

Overview and Analysis of GPU Acceleration for Regular Expressions *

Roman Gajdoš

Slovak University of Technology in Bratislava
Faculty of Informatics and Information Technologies
`xgajdosr@stuba.sk`

25. september 2023

Abstract

Pattern matching is essential across numerous domains and regular expressions are its fundamental component. They are used in various applications, including database queries, network security, bioinformatics and more. It is crucial to ensure efficient use of memory and speed, particularly when handling vast datasets or time-sensitive tasks. Conventional processors face limitations in their ability to execute parallel tasks. However, Graphics Processing Units (GPUs) have become prevalent in many computing systems and offer a promising solution with extensive parallelism, computational capacity, and high memory bandwidth.

This paper provides a comprehensive overview of GPU-accelerated regular expressions, including comparative analysis and reference to previous studies. Our investigation into the acceleration of regular expressions illustrates the potential of the GPU as a practical platform for pre-processing, which outperforms other platforms in terms of speed, but sometimes operates with a lower throughput. Moreover, we demonstrate variations in the performance of finite automata representations and their effect, explaining that the NFA representation is more efficient in terms of memory, whereas the DFA representation is more rapid. We also note the need for further research in this area. There has been minimal recent exploration, specifically regarding the potential combination of different accelerator platforms with GPUs for regular expression acceleration or examination of OpenCL in comparison to CUDA is yet to be undertaken, and would benefit from investigation.

1 Introduction

Pattern matching is widely used in a variety of different domains. Regular expressions have become a prevalent tool for text processing and sanitation due to their flexibility, conciseness, and vast support in most programming languages [5]. They appear in approximately a third of open-source projects [6]. They are employed in technical fields, ranging from database querying [11], texts

*Semestrálny projekt v predmete Metódy inžinierskej práce, ak. rok 2023/24, vedenie: MSc. Mirwais Ahmadzai

editors¹, web scraping [9] to network security, such as deep packet inspection [3], and bioinformatics [10], among others.

Regular expressions are implemented using finite automata, in either deterministic (DFA) or non-deterministic (NFA) form, each with their respective advantages and drawbacks. Each of them has their own advantages and disadvantages [17, 27, 29].

In many applications, regular expressions are applied to large amounts of input data, or require a fast response, or both. It stands to reason that efficiency in both memory and speed is the key to optimal use [24]. Here, the question is how to achieve the greatest possible efficiency for a given problem that is addressed by the regular expressions.

The processor’s capacity to execute multiple expressions simultaneously is notably restricted, even in the current era of multicore processors [13]. However, its frequency and cache memory speed prove excellent for handling small datasets. For tasks that require more extensive parallelism, FPGAs (Field-Programmable Gate Arrays) or ASICs (Application-Specific Integrated Circuits) have been used. The problem is that they are slow to configure [25] and inflexible to change [8, 16].

In recent years, GPUs with their extensive parallelism, computational capabilities, and high memory bandwidth have become prevalent in numerous computing system. They have scaled at a faster rate than CPUs, providing significant computing power [16, 20]. APIs were created to allow General Purpose Graphics Processing Units (GPGPU) to accelerate processing in supported applications, replacing shading languages and simplifying their use for programmers. Two popular APIs are Compute Unified Device Architecture (CUDA) and Open Computing Language (OpenCL) [7].

In this paper, we investigate a variety of GPU-based regular expression execution methods and conduct a comparative analysis of their strengths and weaknesses. Our research begins with a thorough examination of regular expressions in 3. This is followed by a comparison of their representations of finite state automata forms in 3.1. We then move into parallel computing platforms, such as CUDA and OpenCL in 3.2.1.

By combining these findings, our investigation aims to provide a comprehensive overview of GPU-accelerated regular expressions, utilizing previous studies to provide an in-depth comparative analysis in section 4 and 5. Result from this was evaluated in the 6 section. Finally, we conclude with a summary of our findings in section 8.

2 Objective and methodology

At present, we are unaware of any comprehensive analysis of the available materials on the acceleration of regular expressions on GPU. We have therefore decided to review the available data to see if there are any conflicting statements or claims that have evolved over time or are no longer valid. The outcome ought to be an in-depth, up-to-date article showing when to use GPU regular expression acceleration, what its benefits are, and when other solutions are better.

This research compares data from papers on or relevant to the topic and analyses the evolution of key metrics, methodologies, and performance bench-

¹<https://neovim.io/doc/user/change.html#%3Asubstitute>

marks reported in these papers. By considering a variety of sources, we aim to capture the trajectory of GPU regular expression acceleration research over time. Our paper selection criteria include relevance and research quality. The topics of the papers should cover GPU acceleration of regular expressions, or GPU acceleration of finite state machines, and other related topics.

