

# Demystifying Automata Processing: GPUs, FPGAs or Micron's AP?

Marziyeh Nourian<sup>1,3</sup>, Xiang Wang<sup>1</sup>, Xiaodong Yu<sup>2</sup>, Wu-chun Feng<sup>2</sup>, Michela Becchi<sup>1,3</sup>

<sup>1,3</sup>Department of Electrical and Computer Engineering, <sup>2</sup>Department of Computer Science

<sup>1</sup>University of Missouri, <sup>2</sup>Virginia Tech, <sup>3</sup>North Carolina State University

mnouria@ncsu.edu, xw7b4@mail.missouri.edu, xdyu@vt.edu, feng@cs.vt.edu, mbecchi@ncsu.edu

## ABSTRACT

Many established and emerging applications perform at their core some form of pattern matching, a computation that maps naturally onto finite automata abstractions. As a consequence, in recent years there has been a substantial amount of work on high-speed automata processing, which has led to a number of implementations targeting a variety of parallel platforms: CPUs, GPUs, FPGAs, ASICs, and Network Processors. More recently, Micron has announced its Automata Processor (AP), a DRAM-based accelerator of non-deterministic finite automata (NFA). Despite the abundance of work in this domain, the advantages and disadvantages of different automata processing accelerators and the innovation space in this area are still unclear.

In this work we target this problem and propose a toolchain to allow an apples-to-apples comparison of NFA acceleration engines on three platforms: GPUs, FPGAs and Micron's AP. We discuss the automata optimizations that are applicable to these three platforms. We perform an evaluation on large-scale datasets: to this end, we propose an NFA partitioning algorithm that minimizes the number of state replications required to maintain functional equivalence with an unpartitioned NFA, and we evaluate the scalability of each implementation to both large NFAs and large numbers of input streams. Our experimental evaluation covers resource utilization, traversal throughput, and preprocessing overhead and shows that the FPGA provides the best traversal throughputs (on the order of Gbps) at the cost of significant preprocessing times (on the order of hours); GPUs deliver modest traversal throughputs (on the order of Mbps), but offer low preprocessing times (on the order of seconds or minutes) and good pattern densities (they can accommodate large datasets on a single device); Micron's AP delivers throughputs, pattern densities, and preprocessing times that are intermediate between those of FPGAs and GPUs, and it is most suited for applications that use datasets consisting of many small NFAs with a topology that is fixed and known a priori.

## 1 INTRODUCTION

Many established and emerging applications perform at their core some form of pattern matching, a computation that maps

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [Permissions@acm.org](mailto:Permissions@acm.org).

ICS '17, June 14-16, 2017, Chicago, IL, USA  
 © 2017 Association for Computing Machinery.  
 ACM ISBN 978-1-4503-5020-4/17/06...\$15.00  
<http://dx.doi.org/10.1145/3079079.3079100>

naturally onto finite automata abstractions. In biology, for example, several genomics tasks, such as motif discovery, orthology inference, shotgun and de novo assembly, involve string-matching operations on genomics data. In turn, advances in DNA sequencing technology have led to increasingly large volumes of data available for these applications, resulting in a significant increase in their computational requirements. In the networking domain, several applications such as network intrusion detection, content-based routing, and application-level filtering require inspecting network packets for potentially large sets of predefined patterns, and they typically must perform this operation at the rate of packet arrival on the router interface.

Given the number and relevance of applications requiring efficient pattern matching, there has been a substantial amount of work on high-speed automata processing, and this work has originated from different communities: from the networking to the reconfigurable computing and computer architecture to the parallel computing community. These efforts have led to a number of algorithmic [1-9] and architectural solutions targeting different parallel platforms: from CPUs to GPUs [10-12] to FPGAs [13-16] to ASICs [17-19] to Network Processors [20]. More recently, Micron has announced their Automata Processor [21], a DRAM-based accelerator of non-deterministic finite automata (NFA) that has been showcased on a variety of applications: motif discovery in biological sequences [22], association rule mining [23], brill tagging [24], high-speed regular expression matching for network intrusion detection [25], graph processing [26], and sequential pattern mining [27].

