

# Interpreter and Transpiler for simple expressions on Nvidia GPUs using Julia

Daniel Roth



MASTERARBEIT

eingereicht am  
Fachhochschul-Masterstudiengang

Software Engineering

in Hagenberg

im Januar 2025

Advisor:

DI Dr. Gabriel Kronberger

© Copyright 2025 Daniel Roth

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 Roth

# Contents

<b>Declaration</b>	<b>iv</b>
<b>Abstract</b>	<b>vii</b>
<b>Kurzfassung</b>	<b>viii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Background and Motivation . . . . .	1
1.2 Research Question . . . . .	2
1.3 Thesis Structure . . . . .	2
<b>2 Fundamentals and Related Work</b>	<b>4</b>
2.1 Equation learning . . . . .	4
2.2 GPGPU . . . . .	6
2.2.1 Programming GPUs . . . . .	8
2.2.2 PTX . . . . .	16
2.3 Compilers . . . . .	16
2.3.1 Interpreters . . . . .	17
2.3.2 Transpilers . . . . .	17
<b>3 Concept and Design</b>	<b>18</b>
3.1 Requirements . . . . .	18
3.2 Interpreter . . . . .	18
3.2.1 Architecture . . . . .	18
3.2.2 Host . . . . .	18
3.2.3 Device . . . . .	18
3.3 Transpiler . . . . .	18
3.3.1 Architecture . . . . .	19
3.3.2 Host . . . . .	19
3.3.3 Device . . . . .	19
<b>4 Implementation</b>	<b>20</b>
4.1 Technologies . . . . .	20
4.2 Interpreter . . . . .	20
4.3 Transpiler . . . . .	20

<b>5 Evaluation</b>	<b>21</b>
5.1 Test environment . . . . .	21
5.2 Results . . . . .	21
5.2.1 Interpreter . . . . .	21
5.2.2 Transpiler . . . . .	21
5.2.3 Comparison . . . . .	21
<b>6 Conclusion</b>	<b>22</b>
6.1 Future Work . . . . .	22
<b>References</b>	<b>23</b>
Literature . . . . .	23
Online sources . . . . .	26

# 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 or symbolic regression is another field that can heavily benefit from optimisation. One part of equation learning, is to evaluate the expressions generated by a search algorithm which can make up a significant portion of the runtime. This thesis is concerned with optimising the evaluation part to increase the overall performance of equation learning algorithms.

The following expression  $5 - \text{abs}(x_1) * \text{sqrt}(x_2)/10 + 2^x x_3$  which contains simple mathematical operations as well as variables  $x_n$  and parameters  $p_n$  is one example that can be generated by the equation learning algorithm, Usually an equation learning algorithm generates multiple of such expressions per iteration. Out of these expressions all possibly relevant ones have 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.

In his Blog Sutter (2004) described how the free lunch is over in terms of the ever-increasing performance of hardware like the CPU. He states that to gain additional performance, developers need to start developing software for multiple cores and not just hope that on the next generation of CPUs the program magically runs faster. 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 Graphics Processing Units (GPUs) as an easy and affordable option as compared to compute clusters. Especially when talking about performance per dollar, GPUs are very

inexpensive as found by Brodtkorb et al. (2013). Michalakes and Vachharajani (2008) have shown a noticeable speed-up when using GPUs for weather simulation. In addition to computer simulations, GPU acceleration also can be found in other places such as 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 GPUs?
- Under what circumstances is the evaluation of simple arithmetic expressions faster on a GPU 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 Thesis Structure

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. All three evaluators will be compared against each other and the form of the expressions used for the comparisons are outlined. Finally, 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 or symbolic regression to establish common knowledge of the topic and problem this thesis is trying to solve. First the field of equation learning is explored which helps to contextualise the topic of this thesis. The main part of this chapter is split into two sub-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 and where 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. When discussing GPU programming concepts, the terminology used is that of Nvidia and may differ from that used for AMD GPUs.

### 2.1 Equation learning

Equation learning is a field of research that can be used for 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 which is demonstrated by Guillemot (2022) where they explain how validating the models against large amounts of data is a big part in creating such models. Because of this effort, 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. Using genetic programming (GP) for different problems, including symbolic regression, was first described by Koza (1994). He described that finding a computer program to solve a problem for a given input and output, can be done by traversing the search space of all solutions. This fits well for the goal of symbolic regression, where a mathematical

expression needs to be found to describe a problem with specific inputs and outputs. Later, Koza (2010) provided an overview of results that were generated with the help of GP and were competitive with human solutions, showing how symbolic regression is a useful tool. In their book *Symbolic Regression*, Kronberger et al. (2024) show how symbolic regression can be applied for real world scenarios. They also describe symbolic regression in great detail, while being tailored towards beginners and experts.

Keijzer (2004) and Korns (2011) presented ways of improving the quality of symbolic regression algorithms, making symbolic regression more feasible for problem-solving. Bartlett et al. (2024) describe an exhaustive approach for symbolic regression which can find the true optimum for perfectly optimised parameters while retaining simple and interpretable results. Alternatives to GP for symbolic regression also exist with one proposed by Jin et al. (2020). Their approach increased the quality of the results noticeably compared to GP alternatives. Another alternative to heuristics like GP is the usage of neural networks. One such alternative has been introduced by Martius and Lampert (2016) where they used a neural network for their equation learner with mixed results. Later, an extension has been provided by Sahoo et al. (2018). They introduced the division operator, which led to much better results. Further improvements have been described by Werner et al. (2021) with their informed equation learner. By incorporating domain expert knowledge they could limit the search space and find better solutions for particular domains. One drawback of these three implementations is the fact that their neural networks are fixed. An equation learner which can change the network at runtime and therefore evolve over time is proposed by Dong et al. (2024). Their approach further improved the results of neural network equation learners. In their work, Lemos et al. (2022) also used a neural network for symbolic regression. They were able to find an equivalent to Newton's law of gravitation and rediscovered Newton's second and third law only with trajectory data of bodies of our solar system. Although these laws were already known, this research has shown how neural networks and machine learning in general have great potential. An implementation for an equation learner in the physics domain is proposed by Sun et al. (2023). Their algorithm was specifically designed for nonlinear dynamics often occurring in physical systems. When compared to other implementations their equation learner was able to create better results but have the main drawback of high computational cost. As seen by these publications, increasing the quality of generated equations and also increasing the speed of finding these equations is a central part in symbolic regression and equation learning in general.