Furthermore, we will address any conflicting statements or evolving claims within the literature, providing readers with a clear understanding of the current state of knowledge in this domain. The objective is to present an informative and well-rounded assessment of GPU acceleration for regular expressions, shedding light on when, why, and how this approach proves advantageous, as well as delineating scenarios where alternative solutions may be more suitable.

We will also look at the various methods of accelerating regular expressions on GPUs and compare their strengths and weaknesses. We will maintain an unbiased perspective and present information in a clear and concise manner. Technical terminology will be firstly explained in separate section.

3 Background

A regular expression, or regex for short, represents a set of exactly matching strings of characters and special symbols. This set can be infinite. The string of characters is then matched against the pattern to see if it matches. Regular expressions can be constructed in several ways [22]. The most common is to use a formal language, such as the one in POSIX standard². Basic syntax is described as follows: characters of alphabet are matched literally, special symbols are used to match a single/multiple character matches, optional character matches, alternation, any character, line start/end and the empty string. As described in table 1.

Symbols	Meaning
.	Any character
*	Zero or more matches
+	One or more matches
?	Zero or one match
	Alternation
-	Range
\	Escape character
^	Line start
\$	Line end
:	Grouping
[]	Start and end of character class
()	Start and end of group
{ }	Start and end of quantifier

Table 1: Regular expression special symbols, author’s own work

²https://pubs.opengroup.org/onlinepubs/9699919799/basedefs/V1_chap09.html

3.1 Finite Automata

Regular expression matching is performed by using finite automata, a mathematical model of computation that abstracts computations into a finite number of states and transitions [27]. A finite automaton comprises a directed graph in which each node symbolises a state while each edge reflects a state transition. Two widely used representations of finite automata are deterministic finite automata (DFAs) and non-deterministic finite automata (NFAs). While NFAs and DFAs achieve the same outcome, there are some practical differences in terms of resource requirements and traversal behaviour [17]. Although deterministic finite automata (DFAs) have a simpler transition system, their execution is serial and the size of transitions in DFAs may be significantly larger than their equivalent NFAs [15, 16].

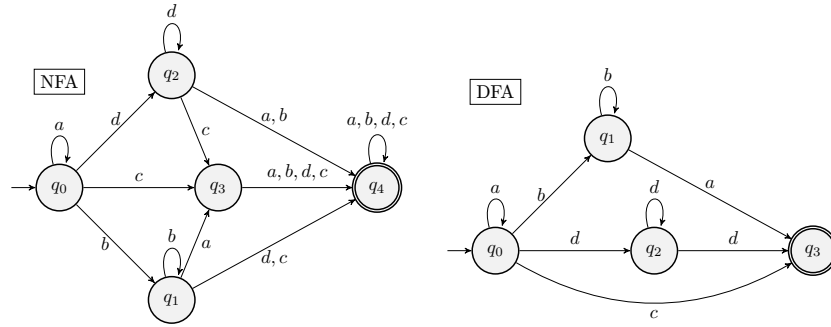


Figure 1: NFA and DFA for regular expression $a^*(b+a|d^*c)$, autor's own work

3.1.1 Automata processing

The regular expression's matching process is equivalent to a finite state machine traversal of the input stream. The process of matching begins with the activation of the initial states. Symbols from the input stream are sequentially utilized by the finite automaton. The process concludes once all the symbols of the input stream are processed. The incoming symbol is matched with the active states, and if it falls within the matchset of an active state, the active state transforms into a matched state [16].

We can guarantee worst-case performance by restricting the processing of each input character. Techniques to limit per-character processing involve enlarging the finite automaton. Therefore, the search space is determined by balancing the size of the automaton and the upper limit of per-character processing [17].

The benefit of using NFA is its capability to construct an NFA consisting of states that are less than or equal to the number of characters in the pattern-set. This makes the representation compact. The primary deficiency of NFAs lies in their traversal, during which the number of active states can vary from iteration to iteration, as well as the amount of work performed [27].