Despite this abundance of work on high-speed automata processing, there is still lack of clarity as to how existing software and hardware solutions are related to and compare with each other. There are several reasons for this.

First, existing solutions are based on different automata models: either non-deterministic or deterministic finite automata (NFAs and DFAs, respectively). While functionally equivalent, NFAs and DFAs have practical differences – in terms of resource requirements and traversal behavior – that are strongly dependent on the characteristics of the underlying pattern set. While there has been a substantial body of work proposing automata designs that trade off the advantages and disadvantages of NFAs and DFAs [1-9], no automata model is preferable on all datasets. This makes it hard to provide a fair comparison between automata processors relying on different automata models.

Second, some automata processing architectures are designed to optimize the peak performance of a single input stream, while others offer better support for stream-level concurrency.

Third, applications relying on finite automata must operate

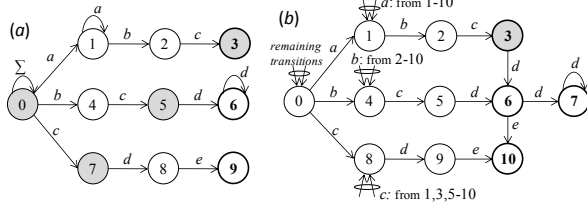


Figure 1: (a) NFA and (b) DFA accepting regular expressions  $a^+bc$ ,  $bcd^+$  and  $cde$ . Accepting states are bold. States active after processing text  $aabc$  are colored gray.

in two steps: in the *preprocessing step*, the required automaton must be generated, optimized, compiled, and loaded onto the target accelerator (through memory configuration and/or place&route operations); in the *traversal step*, the application performs pattern matching by traversing the automaton guided by the content of the input text. Most of the existing automata processing engines have been designed to optimize automata traversal, often at the cost of a significant preprocessing cost. While the preprocessing time is unimportant for some categories of applications (for example, network intrusion detection systems can operate for days or weeks between reconfigurations of their pattern sets), its effect on performance can be significant for other applications with more dynamic pattern sets or traversal times in the order of a few seconds. Unfortunately, the majority of the previous studies on Micron’s AP neglect to report the preprocessing overhead (or part of it) [21-24] or report substantial speedups over preexisting CPU tools (not necessarily based on automata) by comparing the *full* execution time of these tools to *only* the traversal time of the automata-based solution (on the order of seconds or milliseconds), omitting the preprocessing time of the AP-design (on the order of minutes) in the speedup calculation [25]. This can lead to results that are misleading or of limited practical use.

To target these problems and provide an apples-to-apples comparison, we select automata accelerator designs that rely on the same automata model: NFAs. Since NFAs do not suffer from state explosion, their use allows us to perform an evaluation on large-scale datasets without posing any restrictions on the kind of patterns supported. Specifically, we compare GPU- and FPGA-based NFA engines with Micron’s AP. Micron’s AP extends NFAs’ functionality with counters and boolean elements. To ensure functional equivalence and the same degree of programmability across the considered platforms, we extend existing FPGA- and GPU-based designs to support these features, and we adopt the same programming interface for all platforms: namely, Micron’s Automata Network Markup Language (ANML). Different platforms offer different automata density – to take this into account, we perform an analysis on non-trivial dataset sizes, which require partitioning large NFAs across multiple devices. Besides considering peak performance on a single input stream, we evaluate the scalability of the considered automata processor designs to multiple concurrent inputs. Finally, we evaluate the costs of the different preprocessing steps required by the considered architectures, and we study how the size of the automaton and the density of its transitions affect some of the preprocessing stages (e.g., place&route on Micron’s

AP and FPGA).

