

Interpreter and Transpiler for simple expressions on Nvidia GPUs using Julia

Daniel Wiplinger



MASTERARBEIT

eingereicht am
Fachhochschul-Masterstudiengang

Software Engineering

in Hagenberg

im Januar 2025

Advisor:

DI Dr. Gabriel Kronberger

© Copyright 2025 Daniel Wiplinger

This work is published under the conditions of the Creative Commons License *Attribution-NonCommercial-NoDerivatives 4.0 International* (CC BY-NC-ND 4.0)—see <https://creativecommons.org/licenses/by-nc-nd/4.0/>.

Declaration

I hereby declare and confirm that this thesis is entirely the result of my own original work. Where other sources of information have been used, they have been indicated as such and properly acknowledged. I further declare that this or similar work has not been submitted for credit elsewhere. This printed copy is identical to the submitted electronic version.

Hagenberg, January 1, 2025

Daniel Wiplinger

Contents

Declaration	iv
Abstract	vii
Kurzfassung	viii
1 Introduction	1
1.1 Background and Motivation	1
1.2 Research Question	2
1.3 Methodology	2
2 Fundamentals and Related Work	4
2.1 Equation learning	4
2.2 GPGPU	5
2.2.1 Programming GPUs	6
2.2.2 PTX	6
2.3 Compilers	6
2.3.1 Interpreters	7
2.3.2 Transpilers	7
3 Concept and Design	8
3.1 Requirements	8
3.2 Interpreter	8
3.2.1 Architecture	8
3.2.2 Host	8
3.2.3 Device	8
3.3 Transpiler	8
3.3.1 Architecture	9
3.3.2 Host	9
3.3.3 Device	9
4 Implementation	10
4.1 Technologies	10
4.2 Interpreter	10
4.3 Transpiler	10

Contents	vi
5 Evaluation	11
5.1 Test environment	11
5.2 Results	11
5.2.1 Interpreter	11
5.2.2 Transpiler	11
5.2.3 Comparison	11
6 Conclusion	12
6.1 Future Work	12
References	13
Literature	13

Abstract

This should be a 1-page (maximum) summary of your work in English.

Kurzfassung

An dieser Stelle steht eine Zusammenfassung der Arbeit, Umfang max. 1 Seite. ...

Chapter 1

Introduction

This chapter provides an entry point for this thesis. First the motivation of exploring this topic is presented. In addition, the research questions of this thesis are outlined. Lastly the methodology on how to answer these questions will be explained.

1.1 Background and Motivation

Optimisation and acceleration of program code is a crucial part in many fields. For example video games need optimisation to lower the minimum hardware requirements which allows more people to run the game, increasing sales. Another example where optimisation is important are computer simulations. For those, optimisation is even more crucial, as this allows the scientists to run more detailed simulations or get the simulation results faster. Equation learning is another field that can heavily benefit from optimisation. One part of equation learning, is to evaluate the expressions generated by the algorithm which can make up a significant portion of the runtime of the algorithm. This thesis is concerned with optimising the evaluation part to increase the overall performance of the equation learning algorithm.

Considering the following expression $x_1 + 5 - \text{abs}(p_1) * \text{sqrt}(x_2)/10 + 2^3$ which contains simple mathematical operations as well as variables x_n and parameters p_n . This expression is one example that can be generated by the equation learning algorithm and needs to be evaluated for the next iteration. Usually multiple expressions are generated per iteration, which also need to be evaluated. Additionally, multiple different values need to be inserted for all variables and parameters, drastically increasing the amount of evaluations that need to be performed.

The free lunch theorem as described by Adam et al. (2019) states that to gain additional performance, a developer cannot just hope for future hardware to be faster, especially on a single core. Therefore, algorithms need to utilise the other cores on a processor to further acceleration. While this approach means more development overhead, a much greater speed-up can be achieved. However, in some cases the speed-up achieved by this is still not large enough and another approach is needed. One of these approaches is the utilisation of a Graphics Processing Unit (GPU) as an easy and affordable option as compared to compute clusters. Michalakes and Vachharajani (2008) have shown a noticeable speed-up when using the GPU for weather simulation. In ad-

dition to computer simulations GPU acceleration also can be found in other places like networking (S. Han et al., 2010) or structural analysis of buildings (Georgescu et al., 2013).

1.2 Research Question

With these successful implementations of GPU acceleration, this thesis also attempts to improve the performance of evaluating mathematical equations using GPUs. Therefore, the following research questions are formulated:

- How can simple arithmetic expressions that are generated at runtime be efficiently evaluated on graphics cards?
- Under what circumstances is the evaluation of simple arithmetic expressions faster on a graphics card than on a CPU?
- Under which circumstances is the interpretation of the expressions on the GPU or the translation to the intermediate language Parallel Thread Execution (PTX) more efficient?