3.2 Graphics Processing Unit

The GPU, or Graphics Processing Unit, offers significantly higher instruction throughput and memory bandwidth than the CPU at a comparable price and power consumption. While the CPU is optimized for rapid execution of a sequence of operations referred to as a "thread" and can execute dozens of these threads concurrently, the GPU is optimized for thousands of parallel executions (offsetting the slower single-thread performance for increased throughput). This variation in capabilities is a result of differing design objectives for the GPU and CPU.³

All threads within a compute unit share a common instruction counter. Execution of a single compute unit occurs in lock-step, whereby each thread executes the same instruction when directed to do so. When control flow divergence occurs between threads within the same work group, divergent instructions are serialised, which can negatively impact performance. Similarly, an unbalanced workload across threads in a work-group will result in idle threads, which reduces performance [26].

The memory hierarchy of GPUs comprises: Global memory, a larger and slower form of memory accessible by all threads in any compute units. Constant memory, a read-only section of the global memory with a specific cache for faster memory access. Local memory, connected to compute units and shared among threads in a single unit and private memory, exclusively reserved for individual threads [26].

3.2.1 Parallel computing platforms

Due to the significant performance potential of GPUs, their utilization has evolved from high-level shading languages to modern programming languages, reducing time and complexity involved in creating GPU-enabled applications [1]. Two of the most popular APIs for GPU programming are CUDA⁴ (Compute Unified Device Architecture) and OpenCL⁵ (Open Computing Language). CUDA is a proprietary API developed by NVIDIA, while OpenCL is an open royalty-free standard developed by the Khronos Group (an open, member-driven consortium, publishing and maintaining open standards)⁶.

OpenCL provides an efficient and portable way to access the power of different computing platforms. However, when comparing OpenCL to CUDA, there are sometimes performance differences. These differences are often attributed to the portability of OpenCL, which can lead to performance degradation. However, in a fair comparison, there is no inherent reason for OpenCL to perform worse than CUDA. Performance differences are primarily due to programmers and compilers behaviour [7].

4 Related Work

Research into the use of GPUs for regular expression matching started with the publications *Fast Exact String Matching on the GPU* [19] and *A GPU-*

³From <https://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html>

⁴<https://developer.nvidia.com/cuda-toolkit>

⁵<https://www.khronos.org/opencl/>

⁶<https://www.khronos.org/>

based *Multiple-pattern Matching Algorithm for Network Intrusion Detection Systems* [10]. Their pioneering studies introduced a string matching program and multiple-pattern matching algorithm designed to take advantage of GPU processing capabilities. This work was soon followed by *Accelerating Regular Expression Matching Using Hierarchical Parallel Machines on GPU* [14] and *GPU-based NFA Implementation for Memory Efficient High Speed Regular Expression Matching* [29]. These two studies provide solutions to the performance challenges linked to regular expression matching in the context of network intrusion detection and other network functions.

Afterwards, a large, comprehensive study entitled *GPU Acceleration of Regular Expression Matching for Large Datasets: Exploring the Implementation Space* [27] was published. The study examines regular expression matching on GPUs, focusing on practical-sized and complex datasets. It explores the advantages and limitations of different automata representations and various GPU implementation techniques.

In the paper *Demystifying Automata Processing: GPUs, FPGAs or Micron's AP?* [17], the authors compare the performance of GPUs, FPGAs, and Micron's Automata Processor (AP) for regular expression matching.

Later, studies were published dealing with speeding up the processing of finite state machines. Papers *Scaling Out Speculative Execution of Finite-State Machines with Parallel Merge* [24] based on *On-the-Fly Principled Speculation for FSM Parallelization* [28] introduce a speculative execution technique for finite state machines. The latest study *GSpecPal: Speculation-Centric Finite State Machine Parallelization on GPUs* [26] addressed previous problems with DFA parallelization on GPUs.

Why GPUs are Slow at Executing NFAs and How to Make them Faster [15] proposed and evaluated optimisations to improve the throughput of NFA processing on GPUs by addressing tackle suboptimal data movement and under-utilisation. *Asynchronous Automata Processing on GPUs* [16] have developed a lightweight approach to increase the parallelism of automata processing on GPUs by asynchronously searching for patterns in the input stream in parallel.

5 Analysis

5.1 GPU Acceleration of Regular Expression Matching for Large Datasets: Exploring the Implementation Space

This study [27], published in 2013, was the first to make relevant measurements using real-life data. The measurements were conducted using NFA design by [4], along with their optimised NFA; and DFA design by [2], in uncompressed (U-DFA), compressed (C-DFA), and optimized (E-DFA) forms.

Their dataset used both real and synthetic pattern sets. The real patterns were taken from Snort's backdoor and spyware rules. The synthetic sets were generated using a tool described in [3] and tokens from the backdoor rules.