As described earlier, the goal of equation learning is to find an expression that fits a given set of data. The data usually consists of a set of inputs that have been applied to the unknown expression and the output after the input has been applied. An example for such data is described by Werner et al. (2021). In one instance they want to find the power loss formula for an electric machine. They used four inputs, direct and quadratic current as well as temperature and motor speed, and they have an observed output which is the power loss. Now for an arbitrary problem with different input and outputs, the equation learner tries to find an expression that fits this data (Koza, 1994). Fitting in this context means that when the input is applied to the expression, the result will be the same as the observed output. In order to avoid overfitting Bomarito et al. (2022) have proposed a way of using Bayesian model selection to combat overfitting and reduce the complexity of the generated expressions. This also helps with making the expressions

more generalisable and therefore be applicable to unseen inputs. A survey conducted by Dabhi and Chaudhary (2012) shows how overfitting is not desirable and why more generalisable solutions are preferred. To generate an equation, first the operators need to be defined that make up the equation. It is also possible to define a maximum length for an expression as proposed by Bartlett et al. (2024). Expressions also consist of constants as well as variables which represent the inputs. Assuming that a given problem has three variables, the equation learner could generate an expression as seen in 2.1 where  $x_n$  are the variables and  $O$  is the output which should correspond to the observed output for the given variables.

$$O = 5 - \text{abs}(x_1) * \text{sqrt}(x_2)/10 + 2^x x_3 \quad (2.1)$$

A typical equation learner generates multiple expressions at once. If the equation learner generates 300 expressions and each expression needs to be evaluated 50 times to get the best parametrisation for this expression, the total number of evaluations is  $300 * 50 = 15\,000$ . However, it is likely that multiple runs or generations in the context of GP need to be performed. The number of generations is dependent to the problem, but assuming a maximum of 100 generations the total number of evaluations is equal to  $300 * 50 * 100 = 1\,500\,000$ . These values have been taken from the equation learner for predicting discharge voltage curves of batteries as described by Kronberger et al. (2024). Their equation learner converged after 54 generations, resulting in evaluating 800 000 expressions. Depending on the complexity of the generated expressions, performing all of these evaluations takes up a lot of the runtime. Their results took over two days on an eight core desktop CPU. While they did not provide runtime information for all problems they tested, the voltage curve prediction was the slowest. The other problems were in the range of a few seconds and up to a day. Especially the problems that took several hours to days to finish show, that there is still room for performance improvements. While a better CPU with more cores can be used, it is interesting to determine, if using Graphics cards can yield noticeable better performance or not, which is the goal of this thesis.

## 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 (GPGPU). 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 CUDA<sup>1</sup> 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

---

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

development simpler and less error-prone, while retaining the performance of handwritten code. To drastically simplify CUDA development, Besard et al. (2019b) showed that it is possible to develop with CUDA in the high level programming language Julia<sup>2</sup> with similar performance 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 et al. (2019a) have also shown how the combination of Julia and CUDA help in rapidly developing HPC software. While this thesis in general revolves around CUDA, there also exist alternatives by AMD called ROCm<sup>3</sup> and a vendor independent alternative called OpenCL<sup>4</sup>. If not specified otherwise, the following section and its subsections use the information presented by Nvidia (2025b) in their CUDA programming guide.

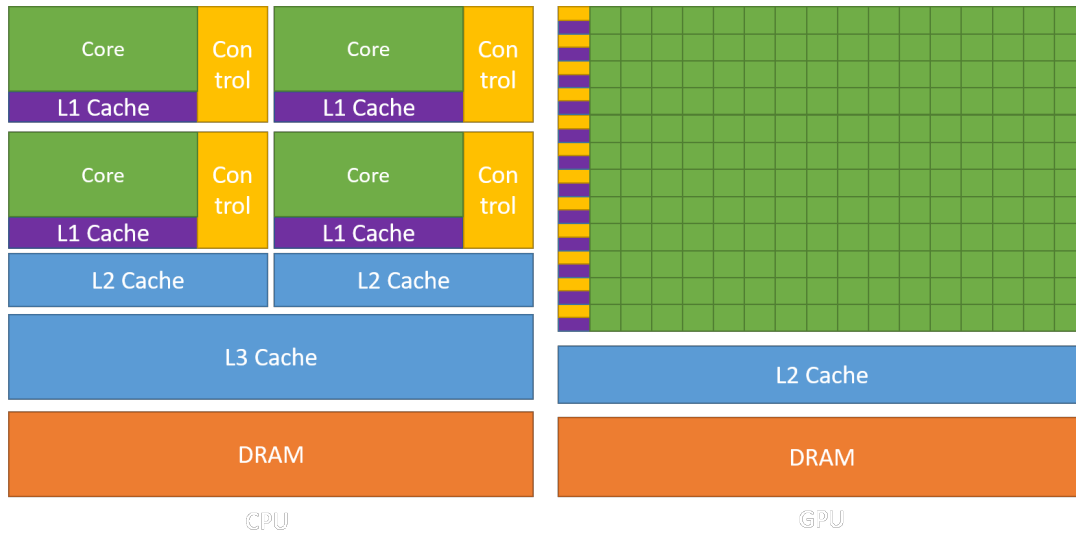
While in the early days of GPGPU programming a lot of research has been done to assess if this approach is feasible, it now seems obvious to use GPUs to accelerate algorithms. GPUs have been used early to speed up weather simulation models. Michalakes and Vachharajani (2008) proposed a method for simulating weather with the Weather Research and Forecast (WRF) model on a GPU. With their approach, they reached a speed-up of the most compute intensive task of 5 to 20, with little GPU optimisation effort. They also found that the GPU usage was 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. (2020b) have developed a way of using adaptive time steps on the GPU to considerably improve the performance of numerical and discrete simulations. In addition to the performance gains they were able to retain the precision and constraint correctness of the simulation. 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 a single GPU can yield similar accuracy at a fraction of the cost. Software network routing can also heavily benefit from GPU acceleration as shown by S. Han et al. (2010), where they achieved a significantly higher 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). However, it also needs to be noted, that GPUs are not always better performing than CPUs as illustrated by Lee et al. (2010), but they still can lead to performance improvements nonetheless.