Answering the first question is necessary to ensure the approach of this thesis is actually feasible. If it is feasible, it is important to evaluate if evaluating the expressions on the GPU actually improves the performance over a parallelised CPU evaluator. To answer if the GPU evaluator is faster than the CPU evaluator, the last research question is important. As there are two major ways of implementing an evaluator on the GPU, they need to be implemented and evaluated to finally state if evaluating expressions on the GPU is faster and if so, which type of implementation results in the best performance.

1.3 Methodology

In order to answer the research questions, this thesis is divided into the following chapters:

Chapter 2: Fundamentals and Related Work

In this chapter, the topic of this thesis is explored. It covers the fundamentals of equation learning and how this thesis fits into this field of research. In addition, the fundamentals of General Purpose GPU computing and how interpreters and transpilers work are explained. Previous research already done within this topic is also explored.

Chapter 3: Concept and Design

Within this chapter, the concepts of implementing the GPU interpreter and transpiler are explained. How these two prototypes can be implemented disregarding concrete technologies is part of this chapter.

Chapter 4: Implementation

This chapter explains the implementation of the GPU interpreter and transpiler. The details of the implementation with the used technologies are covered, such as the interpretation process and the transpilation of the expressions into Parallel Thread Execution (PTX) code.

Chapter 5: Evaluation

The software and hardware requirements and the evaluation environment are introduced in this chapter. Furthermore, the results of the comparison of the GPU and CPU evaluators are presented to show which of these yields the best performance.

Chapter 6: Conclusion

In the final chapter, the entire work is summarised. A brief overview of the implementation as well as the evaluation results will be provided. Additionally, an outlook of possible future research is given.

With this structure the process of creating and evaluating a basic interpreter on the GPU as well as a transpiler for creating PTX code is outlined. Research is done to ensure the implementations are relevant and not outdated. Finally, the evaluation results will answer the research questions and determine if expressions generated at runtime can be evaluated more efficiently on the GPU than on the CPU.

Chapter 2

Fundamentals and Related Work

The goal of this chapter is to provide an overview of equation learning to establish common knowledge of the topic and problem this thesis is trying to solve. The main part of this chapter is split into two parts. The first part is exploring research that has been done in the field of general purpose computations on the GPU (GPGPU) as well as the fundamentals of it. Focus lies on exploring how graphics processing units (GPUs) are used to achieve substantial speed-ups and when they can be effectively employed. The second part describes the basics of how interpreters and compilers are built and how they can be adapted to the workflow of programming GPUs.

2.1 Equation learning

Equation learning is a field of research that aims at understanding and discovering equations from a set of data from various fields like mathematics and physics. Data is usually much more abundant while models often are elusive. Because of this, generating equations with a computer can more easily lead to discovering equations that describe the observed data. Brunton et al. (2016) describe an algorithm that leverages equation learning to discover equations for physical systems. A more literal interpretation of equation learning is demonstrated by Pfahler and Morik (2020). They use machine learning to learn the form of equations. Their aim was to simplify the discovery of relevant publications by the equations they use and not by technical terms, as they may differ by the field of research. However, this kind of equation learning is not relevant for this thesis.

Symbolic regression is a subset of equation learning, that specialises more towards discovering mathematical equations. A lot of research is done in this field. Keijzer (2004) and Korn (2011) presented ways of improving the quality of symbolic regression algorithms, making symbolic regression more feasible for problem-solving. Additionally, Jin et al. (2020) proposed an alternative to genetic programming (GP) for the use in symbolic regression. Their approach increased the quality of the results noticeably compared to GP alternatives. The first two approaches are more concerned with the quality of the output, while the third is also concerned with interpretability and reducing memory consumption. Bartlett et al. (2024) also describe an approach to generate simpler and higher quality equations while being faster than GP algorithms. Heuristics like GP or

neural networks as used by Werner et al. (2021) in their equation learner can help with finding good solutions faster, accelerating scientific progress. As seen by these publications, increasing the quality of generated equations but also increasing the speed of finding these equations is a central part in symbolic regression and equation learning in general. This means research in improving the computational performance of these algorithms is desired.