Measurements were made on system consisting of Xeon E5620 CPU and an NVIDIA GTX 480 GPU, using CentOS 5.5 and CUDA 4.

5.1.1 Study findings

The repetition of wildcards and large character sets can lead to a state explosion when converting an NFA to a DFA. On the chart 2, we can see size of DFAs in MB. In contrast, NFAs are consistently less than 1 MB in size. For synthetic datasets, the uncompressed DFA representation has 40-50 times the memory requirement of its compressed counterpart. At the time of writing, the Dotstar.2 datasets exceed the device memory capacity when uncompressed-DFA representation is utilized.

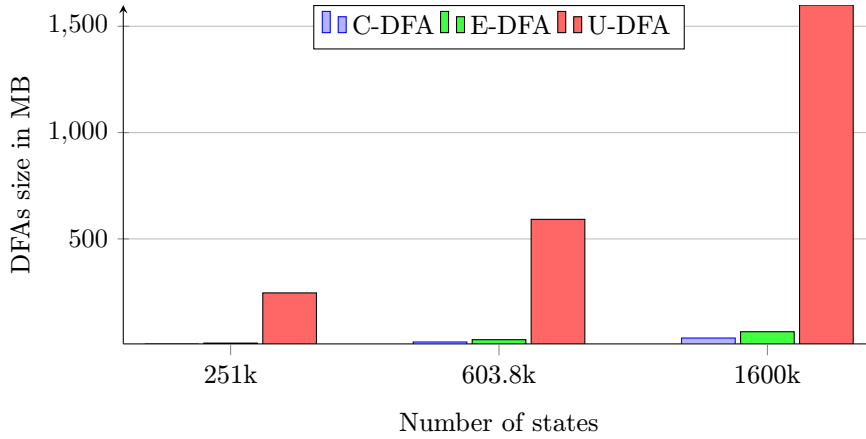


Figure 2: Scaling of DFAs size, author's own work, data from [27]

As shown in chart 3, the U-DFA is the optimal solution for all pattern-sets, when applicable. However, this solution has high memory requirements and cannot be applied to complex pattern-sets due to a lack of graphics card memory. E-DFA offers a reasonable trade-off between U-DFA and C-DFA, resulting in a 3-5 times enhancement in speed compared to C-DFA at the expense of almost 1.5X more memory usage. Moreover, E-DFA performs better than NFA solutions on nearly all datasets.

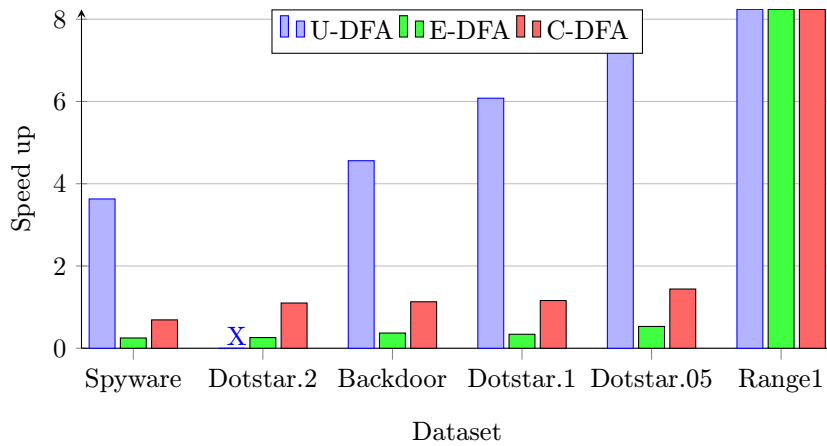


Figure 3: Speedup of DFA's over NFA, author's own work, data from [27]

As shown in chart 4, both GPU implementations based on NFA and DFA algorithms outperformed their CPU equivalent.

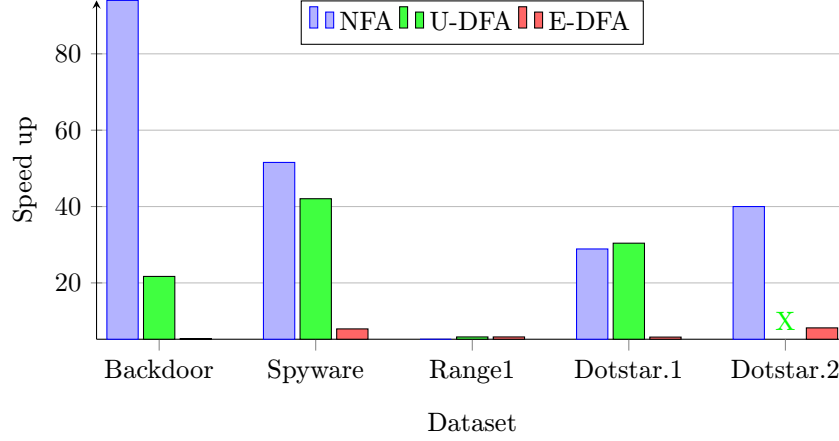


Figure 4: Speed up of GPU accelerated traversals over CPU, author’s own work, data from [27]