---

<sup>2</sup><https://julialang.org/>

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

<sup>4</sup><https://www.khronos.org/opencl/>



**Figure 2.1:** Overview of the architecture of a CPU (left) and a GPU (right). Note the higher number of simpler and smaller cores on the GPU (Nvidia, 2025b).

### 2.2.1 Programming GPUs

The development process on a GPU is vastly different from a CPU. A CPU has tens or hundreds of complex cores with the AMD Epyc 9965<sup>5</sup> having a staggering 192 of those complex cores and twice as many threads. A guide for a simple one core 8-bit CPU has been published by Schuurman (2013). He describes the different and complex parts of a CPU core. Modern CPUs are even more complex, with dedicated fast integer and floating-point arithmetic gates as well as logic gates, sophisticated branch prediction and much more. This makes a CPU perfect for handling complex control flows on a single program strand and on modern CPUs even multiple strands simultaneously. However, as seen in section 2.2, this often isn't enough. On the other hand, a GPU contains thousands or even tens of thousands of cores. For example, the GeForce RTX 5090<sup>6</sup> contains a total of 21760 CUDA cores. To achieve this enormous core count a single GPU core has to be much simpler than one CPU core. As described by Nvidia (2025b) a GPU designates much more transistors towards floating-point computations. This results in less efficient integer arithmetic and control flow handling. There is also less Cache available per core and clock speeds are usually also much lower than those on a CPU. An overview of the differences of a CPU and a GPU architecture can be seen in figure 2.1.

Despite these drawbacks, the sheer number of cores, makes a GPU a valid choice when considering improving the performance of an algorithm. Because of the high number of cores, GPUs are best suited for data parallel scenarios. This is due to the SIMD architecture of these cards. SIMD stands for Single-Instruction Multiple-Data and states that there is a single stream of instructions that is executed on a huge number of data

<sup>5</sup><https://www.amd.com/en/products/processors/server/epyc/9005-series/amd-epyc-9965.html>

<sup>6</sup><https://www.nvidia.com/en-us/geforce/graphics-cards/50-series/rtx-5090/>



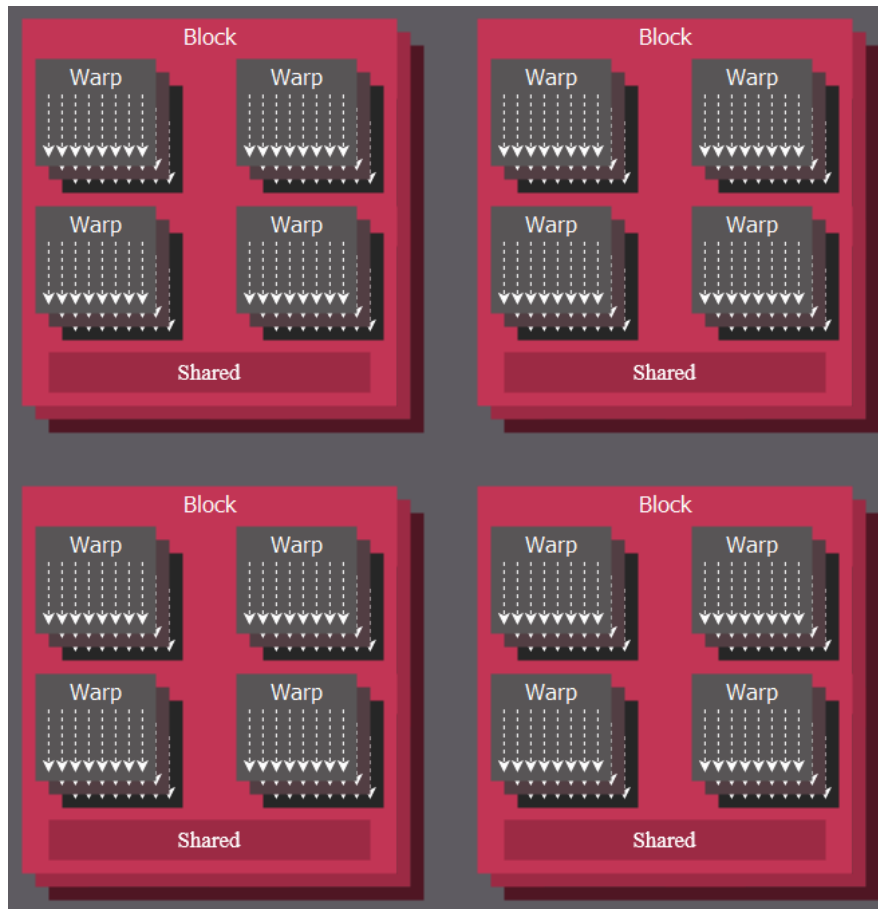
streams. Franchetti et al. (2005) and Tian et al. (2012) describe ways of using SIMD instructions on the CPU. Their approaches lead to noticeable speed-ups of 3.3 and 4.7 respectively by using SIMD instructions instead of serial computations. Extending this to GPUs which are specifically built for SIMD/data parallel calculations shows why they are so powerful despite having less complex and slower cores than a CPU.

### Thread Hierarchy and Tuning

The thousands of cores on a GPU, also called threads, are grouped together in several categories. This is the Thread hierarchy of GPUs. The developer can influence this grouping to a degree which allows them to tune their algorithm for optimal performance. In order to develop a well performing algorithm, it is necessary to know how this grouping works. Tuning the grouping is unique to each algorithm and also dependent on the GPU used, which means it is important to test a lot of different configurations to achieve the best possible result. This section aims at exploring the thread hierarchy and how it can be tuned to fit an algorithm.

At the lowest level of a GPU exists a Streaming Multiprocessor (SM), which is a hardware unit responsible for scheduling and executing threads and also contains the registers used by these threads. An SM is always executing a group of 32 threads simultaneously, and this group is called a warp. The number of threads that can be started is virtually unlimited. However, threads must be grouped in a block, with one block typically containing a maximum of 2048 threads but is often configured to be less. Therefore, if more than 2048 threads are required, more blocks must be created. Blocks can also be grouped thread block clusters which is optional, but can be useful in certain scenarios. All thread blocks or thread block clusters are part of a grid, which manifests as a dispatch of the code run the GPU, also called kernel (AMD, 2025b). All threads in one block have access to some shared memory, which can be used for L1 caching or communication between threads. It is important that the blocks can be scheduled independently, with no dependencies between them. This allows the scheduler to schedule blocks and threads as efficiently as possible. All threads within a warp are guaranteed to be part of the same block, and are therefore executed simultaneously and can access the same memory addresses. Figure 2.2 depicts how threads in a block are grouped into warps for execution and how they share memory.