The expressions generated by an equation learning algorithm can look like this $x_1 + 5 - \text{abs}(p_1) * \text{sqrt}(x_2) / 10 + 2^3$. They consist of several unary and binary operators but also of constants, variables and parameters and expressions mostly differ in length and the kind of terms in the expressions. Per iteration many of these expressions are generated and in addition, matrices of values for the variables and parameters are also created. One row of the variable matrix corresponds to one instantiation of all expressions and this matrix contains multiple rows. This leads to a drastic increase of instantiated expressions that need to be evaluated. Parameters are a bit simpler, as they can be treated as constants for one iteration but can have a different value on another iteration. This means that parameters do not increase the number of expressions that need to be evaluated. However, the increase in evaluations introduced by the variables is still drastic and therefore increases the algorithm runtime significantly.

2.2 General Purpose Computation on Graphics Processing Units

Graphics cards (GPUs) are commonly used to increase the performance of many different applications. Originally they were designed to improve performance and visual quality in games. Dokken et al. (2005) first described the usage of GPUs for general purpose programming. They have shown how the graphics pipeline can be used for GPGPU programming. Because this approach also requires the programmer to understand the graphics terminology, this was not a great solution. Therefore, Nvidia released CUDA¹ in 2007 with the goal of allowing developers to program GPUs independent of the graphics pipeline and terminology. A study of the programmability of GPUs with CUDA and the resulting performance has been conducted by Huang et al. (2008). They found that GPGPU programming has potential, even for non-embarrassingly parallel problems. Research is also done in making the low level CUDA development simpler. T. D. Han and Abdelrahman (2011) have described a directive-based language to make development simpler and less error-prone, while retaining the performance of handwritten code. To drastically simplify CUDA development Besard, Foket, et al. (2019) showed that it is possible to develop with CUDA in the high level programming language Julia² while performing similar to CUDA written in C. In a subsequent study Lin and McIntosh-Smith (2021) found that high performance computing (HPC) on the CPU and GPU in Julia performs similar to HPC development in C. This means that Julia can be a viable alternative to Fortran, C and C++ in the HPC field and has the additional benefit of developer comfort since it is a high level language with modern features such as garbage-collectors. Besard, Churavy, et al. (2019) have also shown how the combination of Julia and CUDA help in rapidly developing HPC software. While this

¹<https://developer.nvidia.com/cuda-toolkit>

²<https://julialang.org/>

thesis in general revolves around CUDA, there also exist alternatives by AMD called ROCm³ and a vendor independent alternative called OpenCL⁴.

While in the early days of GPGPU programming a lot of research has been done to assess if this approach is feasible, it now seems obvious to use GPUs to accelerate algorithms. Weather simulations began using GPUs very early for their models. In 2008 Michalakes and Vachharajani (2008) proposed a method for simulating weather with the WRF model on a GPU. With their approach, they reached a speed-up of the most compute intensive task of 5 to 20, with very little GPU optimisation effort. They also found that the GPU usages was very low, meaning there are resources and potential for more detailed simulations. Generally, simulations are great candidates for using GPUs, as they can benefit heavily from a high degree of parallelism and data throughput. Köster et al. (2020) have developed a way of using adaptive time steps to improve the performance of time step simulations, while retaining their precision and constraint correctness. Black hole simulations are crucial for science and education for a better understanding of our world. Verbraeck and Eisemann (2021) have shown that simulating complex Kerr (rotating) black holes can be done on consumer hardware in a few seconds. Schwarzschild black hole simulations can be performed in real-time with GPUs as described by Hissbach et al. (2022) which is especially helpful for educational scenarios. While both approaches do not have the same accuracy as detailed simulations on supercomputers, they show how single GPUs can yield similar accuracy at a fraction of the cost. Networking can also heavily benefit from GPU acceleration as shown by S. Han et al. (2010), where they achieved a significant increase in throughput than with a CPU only implementation. Finite element structural analysis is an essential tool for many branches of engineering and can also heavily benefit from the usage of GPUs as demonstrated by Georgescu et al. (2013).

2.2.1 Programming GPUs

talk about the fields GPGPU really helped make performance improvements (weather simulations etc). Then describe how it differs from classical programming. talk about architecture (SIMD/SIMT; a lot of “slow” cores).

starting from here I can hopefully incorporate more images to break up these walls of text

2.2.2 Parallel Thread Execution

Describe what PTX is to get a common ground for the implementation chapter. Probably a short section

2.3 Compilers

brief overview about compilers (just setting the stage for the subsections basically). Talk about register management and these things

³<https://www.amd.com/de/products/software/rocm.html>

⁴<https://www.khronos.org/opencl/>

2.3.1 Interpreters

What are interpreters; how they work; should mostly contain/reference gpu interpreters

2.3.2 Transpilers

talk about what transpilers are and how to implement them. If possible also gpu specific transpilation.