5.2 Demystifying Automata Processing: GPUs, FPGAs or Micron’s AP?

This 2017 study [17] investigated the acceleration of Finite State Automata on GPU, FPGA and Micron’s AP. The automata model employed was NFA, and the algorithm used for its traversal engine was described in [21].

Their dataset consists of small NIDS, bioinformatics and synthetic pattern sets. The measurements were performed on a system that consisted of a dual 6-core Intel Xeon CPU and an Nvidia Titan X GPU (Maxwell architecture), using CentOS 6.4 and CUDA 7. The FPGA Xilinx XC6VLX130T was chosen as it fell in the same price range as the GPU (approximately \$1,200 at that time).

5.2.1 Study findings

The NFA format for GPU has a very small global memory utilisation, and even the largest dataset only occupies up to 133MB of global memory, which means we are able to fit any size dataset into a single GPU device, unlike FPGA and Micron’s AP, which required more devices, as shown in the graph 5. However, the GPU records the lowest throughput data, with a 20 times lower throughput on smaller datasets up to 1000 times lower on the largest datasets.

The optimization and partitioning steps for the NFA, which are common for all platforms, take between 3 and 249 seconds, varying with the size of the dataset. The GPU preprocessing task is primarily concerned with parsing the NFA partition files. As clearly demonstrated in Figure 6, GPUs provide considerably shorter preprocessing times.

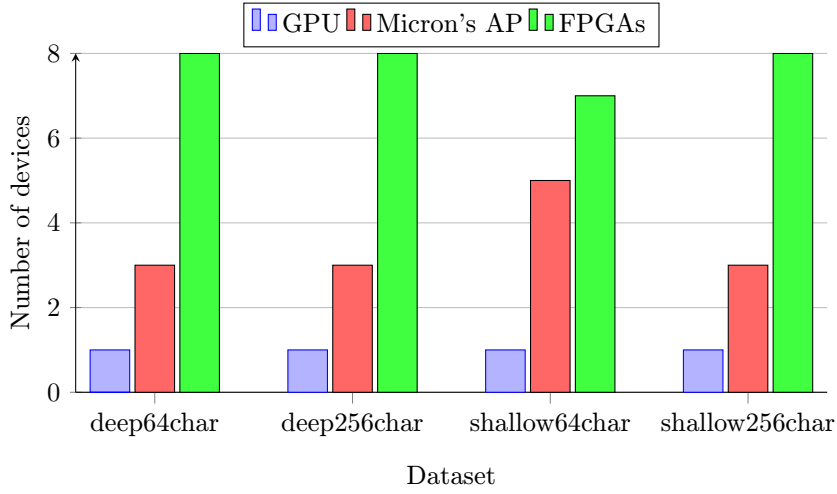


Figure 5: Number of devices, part of figure 5 from [17]

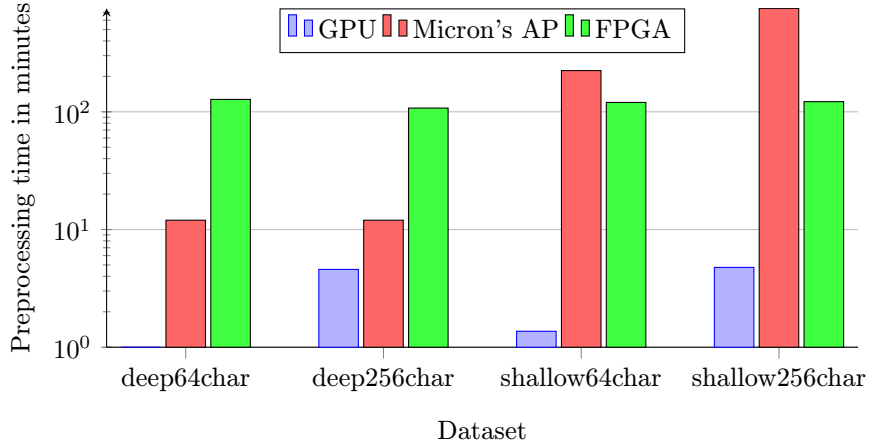


Figure 6: Preprocessing time, part of figure 5 from [17]

5.3 Papers on speeding up FSA processing on GPUs

On-the-Fly Principled Speculation for FSM Parallelization [28], published in 2015, researched speculative execution of finite state machines. Their research was focused on DFAs. They have employed techniques to accelerate DFA executions, which have subsequently been utilised to increase GPU acceleration of DFAs.

