

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	3
2.1 Equation learning	3
2.2 GPGPU	3
2.2.1 PTX	3
2.3 GPU Interpretation	3
2.4 Transpiler	3
3 Concept and Design	4
3.1 Requirements	4
3.2 Interpreter	4
3.2.1 Architecture	4
3.2.2 Host	4
3.2.3 Device	4
3.3 Transpiler	4
3.3.1 Architecture	5
3.3.2 Host	5
3.3.3 Device	5
4 Implementation	6
4.1 Technologies	6
4.2 Interpreter	6
4.3 Transpiler	6
5 Evaluation	7
5.1 Test environment	7

Contents	vi
5.2 Results	7
5.2.1 Interpreter	7
5.2.2 Transpiler	7
5.2.3 Comparison	7
6 Conclusion	8
6.1 Future Work	8
References	9
Literature	9

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 of program code is a crucial part in many different fields. For example video games need a lot of optimisation to lower the minimum hardware requirements which allows more people to run the game. Another example for optimisation 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. This thesis is concerned with optimising that part to increase the overall performance of the equation learning algorithm. The free lunch theorem as described by Adam et al. (2019) states that optimising a program which runs on a single core will eventually lead to a performance plateau as no further optimisations can be done. Therefore, these algorithms need to utilise the other cores on a processor. 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) to further increase the performance. Michalakes and Vachharajani (2008) have shown a noticeable speed-up when using the GPU for weather simulation. In addition to simulation GPU acceleration also can be found in other places like networking (Han et al., 2010) or structural analysis of buildings (Georgescu et al., 2013).

With these successful implementations of GPU acceleration, this thesis also attempts to improve the performance of evaluating mathematical equations using GPUs. The baseline to compare against is an expression evaluator running in parallel on the CPU. (talk a bit more about what will be attempted I think. Look at other papers to see what they write in this section.)

1.2 Research Question

What are the research questions and how they will be answered

1.3 Methodology

Will give an overview of the chapters and what to expect

Chapter 2

Fundamentals and Related Work

2.1 Equation learning

Section describing what equation learning is and why it is relevant for the thesis

2.2 General Purpose Computation on Graphics Processing Units

Describe what GPGPU is and how it differs from classical programming. talk about architecture (SIMD) and some scientific papers on how they use GPUs to accelerate tasks

2.2.1 Parallel Thread Execution

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

2.3 GPU Interpretation

Different sources on how to do interpretation on the gpu (and maybe interpretation in general too?)

2.4 Transpiler

talk about what transpilers are and how to implement them. If possible also gpu specific transpilation. Also talk about compilation and register management. and probably find a better title

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

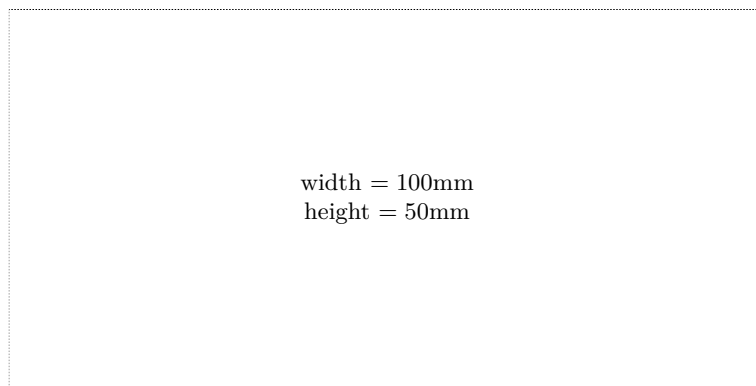
References

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)
- 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 p. 1)
- 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 p. 1)
- 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 p. 1)

Check Final Print Size

— Check final print size! —



— Remove this page after printing! —