Chapter 3

Concept and Design

introduction to what needs to be done. also clarify terms “Host” and “Device” here

3.1 Requirements and Data

short section. Multiple expressions; vars for all expressions; params unique to expression; operators that need to be supported

3.2 Interpreter

as introduction to this section talk about what “interpreter” means in this context. so “gpu parses expr and calculates”

3.2.1 Architecture

talk about the coarse grained architecture on how the interpreter will work. (.5 to 1 page probably)

3.2.2 Host

talk about the steps taken to prepare for GPU interpretation

3.2.3 Device

talk about how the actual interpreter will be implemented

3.3 Transpiler

as introduction to this section talk about what “transpiler” means in this context. so “cpu takes expressions and generates ptx for gpu execution”

3.3.1 Architecture

talk about the coarse grained architecture on how the transpiler will work. (.5 to 1 page probably)

3.3.2 Host

talk about how the transpiler is implemented

3.3.3 Device

talk about what the GPU does. short section since the gpu does not do much

Chapter 4

Implementation

4.1 Technologies

Short section; CUDA, PTX, Julia, CUDA.jl

Probably reference the performance evaluation papers for Julia and CUDA.jl

4.2 Interpreter

Talk about how the interpreter has been developed.

4.3 Transpiler

Talk about how the transpiler has been developed

Chapter 5

Evaluation

5.1 Test environment

Explain the hardware used, as well as the actual data (how many expressions, variables etc.)

5.2 Results

talk about what we will see now (results only for interpreter, then transpiler and then compared with each other and a CPU interpreter)

5.2.1 Interpreter

Results only for Interpreter

5.2.2 Transpiler

Results only for Transpiler

5.2.3 Comparison

Comparison of Interpreter and Transpiler as well as Comparing the two with CPU interpreter

Chapter 6

Conclusion and Future Work

Summarise the results

6.1 Future Work

talk about what can be improved

References

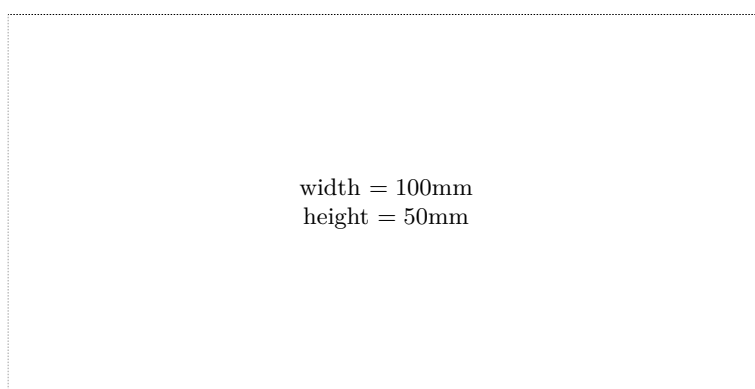
Literature

- Adam, S. P., Alexandropoulos, S.-A. N., Pardalos, P. M., & Vrahatis, M. N. (2019). No free lunch theorem: A review. In I. C. Demetriou & P. M. Pardalos (Eds.), *Approximation and optimization : Algorithms, complexity and applications* (pp. 57–82). Springer International Publishing. https://doi.org/10.1007/978-3-030-12767-1_5. (Cit. on p. 1)
- Bartlett, D. J., Desmond, H., & Ferreira, P. G. (2024). Exhaustive symbolic regression [Conference Name: IEEE Transactions on Evolutionary Computation]. *IEEE Transactions on Evolutionary Computation*, *28*(4), 950–964. <https://doi.org/10.1109/TEVC.2023.3280250> (cit. on p. 4)
- Besard, T., Churavy, V., Edelman, A., & Sutter, B. D. (2019). Rapid software prototyping for heterogeneous and distributed platforms. *Advances in Engineering Software*, *132*, 29–46. <https://doi.org/10.1016/j.advengsoft.2019.02.002> (cit. on p. 5)
- Besard, T., Foket, C., & De Sutter, B. (2019). Effective extensible programming: Unleashing julia on GPUs. *IEEE Transactions on Parallel and Distributed Systems*, *30*(4), 827–841. <https://doi.org/10.1109/TPDS.2018.2872064> (cit. on p. 5)
- Brunton, S. L., Proctor, J. L., & Kutz, J. N. (2016). Discovering governing equations from data by sparse identification of nonlinear dynamical systems [Publisher: Proceedings of the National Academy of Sciences]. *Proceedings of the National Academy of Sciences*, *113*(15), 3932–3937. <https://doi.org/10.1073/pnas.1517384113> (cit. on p. 4)
- Dokken, T., Hagen, T. R., & Hjelmervik, J. M. (2005). The GPU as a high performance computational resource. *Proceedings of the 21st Spring Conference on Computer Graphics*, 21–26. <https://doi.org/10.1145/1090122.1090126> (cit. on p. 5)
- Georgescu, S., Chow, P., & Okuda, H. (2013). GPU acceleration for FEM-based structural analysis. *Archives of Computational Methods in Engineering*, *20*(2), 111–121. <https://doi.org/10.1007/s11831-013-9082-8> (cit. on pp. 2, 6)
- Han, S., Jang, K., Park, K., & Moon, S. (2010). PacketShader: A GPU-accelerated software router. *SIGCOMM Comput. Commun. Rev.*, *40*(4), 195–206. <https://doi.org/10.1145/1851275.1851207> (cit. on pp. 2, 6)
- Han, T. D., & Abdelrahman, T. S. (2011). hiCUDA: High-level GPGPU programming [Conference Name: IEEE Transactions on Parallel and Distributed Systems]. *IEEE Transactions on Parallel and Distributed Systems*, *22*(1), 78–90. Retrieved