A piece of code that is executed on a GPU is written as a kernel which can be configured. The most important configuration is how threads are grouped into blocks. The GPU allows the kernel to allocate threads and blocks and block clusters in up to three dimensions. This is often useful because of the already mentioned shared memory, which will be explained in more detail in section 2.2.1. Considering the case where an image needs to be blurred, it not only simplifies the development if threads are arranged in a 2D grid, it also helps with optimising memory access. As the threads in a block, need to access a lot of the same data, this data can be loaded in the shared memory of the block. This allows the data to be accessed much quicker compared to when threads are allocated in only one dimension. With one dimensional blocks it is possible that threads assigned to nearby pixels, are part of a different block, leading to a lot of duplicate data transfer. Although the size in each dimension of the blocks can be almost arbitrary, blocks that are too large might lead to other problems which are described in more

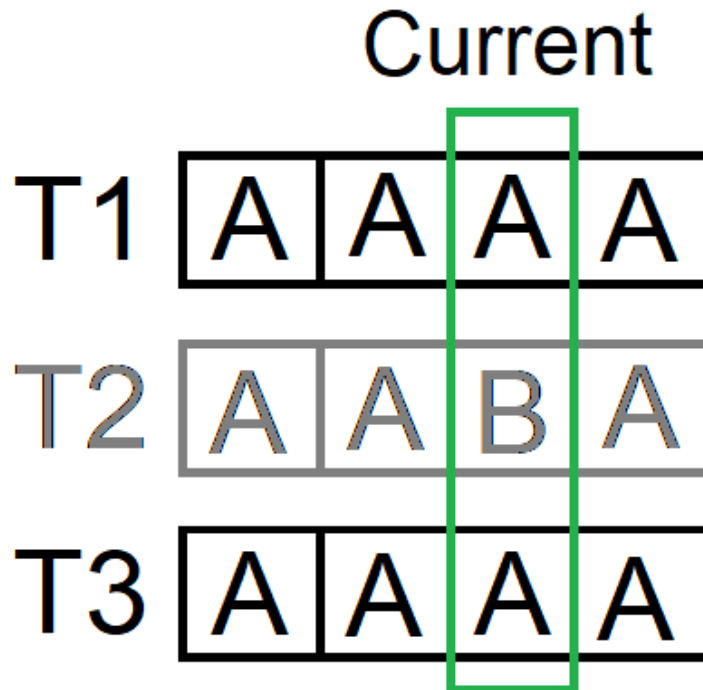


**Figure 2.2:** An overview of the thread hierarchy with blocks being split into multiple warps and their shared memory (AMD, 2025b).

detail in section 2.2.1.

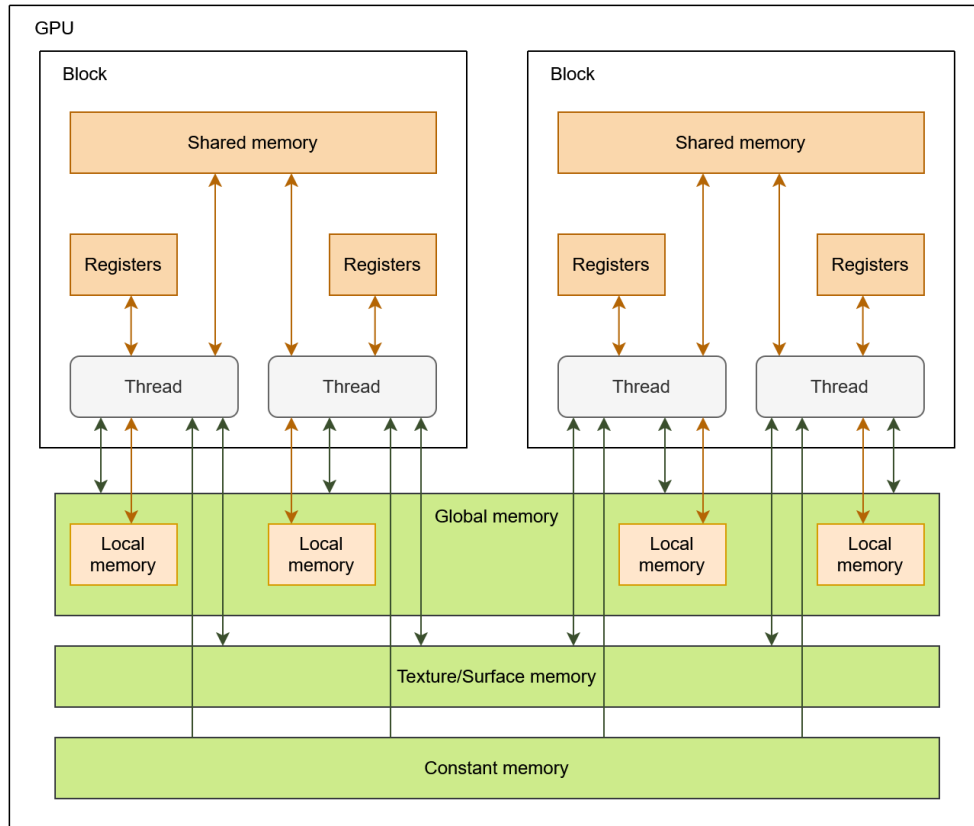
All threads in a warp start at the same point in a program, they have their own instruction address, allowing them to work independently. Because of the SIMD architecture, all threads in a warp must execute the same instructions and if threads start diverging, the SM must pause threads with different instructions and execute them later. Figure 2.3 shows how such divergences can impact performance. The situation described by the figure also shows, that after the divergent thread would reconverge, this does not happen and leads to T2 being executed after T1 and T3 are finished. In situations where a lot of data dependent thread divergence happens, most of the benefits of using a GPU have vanished.