The paper *Scaling Out Speculative Execution of Finite State Machines with Parallel Merge* [24], published in 2020, develops optimisation techniques to improve the performance of DFAs execution on GPUs by extending existing work on speculative execution of FSMs and adapting it to GPU design. Their findings demonstrate significant enhancements on a scale ranging from tens to hundreds.

The latest paper on the execution of DFAs *GSpecPal: Speculation-Centric Finite State Machine Parallelization on GPUs* [23], published in 2022, have

addressed the existing bottlenecks, achieving an average speedup of 7.2x, but although it works with regular expressions, this approach has not been tested on any real dataset to verify its effectiveness in a real-world environment.

Why GPUs are Slow at Executing NFAs and How to Make Them Faster [15], published in 2020, improves the implementation of NFAs. The authors conducted measurements on the NVIDIA Quadro P6000 GPU and demonstrated more than a 25 times speedup in several applications, outperforming the AP chip in some real-world datasets, such as snort or YARA. However, it should be noted that their approach has limitations and is most effective in applications with a large number of states and is not usable with large NFA's that may occur.

The most recent research on NFAs, *Asynchronous Automata Processing on GPUs* [16] from 2023, improved upon previous work and achieved a system that achieves a speedup of 2-50 times depending on the degree of parallelization in previous solutions. No measurements have been taken on real-world datasets, and there is no information available on how this compares to other platforms.

5.4 OpenCL and CUDA implementations

All the studies, except for [18] which utilized OpenCL, employed CUDA for implementing their approaches. Although the research by [7] suggests that there is no significant difference in performance, writing optimal code appears more convenient in CUDA.

5.5 Heterogeneous computing

Although research has been conducted on combining multiple accelerator platforms [12], no work has been undertaken to utilize this method for accelerating regular expressions.

6 Results

Our analysis of the literature has demonstrated that the GPU is a viable platform for accelerating regular expressions. However, their use depends on the specific problem being addressed, and the optimal use case for GPUs remains unclear.

6.1 Finite Automata

The study [27] shows that DFA is the fastest solution for regular expression matching. However, a limiting factor can be the size of the DFA, especially on large datasets with multiple states, which can exceed the GPU's memory capacity. This can be alleviated by using a compressed DFA, which is slower than the uncompressed version, but faster than Non-Deterministic Finite Automata (NFA). Numerous studies have been undertaken to enhance the performance of DFA on GPU, with reference to [23, 24, 28]. The state of the art solution [23] should offer the fastest speed for some applications, but this solution has not been tested against current FPGA or AP solutions and is not known preprocessing time.

NFA has a smaller size and a much faster preprocessing time compared to DFA, although it is slower in operation, as stated in [17]. Since the time of

measuring in [27], several studies [15, 16] have been carried out to boost NFA performance on GPUs. The latest solution [16] provides a speedup over AP chip in some real-life datasets.

No current study has yet offered data on when it is optimal to use DFA or NFA.

6.2 Other platforms

The research paper [17] demonstrates that for regular expression matching, FPGA and AP produce faster results compared to GPU. However, current state-of-the-art GPU solutions offer speeds that are comparable to FPGA and AP, while having significantly faster preprocessing time, as illustrated in Figure 6. While FPGA and AP remain superior solutions for datasets that are static and for businesses that can afford specialized solutions, GPUs are present in most devices and can be utilized for a multitude of tasks.

Currently, no studies have been conducted on the combination of multiple accelerator platforms for the acceleration of regular expressions. However, prior research has examined the potential of integrating multiple accelerator platforms [12]. It is our belief that such an approach could be utilised to accelerate regular expressions, but further research is required.

6.3 OpenCL and CUDA implementations

All the aforementioned studies, with the exception of [18] which used OpenCL, used CUDA to implement their approaches. According to research by [7], there is no significant difference in performance, although writing optimal code in CUDA seems to be more convenient, but further research is needed. However, using the CUDA implementation limits the ability to run the code on anything other than Nvidia GPUs.