- March 1, 2025, from <https://ieeexplore.ieee.org/abstract/document/5445082> (cit. on p. 5)
- Hissbach, A.-M., Dick, C., & Lawonn, K. (2022). *An overview of techniques for egocentric black hole visualization and their suitability for planetarium applications*. The Eurographics Association. Retrieved March 2, 2025, from <https://doi.org/10.2312/vmv.20221207>. (Cit. on p. 6)
- Huang, Q., Huang, Z., Werstein, P., & Purvis, M. (2008). GPU as a general purpose computing resource [ISSN: 2379-5352]. *2008 Ninth International Conference on Parallel and Distributed Computing, Applications and Technologies*, 151–158. <https://doi.org/10.1109/PDCAT.2008.38> (cit. on p. 5)
- Jin, Y., Fu, W., Kang, J., Guo, J., & Guo, J. (2020, January 16). Bayesian symbolic regression. <https://doi.org/10.48550/arXiv.1910.08892>. (Cit. on p. 4)
- Keijzer, M. (2004). Scaled symbolic regression. *Genetic Programming and Evolvable Machines*, 5(3), 259–269. <https://doi.org/10.1023/B:GENP.0000030195.77571.f9> (cit. on p. 4)
- Korns, M. F. (2011). Accuracy in symbolic regression. In R. Riolo, E. Vladislavleva, & J. H. Moore (Eds.), *Genetic programming theory and practice IX* (pp. 129–151). Springer. https://doi.org/10.1007/978-1-4614-1770-5_8. (Cit. on p. 4)
- Köster, M., Groß, J., & Krüger, A. (2020). High-performance simulations on GPUs using adaptive time steps. In M. Qiu (Ed.), *Algorithms and architectures for parallel processing* (pp. 369–385). Springer International Publishing. https://doi.org/10.1007/978-3-030-60245-1_26. (Cit. on p. 6)
- Lin, W.-C., & McIntosh-Smith, S. (2021). Comparing julia to performance portable parallel programming models for HPC. *2021 International Workshop on Performance Modeling, Benchmarking and Simulation of High Performance Computer Systems (PMBS)*, 94–105. <https://doi.org/10.1109/PMBS54543.2021.00016> (cit. on p. 5)
- Michalakes, J., & Vachharajani, M. (2008). GPU acceleration of numerical weather prediction [ISSN: 1530-2075]. *2008 IEEE International Symposium on Parallel and Distributed Processing*, 1–7. <https://doi.org/10.1109/IPDPS.2008.4536351> (cit. on pp. 1, 6)
- Pfahler, L., & Morik, K. (2020). Semantic search in millions of equations. *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 135–143. <https://doi.org/10.1145/3394486.3403056> (cit. on p. 4)
- Verbraeck, A., & Eisemann, E. (2021). Interactive black-hole visualization [Conference Name: IEEE Transactions on Visualization and Computer Graphics]. *IEEE Transactions on Visualization and Computer Graphics*, 27(2), 796–805. <https://doi.org/10.1109/TVCG.2020.3030452> (cit. on p. 6)
- Werner, M., Junginger, A., Hennig, P., & Martius, G. (2021, May 13). Informed equation learning. <https://doi.org/10.48550/arXiv.2105.06331>. (Cit. on p. 5)

Check Final Print Size

— Check final print size! —



— Remove this page after printing! —