To summarize, we make the following contributions:

- We extend existing FPGA- and GPU-based automata processing designs to support Micron’s AP counters and boolean elements, and we propose a compiler toolchain to automatically deploy extended NFAs (in ANML form) onto these three platforms.
- We propose an NFA partitioning scheme aimed at minimizing the amount of state replication required to handle large NFAs while preserving functional equivalence with a single unpartitioned NFA.
- For GPU deployment, we explore different state layouts and kernels suited to NFAs with varying characteristics.
- We perform an apples-to-apples comparison between Micron’s AP, GPU- and FPGA-based NFA accelerator designs on large-scale datasets. Our evaluation covers resource utilization, throughput and preprocessing costs for real-world NFAs used in networking and bioinformatics applications, as well as synthetic datasets covering regular expressions datasets with various characteristics.

## 2 BACKGROUND AND RELATED WORK

### 2.1 Background on Automata Processing

Regular expression matching has traditionally been implemented by representing the pattern-set through finite automata (FA) [28]. The matching operation is equivalent to a FA traversal guided by the content of the input stream. Worst-case performance guarantees can be offered by bounding the amount of processing performed per input character. However, techniques to keep per-character processing low involve increasing the size of the finite automaton, the basic data structure in the regular expression matching engine. As the size of pattern-sets and the expressiveness of individual patterns increase, limiting the size of the automaton to fit on reasonably provisioned hardware platforms becomes challenging. Thus, the exploration space is characterized by a trade-off between the size of the automaton and the worst-case bound on the amount of per character processing.

NFAs and DFAs are at the two extremes in this exploration space. NFAs have a limited size but can require expensive per-character processing, whereas DFAs offer limited per-character processing at the cost of a possibly large automaton. In Figure 1 we show the NFA and DFA accepting three simple patterns ( $a^+bc$ ,  $bcd^+$  and  $cde$ ). In the figure, states active after processing text  $aabc$  are colored gray. In the NFA, the number of states and transitions is limited by the number of symbols in the pattern-set. In the DFA, every state presents one transition for each character in the alphabet ( $\Sigma$ ). Each DFA state corresponds to a set of NFA states that can be simultaneously active [28]; therefore, the number of states in a DFA equivalent to an  $N$ -state NFA can potentially be  $2^N$ . In practice, previous work [2, 5, 29] has shown that this so-called “state explosion” happens only in the presence of complex patterns (typically those containing repetitions of large character sets). Since each DFA state corresponds to a set of simultaneously active NFA states, DFAs

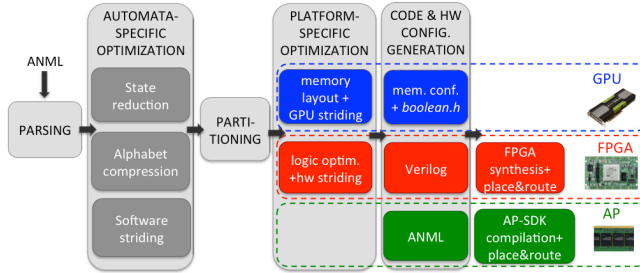


Figure 2: Our toolchain

ensure minimal per-character processing (only one state transition is taken for each input character).

From an implementation perspective, existing regular expression matching engines can be classified into two categories: memory-based [1-12, 17, 19], and logic-based [13-16]. Within the former, the FA is stored in memory; within the latter, it is stored in combinational and sequential logic. Memory-based implementations can be deployed on various platforms (GPUs, network processors, ASICs, FPGAs); logic-based implementations typically target FPGAs. In a memory-based implementation, design goals are the minimization of the memory size needed to store the automaton and of the memory bandwidth needed to operate it. Similarly, in a logic-based implementation the design aims at minimizing the logic utilization while allowing fast operation (that is, a high clock frequency). Existing proposals targeting DFA-based, memory-centric solutions have focused on designing compression mechanisms to reduce the DFA memory footprint and novel automata to alleviate the state explosion problem [1-9]. Despite the complexity of their design, memory-centric solutions have three advantages: fast reconfigurability, low power consumption, and scalability in the number of input streams. On the other hand, logic-centric solutions allow for easily achieving peak worst-case performance on a single input stream, at the expense of lack of scalability in the number of concurrent inputs.