7 Discussion and Future work

Our analysis didn't provide a definitive answer to the question of when to use GPU acceleration of regular expressions. Although work has been done to optimise the performance of regular expressions on GPUs, there's still a gap in terms of their use in real-world applications. The current state of the art solutions offer theoretically comparable performance to FPGA and AP, but the lack of real-world data makes it difficult to determine the optimal use case for GPUs and whether it is worth the effort to implement GPU acceleration. The lack of current research indicates that the topic is not currently a priority for researchers, most likely due to the fact that current solutions are sufficient for most use cases. However, there are still some areas that could be explored.

CUDA was the most widely used platform in the research, but there is a lack of research into the use of OpenCL. Using OpenCL would allow the code to run on a wider range of GPUs, but it is not known whether the performance or ease of implementation would be comparable to CUDA. Another area that could be explored is the combination of different accelerator platforms. The research shows that it is possible to combine multiple accelerator platforms, but it is not known whether this would be beneficial for the same use cases.

8 Conclusions

Our study of the acceleration of regular expressions demonstrates that the GPU shows potential as a viable platform, depending on specific use cases and computational requirements. Deterministic Finite Automata (DFA) are able to perform at high speeds, although they encounter difficulties when processing large datasets due to their size. Non-Deterministic Finite Automata (NFA) present a feasible alternative since they are smaller in size and have a faster pre-processing, but the optimal circumstances for using either DFA or NFA remain unclear.

Field-Programmable Gate Arrays (FPGA) and Application Processors (AP) are superior to GPUs in terms of regular expression matching speed. Nevertheless, GPUs currently compete with FPGA and AP speeds with notably expedited preprocessing.

The possibility of combining multiple accelerator platforms with GPUs for regular expression acceleration, including an OpenCL comparison with CUDA, remains unexplored and warrants further investigation.

References

- [1] A. Asaduzzaman, A. Trent, S. Osborne, C. Aldershof, and F. N. Sibai. Impact of cuda and opencl on parallel and distributed computing. In *2021 8th International Conference on Electrical and Electronics Engineering (ICEEE)*, pages 238–242, 2021.
- [2] M. Becchi and P. Crowley. An improved algorithm to accelerate regular expression evaluation. In *Proceedings of the 3rd ACM/IEEE Symposium on Architecture for networking and communications systems*, pages 145–154, 2007.
- [3] M. Becchi, M. Franklin, and P. Crowley. A workload for evaluating deep packet inspection architectures. In *2008 IEEE International Symposium on Workload Characterization*, pages 79–89. IEEE, 2008.
- [4] N. Cascarano, P. Rolando, F. Risso, and R. Sisto. infant: Nfa pattern matching on gpgpu devices. *ACM SIGCOMM Computer Communication Review*, 40(5):20–26, 2010.
- [5] C. Chapman and K. T. Stolee. Exploring regular expression usage and context in python. In *Proceedings of the 25th International Symposium on Software Testing and Analysis*, ISSTA 2016, page 282–293, New York, NY, USA, 2016. Association for Computing Machinery.
- [6] J. C. Davis, L. G. Michael IV, C. A. Coghlan, F. Servant, and D. Lee. Why aren’t regular expressions a lingua franca? an empirical study on the re-use and portability of regular expressions. In *Proceedings of the 2019 27th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, ESEC/FSE 2019, page 443–454, New York, NY, USA, 2019. Association for Computing Machinery.

- [7] J. Fang, A. L. Varbanescu, and H. Sips. A comprehensive performance comparison of cuda and opencl. In *2011 International Conference on Parallel Processing*, pages 216–225, 2011.
- [8] A. Fuchs and D. Wentzlaff. The accelerator wall: Limits of chip specialization. In *2019 IEEE International Symposium on High Performance Computer Architecture (HPCA)*, pages 1–14. IEEE, 2019.
- [9] R. Gunawan, A. Rahmatulloh, I. Darmawan, and F. Firdaus. Comparison of web scraping techniques : Regular expression, html dom and xpath. In *Proceedings of the 2018 International Conference on Industrial Enterprise and System Engineering (IcoIESE 2018)*, pages 283–287. Atlantis Press, 2019/03.
- [10] N.-F. Huang, H.-W. Hung, S.-H. Lai, Y.-M. Chu, and W.-Y. Tsai. A gpu-based multiple-pattern matching algorithm for network intrusion detection systems. In *22nd International Conference on Advanced Information Networking and Applications-Workshops (aina workshops 2008)*, pages 62–67. IEEE, 2008.
- [11] Z. István, D. Sidler, and G. Alonso. Runtime parameterizable regular expression operators for databases. In *2016 IEEE 24th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM)*, pages 204–211, 2016.
- [12] R. Kobayashi, N. Fujita, Y. Yamaguchi, A. Nakamichi, and T. Boku. Gpu-fpga heterogeneous computing with opencl-enabled direct memory access. In *2019 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW)*, pages 489–498. IEEE, 2019.
- [13] V. W. Lee, C. Kim, J. Chhugani, M. Deisher, D. Kim, A. D. Nguyen, N. Satish, M. Smelyanskiy, S. Chennupaty, P. Hammarlund, R. Singhal, and P. Dubey. Debunking the 100x gpu vs. cpu myth: An evaluation of throughput computing on cpu and gpu. *SIGARCH Comput. Archit. News*, 38(3):451–460, jun 2010.
- [14] C.-H. Lin, C.-H. Liu, and S.-C. Chang. Accelerating regular expression matching using hierarchical parallel machines on gpu. In *2011 IEEE Global Telecommunications Conference - GLOBECOM 2011*, pages 1–5, 2011.
- [15] H. Liu, S. Pai, and A. Jog. Why gpus are slow at executing nfas and how to make them faster. In *Proceedings of the Twenty-Fifth International Conference on Architectural Support for Programming Languages and Operating Systems*, ASPLOS '20, page 251–265, New York, NY, USA, 2020. Association for Computing Machinery.
- [16] H. Liu, S. Pai, and A. Jog. Asynchronous automata processing on gpus. *Proc. ACM Meas. Anal. Comput. Syst.*, 7(1), mar 2023.
- [17] M. Nourian, X. Wang, X. Yu, W.-c. Feng, and M. Becchi. Demystifying automata processing: Gpus, fpgas or micron’s ap? In *Proceedings of the International Conference on Supercomputing*, ICS '17, New York, NY, USA, 2017. Association for Computing Machinery.

- [18] P. Rościszewski, J. Cychnerski, and A. Brzeski. A regular expression matching application with configurable data intensity for testing heterogeneous hpc systems. *Information Systems Architecture and Technology*, page 49, 2014.
- [19] M. C. Schatz and C. Trapnell. Fast exact string matching on the gpu. *Center for Bioinformatics and Computational Biology*, 2007.
- [20] Y. Sun, N. B. Agostini, S. Dong, and D. Kaeli. Summarizing cpu and gpu design trends with product data. *arXiv preprint arXiv:1911.11313*, 2019.
- [21] J. Van Lunteren, C. Hagleitner, T. Heil, G. Biran, U. Shvadron, and K. Atasu. Designing a programmable wire-speed regular-expression matching accelerator. In *2012 45th Annual IEEE/ACM International Symposium on Microarchitecture*, pages 461–472. IEEE, 2012.
- [22] X. Wang. *Techniques for efficient regular expression matching across hardware architectures*. University of Missouri-Columbia, 2014.
- [23] Y. Wang, R. Watling, J. Qiu, and Z. Wang. Gspecpal: Speculation-centric finite state machine parallelization on gpus. In *2022 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, pages 481–491. IEEE, 2022.
- [24] Y. Xia, P. Jiang, and G. Agrawal. Scaling out speculative execution of finite-state machines with parallel merge. In *Proceedings of the 25th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming, PPOPP '20*, page 160–172, New York, NY, USA, 2020. Association for Computing Machinery.
- [25] C. Xu, J. Su, and S. Chen. Exploring efficient grouping algorithms in regular expression matching. *PloS one*, 13(10):e0206068, 2018.
- [26] V. Yaneva, A. Rajan, and C. Dubach. Gpu acceleration of finite state machine input execution: Improving scale and performance. *Software Testing, Verification and Reliability*, 32(1):e1796, 2022.
- [27] X. Yu and M. Becchi. Gpu acceleration of regular expression matching for large datasets: Exploring the implementation space. In *Proceedings of the ACM International Conference on Computing Frontiers, CF '13*, New York, NY, USA, 2013. Association for Computing Machinery.
- [28] Z. Zhao and X. Shen. On-the-fly principled speculation for fsm parallelization. *ACM SIGPLAN Notices*, 50(4):619–630, 2015.
- [29] Y. Zu, M. Yang, Z. Xu, L. Wang, X. Tian, K. Peng, and Q. Dong. Gpu-based nfa implementation for memory efficient high speed regular expression matching. *SIGPLAN Not.*, 47(8):129–140, feb 2012.