Threads not executing the same instruction is against the SIMD principle but can happen in reality, due to data dependent branching. Consequently, this leads to bad resource utilisation, which in turn leads to worse performance. Another possibility of threads being paused (inactive threads) is the fact that sometimes, the number of threads started is not divisible by 32. In such cases, the last warp still contains 32 threads but only the threads with work are executed.



**Figure 2.3:** Thread T2 wants to execute instruction B while T1 and T3 want to execute instruction A. Therefore T2 will be an inactive thread this cycle and active once T1 and T3 are finished. This means that now the divergent threads are serialised.

Modern GPUs implement the so called Single-Instruction Multiple-Thread (SIMT) architecture. In many cases a developer does not need to know the details of SIMT and can develop fast and correct programs with just the SIMD architecture in mind. However, leveraging the power of SIMT can yield substantial performance gains by re-converging threads once data dependent divergence occurred. A proposal for a re-convergence algorithm was proposed by Collange (2011) where they have shown that these approaches help with hardware occupation, resulting in improved performance as threads are now no longer fully serialised. Another approach for increasing occupancy using the SIMT architecture is proposed by Fung and Aamodt (2011). They introduced a technique for compacting thread blocks by moving divergent threads to new warps until they reconverge. This approach resulted in a noticeable speed-up between 17% and 22%. Another example where a SIMT aware algorithm can perform better was proposed by Köster et al. (2020a). While they did not implement techniques for thread re-convergence, they implemented a thread compaction algorithm. On data-dependent divergence it is possible for threads to end early, leaving a warp with only partial active threads. This means the deactivated threads are still occupied and cannot be used for other work. Their thread compaction tackles this problem by moving active threads into a new thread block, releasing the inactive threads to perform other work. With this they were able to gain a speed-up of roughly 4 times compared to previous implementations.



**Figure 2.4:** The layout of the memory in the GPU. The connections between the memory regions can be seen as well as the different kinds of memory available.

### Memory Model

On a GPU there are two parts that contribute to the performance of an algorithm. The one already looked at is the compute-portion of the GPU. This is necessary because if threads are serialised or run inefficiently, there is nothing that can make the algorithm execute faster. However, algorithms run on a GPU usually require huge amounts of data to be processed, as they are designed for exactly that purpose. The purpose of this section is to explain how the memory model of the GPU works and how it can influence the performance of an algorithm. In figure 2.4 the memory layout and the kinds of memory available are depicted. The different parts will be explained in this section.

On a GPU there are multiple levels and kinds of memory available. All these levels and kinds have different purposes they are optimised for. This means that it is important to know what they are and how they can be best used for specific tasks. On the lowest level threads have registers and local memory available. Registers is the fastest way to access memory but is also the least abundant memory with up to a maximum of 255 32-Bit registers per thread on Nvidia and 256 on AMD (AMD, 2025a). However, using all registers of a thread can lead to other problems which are described in more detail in section 2.2.1. On the other side, the thread local memory is significantly slower

than registers. This is due to the fact, that local memory is actually stored in global memory and therefore has the same limitations which are explained later. This means it is important to try and avoid local memory as much as possible. Local memory is usually only used when a thread uses too many registers. The compiler will then spill the remaining data into local memory and loads it into registers once needed, drastically slowing down the application.

Shared memory is the next tier of memory on a GPU. Unlike local memory and registers, shared memory is shared between all threads inside a block. The amount of shared memory is depending on the GPU architecture but for Nvidia it hovers at around 100 Kilobyte (KB) per block. While this memory is slower than registers, its primary use-case is communicating and sharing data between threads in a block. It is advised that all threads in a block access a lot of overlapping data, as then data from global memory can be loaded into faster shared memory once and then accessed multiple times further increasing performance. Loading data into shared memory and accessing that data has to be done manually. Because shared memory is part of the unified data cache, it can either be used as a cache or for manual use, meaning a developer can allocate more shared memory towards caching if needed. Another feature of shared memory are the so-called memory banks. Shared memory is always split into 32 equally sized memory modules also called memory banks. All available memory addresses lie in one of these banks. This means if two threads access two different memory addresses which lie in different banks, the access can be performed simultaneously, increasing the throughput.

The most abundant and slowest memory is the global memory and resides in device memory. A key constraint of device memory and therefore global memory is, that is accessed in either 32, 64 or 128 byte chunks. This means if a thread wants to access 8 bytes from global memory, alongside the 8 bytes, the 24 bytes after the requested 8 bytes are also transferred. As a result, the throughput is only a fourth of the theoretical maximum. Therefore, it is important to follow optimal access patterns. What these optimal patterns are, are architecture dependent and are described in the according sections in the CUDA programming guide.

A small portion of device memory is allocated to constant memory. Constant memory is accessible by all threads and as the name implies, can not be written to by threads. It can be initialised by the CPU when starting a kernel if needed. As constant memory has a separate cache, it can be used to speed-up data access for constant and frequently accessed data.

Another special kind of memory is the texture and surface memory. According to AMD (2025b) texture memory is read-only memory, while surface memory can also be written to, which is the only difference between these two kinds of memory. Nvidia does not explicitly state this behaviour, but due to the fact that accessing textures is only performed via caches, it is implied that on Nvidia GPUs, texture memory is also read-only. As the name implies, this kind of memory is optimised for accessing textures. This means that threads of the same warp, accessing data which is spatially close together, will result in increased performance. As already mentioned, surface memory works the same way, with the difference, that it can be written to. It is therefore well suited for manipulating two- or three-dimensional data.

Compute Capability	8.9	10.x
Max. number of threads per block	1 024	
Warp size	32 threads	
Max. number of warps per SM	48	64
Max. number of blocks per SM	24	32
Max. number of threads per SM	1 536	2 048
Number of 32-bit registers per SM	64 000	
Max. number of 32-bit registers per block	64 000	
Max. number of 32-bit registers per thread	255	
Max. amount of shared memory per SM	100 Kilobytes	228 Kilobytes
Max. amount of shared memory per block	99 Kilobytes	227 Kilobytes

**Table 2.1:** A simplified version of the technical specifications for the Compute Capabilities 8.9 and 10.x (Nvidia, 2025b). These correspond to the Nvidia Ada Lovelace and Blackwell microarchitectures.

## Occupancy

Occupancy describes the utilisation of a GPU. A high occupancy means, that there are Warps executing, or in other words, the cores are occupied with work. This is important, as a low occupancy means that the GPU is waiting for work to be scheduled and is therefore idle. As a result, it is important to achieve high occupancy in order to increase the performance of an algorithm. It needs to be noted, that occupancy is not the only option for improving performance. As it is possible for the GPU to have a high occupancy while performing a lot of unnecessary or redundant work or utilising compute-resources that are slower. An example for the latter would be developing an algorithm that uses 64-bit floating point (FP64) numbers while 32-bit floating point (FP32) numbers would have sufficient accuracy. Because GPUs tend to have fewer FP64 compute-resources than they have FP32 compute-resources, performing FP64 operations will take longer. However, despite these drawbacks, having low occupancy will very likely result in performance degradation while high occupancy will either improve performance or do no harm otherwise. Ways of achieving high occupancy will be outlined in this section as most other performance problems can be solved algorithmically.

When starting a kernel, the most important configuration is the number of threads and thread blocks that need to be started. This is important, as this has other effects on occupancy as well. In table 2.1 the most notable limitations are presented that can affect occupancy. These limitations need to be considered when choosing a kernel configuration. It is important to note, that depending on the GPU and problem, the occupancy tuning might differ, and the same approach might perform well on one GPU but perform poorly on another GPU. Therefore, the things discussed here are only guidelines and tools like Nvidia Nsight Compute<sup>7</sup> and Nsight Systems<sup>8</sup> are essential for performance tuning. Nsight compute also contains an occupancy calculator which takes a kernel and computes how the configuration performs in terms of occupancy and also

<sup>7</sup><https://developer.nvidia.com/nsight-compute>

<sup>8</sup><https://developer.nvidia.com/nsight-systems>

lets the developer try out different configurations (Nvidia, 2025c).

In general, it is important to have as many warps as possible ready for execution. While this means that a lot of warps could be executed but are not, this is actually desired. A key feature of GPUs is so-called latency hiding, meaning that while a warp waits for data to be retrieved for example, another warp ready for execution can now be run. With low occupancy, and therefore little to no warps waiting for executing, latency hiding does not work, as now the hardware is idle. As a result, the runtime increases which also explains why high occupancy is not guaranteed to result in performance improvements while low occupancy can and often will increase the runtime.

As seen in table 2.1, there exist different limitations that can impact occupancy. The number of warps per SM is important, as this means this is the degree of parallelism achievable per SM. If due to other limitations, the number of warps per SM is below the maximum, there is idle hardware. One such limitation is the number of registers per block and SM. In the case of compute capability 8.9, one SM can handle  $32 * 48 = 1\,536$  threads. This leaves  $64\,000/1\,536 \approx 41$  registers per thread, which is lower than the theoretical maximum of 255 registers per thread. Typically, one register is mapped to one variable in the kernel code, meaning a developer can use up to 41 variables in their code. However, if the variable needs 64 bits, the register usage doubles, as all registers on a GPU are 32-bit. On a GPU with compute capability 10.x a developer can use up to  $64\,000/2\,048 \approx 31$  registers. Of course a developer can use more registers, but this results in less occupancy. However, depending on the algorithm using more registers might be more beneficial to performance than the lower occupancy, in which case occupancy is not as important. If a developer needs more than 255 registers for their variables the additional variables will spill into local memory which is, as described in section 2.2.1, not desirable.

Additionally, shared memory consumption can also impact the occupancy. If for example a block needs all the available shared memory, which is almost the same as the amount of shared memory per SM, this SM can only serve this block. On compute capability 10.x, this would mean that occupancy would be at maximum 50% as a block can have up to 1 024 threads and an SM supports up to 2 048 threads. Again, in such cases it needs to be determined, if the performance gain of using this much shared memory is worth the lower occupancy.

Balancing these limitations and therefore the occupancy and performance often requires a lot of trial and error with help of the aforementioned tools. In cases where occupancy is already high and the amount of warps ready for execution is also high, other areas for performance improvements need to be explored. Algorithmic optimisation is always a good idea. Some performance improvements can be achieved by altering the computations to use different parts of the GPU. One of such optimisations is using FP32 operations wherever possible. Another well suited optimisation is to rewrite the algorithm to use as many Fused Multiply-Add (FMA) instructions. FMA is a special floating point instruction, that multiplies two values and adds a third, all in a single clock cycle (Nvidia, 2025a). However, the result might slightly deviate from performing these two operations separately, which means in accuracy sensitive scenarios, this instruction should be avoided. If the compiler detects a floating point operation with the FMA structure, it will automatically be compiled to an FMA instruction. To prevent this, in C++ the developer can call the functions `__fadd__` and `__fmul__` for addition

and multiplication respectively.

### 2.2.2 Parallel Thread Execution

While in most cases a GPU in a higher level language like C++ or even Julia<sup>9</sup>, it is also possible to program GPUs with the low level language Parallel Thread Execution (PTX) developed by Nvidia. A brief overview of what PTX is and how it can be used to program GPUs is given in this section. Information in this section is taken from the PTX documentation (Nvidia, 2025d) if not stated otherwise.

PTX defines a virtual machine with an own instruction set architecture (ISA) and is designed for data-parallel processing on a GPU. It is an abstraction of the underlying hardware instruction set, allowing PTX code to be portable across Nvidia GPUs. In order for PTX code to be usable for the GPU, the compiler is responsible for compiling the code to the hardware instruction set of the GPU it is run on. A developer typically writes a kernel in CUDA using C++, for example, and the Nvidia compiler generates the PTX code for that kernel. The concepts for programming the GPU with PTX and CUDA are the same, apart from the terminology which is slightly different. For consistency, the CUDA terminology will continue to be used.