## 2.2 Micron’s Automata Processor Overview

Micron’s Automata Processor [21] is a DRAM-based, reconfigurable accelerator that simulates NFA traversal at high speed. The AP includes three kinds of programmable elements stored in SDRAM: State Transition Elements (STE), Counter Elements (CE) and Boolean Elements (BE), which implement states/transitions, counters and logical operators between states, respectively. Each STE includes a 256-bit mask (one bit per ASCII symbol), and symbols triggering state transitions are associated to states (and encoded into STEs) rather than to transitions. Transitions between states are then implemented through a routing matrix consisting of programmable switches, buffers, routing lines, and cross-point connections. The routing capacity is limited by tradeoffs between clock rate, propagation delays and power consumption, and these constraints influence place&route of automata onto the AP hardware. Micron’s current generation of AP board (AP-D480) includes 16 or 32 chips organized into two to four ranks (8 chips per rank), and its design can scale up to 48 chips. Each AP chip consists of two half-cores. There are no routes either between half-cores or

inter-chips, which implies that NFA transitions across half-cores and chips are not possible. Programmable elements are organized in blocks: each block consists of 16 rows, where a row includes eight groups of two STEs and one special purpose element (CE or BE). Each chip contains a total of 49,152 STEs, 768 CE and 2,304 BE, organized in 192 blocks and equally residing in both half-cores. Current boards allow up to 6,144 elements per chip to be set as report elements.

AP automata can be described in ANML (an XML-based language). Recently proposed high-level programming languages for the AP are mapped and compiled into ANML [30]. Micron’s SDK includes a toolchain that parses ANML designs, compiles them into internal objects consisting of subgraphs, places and routes these subgraphs onto the AP hardware, and finally generates a binary image that can be used to program the AP memory and routing matrix. Once the AP has been programmed, it will be able to simulate the NFA traversal. AP chips can be grouped into logical cores of 2, 4 or 8, each processing a stream of 8-bit input characters [25]. The AP nominally operates at a 133MHz frequency, and, in absence of matches, it processes one input character per clock cycle from all input streams. Once matches occur, AP generates reporting events in vector format and stores them in an output-buffer; reporting matches to the host system requires from 91 to 291 clock cycles.

## 3 TOOLCHAIN

### 3.1 Overall design

Figure 2 shows the toolchain designed to deploy ANML specifications on GPU, FPGA and Micron’s AP. In the figure, grey boxes represent the software components that we have designed and implemented. The last two modules leading to FPGA and AP are Xilinx and Micron’s software development kits used for the final synthesis/compilation, map, and place&route on these two devices.

The input to the toolchain is an ANML file that contains one or more automata networks (each including one or more NFAs). We don’t impose any constraints on these networks: in other words, they don’t need to be designed to fit a particular device or optimized for it. Once parsed, these networks are stored in our toolchain using an internal representation for later processing and optimization. We distinguish two categories of optimizations: *automata-specific* and *platform-specific*. Since the GPU, FPGA and AP are used as NFA traversal accelerators, optimizations to the automaton apply to all platforms. In our previous work [14], we have described several NFA optimizations (*state reduction*, *alphabet compression* and *software striding*) and put them to practice on FPGA; these optimizations apply to GPUs and AP as well. Automata-specific optimizations can be selectively enabled and disabled. Platform-specific optimizations are related to the way the NFA is encoded for the particular target device; these optimizations include compact and efficient memory encodings, logic utilizations, and striding mechanisms that are specific to a particular hardware platform. Since the internals of the operation of the AP hardware and its software stack (including the compilation, map and place&route processes) are proprietary, AP-specific optimizations are

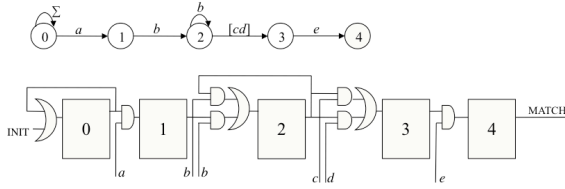


Figure 3: NFA accepting regular expressions  $ab+[cd]e$  and corresponding one-hot encoding representation

deferred to the AP SDK tools (last phase of the toolchain). The *partitioning* step, that takes a potentially large network and breaks it into multiple NFA partitions so that each of them can fit the target hardware, is performed after the automata-specific optimization step. This allows partitioning to be done on an already optimized NFA. Our partitioning algorithm is platform-independent, but its configuration depends on the target platform. The *code and configuration generation step* produces the files required for the final deployment of the automata network on the hardware. For GPUs, all is needed is a configuration file that includes the information necessary to load the NFA partitions into memory, and a header file with the definition of the boolean connectors in the ANML specification. FPGAs are configured through a Verilog file describing the NFA network and its interface. The AP is configured through an ANML file; this output file differs from the input file in that it contains a partitioned and optimized automata network.

### 3.2 GPU implementation

We reuse and extend iNFAnt [18], an NFA-traversal engine for GPUs. iNFAnt stores the NFA in device memory, and encodes the transition table as set of (source, destination) pairs indexed by the input character. In order to allow efficient execution, iNFAnt stores the set of active states in shared memory in bit-vector form. For each input character, iNFAnt retrieves from memory all the transitions on that symbol, and, if their source state is active, the engine updates the active state vector with the destination state information. In iNFAnt, each thread-block is assigned an input stream, and threads within a block process the state transitions and update the state vector cooperatively.

We extended iNFAnt with the following functionalities:

*Support for multiple NFA partitions* – We map each NFA partition to a thread-block, allowing multiple blocks to process the same input stream on different partitions. The transition lists corresponding to different partitions are laid out sequentially, and an indexing array maps each partition to the proper set of thread-blocks, each operating on a different input stream.

*Traversal kernels based on compressed sparse row (CSR) layout* – We consider an alternative memory layout where transitions represented as (input symbol, destination) pairs are indexed by the source state. For each input symbol, this layout allows processing only the transitions that originate from active states. We store the identifiers of the active states in a queue in global memory. We consider two variants of this kernel: *CSR-state* and *CSR-tx*, the former mapping active states to threads, and the latter mapping outgoing transitions from active states to threads.

*Support for counters and boolean elements* – We associate a

special state to each counter and boolean element, and store these special states at the end of the state vector. The activation of special states triggers code implementing the operation of the particular counter or boolean element. Boolean operators are also associated combinational code that is stored in an automatically generated header file.

### 3.3 FPGA implementation

On FPGA, NFA processing can be realized in two ways: either by implementing a traversal engine that accesses the NFA stored in memory, or by directly encoding the NFA in logic. Most logic-based NFA implementations are based on the one-hot encoding scheme [13], in which states are represented as flip-flops while transitions are implemented by and-ing and or-ing the outputs of the flip-flops with the decoded input character. For example, Figure 3 shows the one-hot encoding representation of the NFA accepting regular expression  $ab^+[cd]e$ . The main advantage of this scheme is that it limits the traversal time to one clock cycle per input character independent of the number of states that are active (this property is shared by Micron’s AP). On the other hand, this implementation suffers from two limitations: first, updating the NFA requires reprogramming the device; second, multiple input stream support requires logic replication. The pros and cons of a memory-based FPGA design are comparable to those of a GPU solution: easy support for multiple input streams at the cost of irregular and unpredictable memory access patterns, leading to dataset dependent performance. In this paper we use the optimized logic-based implementation that we have described in our previous work [14], and extend it to support counters and boolean elements (a trivial extension).

### 3.4 Automata-specific optimizations

Our toolchain includes three automata-specific optimizations: state reduction, alphabet compression and software striding [14]. