be the automatic trace collection, so the user wouldn't have to collect and analyze the data manually. This way, proportionally more tests could be run, providing more reliable results with less effort than before. Apart from making extensive testing and experimentation more straightforward for the developers, a functioning test and validation framework would make it easier to reproduce SMPI results, the importance of which is discussed above.

Apart from wanting to simplify this fairly convoluted process of getting test results, there are other problems that need to be addressed.

3.4. Instrumentation overhead

3.5. Clock synchronization

4. fejezet

Implementation

5. fejezet

Evaluation Plan

6. fejezet

Results

7. fejezet

Conclusion

A. Függelék

Appendix Title Here

Write your Appendix content here.

Irodalomjegyzék

- [1] Message Passing Interface, Forum. *MPI: A Message-Passing Interface Standard Version 3.0*, September 2012.
- [2] Gabriel, E., Fagg G., E., Bosilca, G., Angskun, T., Dongarra J., J., Squyres J., M., Sahay, V., Kambadur, P., Barrett, B., Lumsdaine, A., Castain R., H., Daniel D., J., Graham R., L., and Woodall T., S. Open mpi: Goals and concept, and design of a next generation mpi implementation. In *Proceedings, 11th European PVM/MPI Users' Group Meeting*, pages 97–104, Budapest, Hungary, September 2004.
- [3] , MPICH. Mpich overview, 2012. URL <http://www.mpich.org/about/overview/>.
- [4] Thakur, R., Rabenseifner, R., and Gropp, W. Optimization of collective communication operations in mpich. *Int'l Journal of High Performance Computing Applications*, pages 49–66, 2005.
- [5] Buntinas D. Mercier, D. and Gropp, W. Implementation and evaluation of shared-memory communication and synchronization operations in mpich2 using the nemezis communication subsystem. *Parallel Computing*, 33:634–644, September 2007. URL <http://linkinghub.elsevier.com/retrieve/pii/S0167819107000786>.
- [6] Jane, Hillston. Performance modelling and lecture 1 ("modelling and simulation"), 2012. URL <http://www.inf.ed.ac.uk/teaching/courses/pm/PM-lecture1.pdf>.
- [7] Casanova, H., Legrand, A., and Quinson, M. Simgrid: a generic framework for large-scale distributed experiments. In *Proceedings of the Tenth International Event on Computer Modeling and Simulation*, pages 126–131, 2008.
- [8] Clauss, P.-N., Stillwell, M., Genaud, S., Suter, F., Casanova, H., and Quinson, M. Single node on-line simulation of mpi applications with smpi. In *International Parallel & Distributed Processing Symposium*, pages 664–675, May 2011.
- [9] Jane, Hillston. Performance modelling and lecture 13 ("simulation models: Introduction and motivation"), 2012. URL <http://www.inf.ed.ac.uk/teaching/courses/pm/PM-lecture1.pdf>.

- [10] Faraj, A., Yuan, X., and Lowenthal D., K. Star-mpi: Self tuned adaptive routines for mpi collective operations. In *The 20th ACM International Conference on Supercomputing*, pages 199–208, June 2006.
- [11] Bechhofer, S., Ainsworth, J., Bhagat, J., Buchan, I., Couch, P., Cruickshank, D., Roure D., D., Delderfield, M., Dunlop, I., and Gamble, M. Why linked data is not enough for scientists. In *e-Science and 2010 IEEE Sixth International Conference*, pages 300–307, 2010. URL http://ieeexplore.ieee.org/xpls/abs_all.jsp?arnumber=5693931.
- [12] Mesirov J., P. Accessible reproducible research. *Science Magazine*, 327:415–416, 2010. URL <http://statlab.bio5.org/foswiki/pub/Main/PapersForClassCPH685/science-reproducible-research.pdf>.
- [13] Bédaride, P., Degomme, A., Genaud, S., Legrand, A., Markomanolis G., S., Quinson, M., Stillwell, M., Suter, F., and Videau, B. Improving simulations of mpi applications using a hybrid network model with topology and contention support. 2013.
- [14] G. S., Markomanolis and F., Suter. *Time-Independent Trace Acquisition Framework – A Grid’5000 How-to*, 2011. URL <http://hal.inria.fr/inria-00593842>.
- [15] S. S., Shende and A. D., Malony. The tau parallel performance system. *International Journal of High Performance Computing Applications*, 20:287–311, 2006. URL <http://statlab.bio5.org/foswiki/pub/Main/PapersForClassCPH685/science-reproducible-research.pdf>.