Syntactically PTX resembles Assembly style code. Every PTX code must have a `.version` directive which indicates the PTX version and an optional `.target` directive which indicates the compute capability. If the program works in 64 bit addresses, the optional `.address_size` directive can be used to indicate that, which simplifies the code for such applications. After these directives, the actual code is written. As each PTX code needs an entry point (the kernel) the `.entry` directive indicates the name of the kernel and the parameters needed. It is also possible to write helper functions with the `.func` directive. Inside the kernel or a helper function, normal PTX code can be written. Because PTX is very low level, it assumes an underlying register machine, therefore a developer needs to think about register management. This includes loading data from global or shared memory into registers if needed. Code for manipulating data like addition and subtraction generally follow the structure `operation.datatype` followed by three parameters for that operation. For adding two FP32 values together and storing them in the register `%n`, the code looks like the following:

```
add.f32    \%n, 0.1, 0.2;
```

Loops in the classical sense do not exist in PTX. Alternatively a developer needs to define jump targets for the beginning and end of the loop. The code in 2.1 shows how a function with simple loop can be implemented. The loop counts down to zero from the passed parameter  $N$  which is loaded into the register `%n` in line 6. If the value in the register `%n` reached zero the loop branches at line 9 to the jump target at line 12 and the loop has finished. All other used directives and further information on writing PTX code can be taken from the PTX documentation (Nvidia, 2025d).

## 2.3 Compilers

Maybe even move this entire section to “Concept and Design”?

---

<sup>9</sup><https://juliagpu.org/>



```
1 .func loop(.param .u32 N)
2 {
3   .reg .u32 \%n;
4   .reg .pred \%p;
5
6   ld.param.u32 \%n, [N];
7 Loop:
8   setp.eq.u32 \%p, \%n, 0;
9   @\%p bra    Done;
10  sub.u32    \%n, \%n, 1;
11  bra      Loop;
12 Done:
13 }
```

**Program 2.1:** A PTX program fragment depicting how loops can be implemented.

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

### 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

- Bartlett, D. J., Desmond, H., & Ferreira, P. G. (2024). Exhaustive symbolic regression. *IEEE Transactions on Evolutionary Computation*, 28(4), 950–964. <https://doi.org/10.1109/TEVC.2023.3280250> (cit. on pp. 5, 6)
- Besard, T., Churavy, V., Edelman, A., & Sutter, B. D. (2019a). 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. 7)
- Besard, T., Foket, C., & De Sutter, B. (2019b). 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. 7)
- Bomarito, G. F., Leser, P. E., Strauss, N. C. M., Garbrecht, K. M., & Hochhalter, J. D. (2022). Bayesian model selection for reducing bloat and overfitting in genetic programming for symbolic regression. *Proceedings of the Genetic and Evolutionary Computation Conference Companion*, 526–529. <https://doi.org/10.1145/3520304.3528899> (cit. on p. 5)
- Brodtkorb, A. R., Hagen, T. R., & Sætra, M. L. (2013). Graphics processing unit (GPU) programming strategies and trends in GPU computing. *Journal of Parallel and Distributed Computing*, 73(1), 4–13. <https://doi.org/10.1016/j.jpdc.2012.04.003> (cit. on p. 2)
- Brunton, S. L., Proctor, J. L., & Kutz, J. N. (2016). Discovering governing equations from data by sparse identification of nonlinear dynamical systems. *Proceedings of the National Academy of Sciences*, 113(15), 3932–3937. <https://doi.org/10.1073/pnas.1517384113> (cit. on p. 4)
- Collange, C. (2011, September). *Stack-less SIMT reconvergence at low cost* (Research Report). ENS Lyon. <https://hal.science/hal-00622654>. (Cit. on p. 11)
- Dabhi, V. K., & Chaudhary, S. (2012). A survey on techniques of improving generalization ability of genetic programming solutions. <https://doi.org/10.48550/ARXIV.1211.1119> (cit. on p. 6)
- 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. 6)
- Dong, J., Zhong, J., Liu, W.-L., & Zhang, J. (2024). Evolving equation learner for symbolic regression. *IEEE Transactions on Evolutionary Computation*, 1–1. <https://doi.org/10.1109/TEVC.2024.3404650> (cit. on p. 5)

- Franchetti, F., Kral, S., Lorenz, J., & Ueberhuber, C. (2005). Efficient utilization of SIMD extensions. *Proceedings of the IEEE*, *93*(2), 409–425. <https://doi.org/10.1109/JPROC.2004.840491> (cit. on p. 9)
- Fung, W. W. L., & Aamodt, T. M. (2011). Thread block compaction for efficient SIMT control flow. *2011 IEEE 17th International Symposium on High Performance Computer Architecture*, 25–36. <https://doi.org/10.1109/HPCA.2011.5749714> (cit. on p. 11)
- 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, 7)
- Guillemot, H. (2022, December 31). Climate models. In K. De Pryck (Ed.), *A critical assessment of the intergovernmental panel on climate change* (1st ed., pp. 126–136). Cambridge University Press. <http://dx.doi.org/10.1017/9781009082099.018>. (Cit. on p. 4)
- 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, 7)
- Han, T. D., & Abdelrahman, T. S. (2011). hiCUDA: High-level GPGPU programming. *IEEE Transactions on Parallel and Distributed Systems*, *22*(1), 78–90. <https://doi.org/10.1109/TPDS.2010.62> (cit. on p. 6)
- Hissbach, A.-M., Dick, C., & Lawonn, K. (2022). An overview of techniques for egocentric black hole visualization and their suitability for planetarium applications. In J. Bender, M. Botsch, & D. A. Keim (Eds.), *Vision, modeling, and visualization*. The Eurographics Association. <https://doi.org/10.2312/vmv.20221207>. (Cit. on p. 7)
- Huang, Q., Huang, Z., Werstein, P., & Purvis, M. (2008). GPU as a general purpose computing resource. *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. 6)
- 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. 5)
- 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. 5)
- 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](https://doi.org/10.1007/978-1-4614-1770-5_8). (Cit. on p. 5)
- Köster, M., Groß, J., & Krüger, A. (2020a). Massively parallel rule-based interpreter execution on GPUs using thread compaction. *International Journal of Parallel Programming*, *48*(4), 675–691. <https://doi.org/10.1007/s10766-020-00670-2> (cit. on p. 11)
- Köster, M., Groß, J., & Krüger, A. (2020b). 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](https://doi.org/10.1007/978-3-030-60245-1_26). (Cit. on p. 7)



- Koza, J. R. (2010). Human-competitive results produced by genetic programming. *Genetic Programming and Evolvable Machines*, 11(3), 251–284. <https://doi.org/10.1007/s10710-010-9112-3> (cit. on p. 5)
- Koza, J. (1994). Genetic programming as a means for programming computers by natural selection. *Statistics and Computing*, 4(2). <https://doi.org/10.1007/BF00175355> (cit. on pp. 4, 5)
- Kronberger, G., Burlacu, B., Kommenda, M., Winkler, S. M., & Affenzeller, M. (2024, July). *Symbolic regression*. Chapman; Hall/CRC. <http://dx.doi.org/10.1201/9781315166407>. (Cit. on pp. 5, 6)
- Lee, V. W., Kim, C., Chhugani, J., Deisher, M., Kim, D., Nguyen, A. D., Satish, N., Smelyanskiy, M., Chennupati, S., Hammarlund, P., Singhal, R., & Dubey, P. (2010). Debunking the 100x GPU vs. CPU myth: An evaluation of throughput computing on CPU and GPU. *Proceedings of the 37th annual international symposium on Computer architecture*, 451–460. <https://doi.org/10.1145/1815961.1816021> (cit. on p. 7)
- Lemos, P., Jeffrey, N., Cranmer, M., Ho, S., & Battaglia, P. (2022, February 4). Rediscovering orbital mechanics with machine learning. <https://doi.org/10.48550/arXiv.2202.02306>. (Cit. on p. 5)
- 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. 7)
- Martius, G., & Lampert, C. H. (2016). Extrapolation and learning equations [Version Number: 1]. <https://doi.org/10.48550/ARXIV.1610.02995>. (Cit. on p. 5)
- Michalakes, J., & Vachharajani, M. (2008). GPU acceleration of numerical weather prediction. *2008 IEEE International Symposium on Parallel and Distributed Processing*, 1–7. <https://doi.org/10.1109/IPDPS.2008.4536351> (cit. on pp. 2, 7)
- 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)
- Sahoo, S. S., Lampert, C. H., & Martius, G. (2018). Learning equations for extrapolation and control. <https://doi.org/10.48550/ARXIV.1806.07259>. (Cit. on p. 5)
- Schuurman, D. C. (2013). Step-by-step design and simulation of a simple CPU architecture. *Proceeding of the 44th ACM technical symposium on Computer science education*, 335–340. <https://doi.org/10.1145/2445196.2445296> (cit. on p. 8)
- Sun, F., Liu, Y., Wang, J.-X., & Sun, H. (2023, February 2). Symbolic physics learner: Discovering governing equations via monte carlo tree search. <https://doi.org/10.48550/arXiv.2205.13134>. (Cit. on p. 5)
- Tian, X., Saito, H., Girkar, M., Preis, S. V., Kozhukhov, S. S., Cherkasov, A. G., Nelson, C., Panchenko, N., & Geva, R. (2012). Compiling c/c++ SIMD extensions for function and loop vectorizaion on multicore-SIMD processors. *2012 IEEE 26th International Parallel and Distributed Processing Symposium Workshops & PhD Forum*, 2349–2358. <https://doi.org/10.1109/IPDPSW.2012.292> (cit. on p. 9)

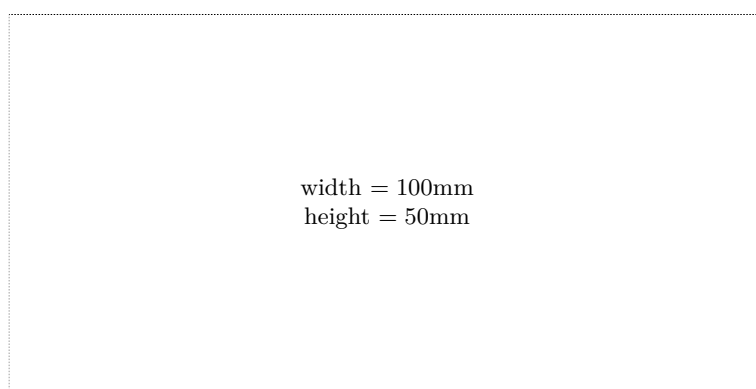
- Verbraeck, A., & Eisemann, E. (2021). Interactive black-hole visualization. *IEEE Transactions on Visualization and Computer Graphics*, 27(2), 796–805. <https://doi.org/10.1109/TVCG.2020.3030452> (cit. on p. 7)
- Werner, M., Junginger, A., Hennig, P., & Martius, G. (2021, May 13). Informed equation learning. <https://doi.org/10.48550/arXiv.2105.06331>. (Cit. on p. 5)

## Online sources

- AMD. (2025a, February). *Hardware features — HIP 6.3.42134 documentation*. Retrieved March 15, 2025, from [https://rocm.docs.amd.com/projects/HIP/en/latest/reference/hardware\\_features.html](https://rocm.docs.amd.com/projects/HIP/en/latest/reference/hardware_features.html). (Cit. on p. 12)
- AMD. (2025b, February). *HIP programming model — HIP 6.3.42134 documentation*. Retrieved March 9, 2025, from [https://rocm.docs.amd.com/projects/HIP/en/latest/understand/programming\\_model.html](https://rocm.docs.amd.com/projects/HIP/en/latest/understand/programming_model.html). (Cit. on pp. 9, 10, 13)
- Nvidia. (2025a, March). *CUDA c++ best practices guide 12.8 documentation*. Retrieved March 16, 2025, from <https://docs.nvidia.com/cuda/cuda-c-best-practices-guide/index.html>. (Cit. on p. 15)
- Nvidia. (2025b, March). *CUDA c++ programming guide*. Retrieved November 22, 2024, from <https://docs.nvidia.com/cuda/cuda-c-programming-guide/>. (Cit. on pp. 7, 8, 14)
- Nvidia. (2025c, March). *Nsight compute — NsightCompute 12.8 documentation*. Retrieved March 16, 2025, from <https://docs.nvidia.com/nsight-compute/NsightCompute/index.html#occupancy-calculator>. (Cit. on p. 15)
- Nvidia. (2025d, March). *Parallel thread execution ISA version 8.7*. Retrieved March 15, 2025, from <https://docs.nvidia.com/cuda/parallel-thread-execution/>. (Cit. on p. 16)
- Sutter, H. (2004, December). *The free lunch is over: A fundamental turn toward concurrency in software*. Retrieved March 13, 2025, from <http://www.gotw.ca/publications/concurrency-ddj.htm>. (Cit. on p. 1)

# Check Final Print Size

— Check final print size! —



— Remove this page after printing! —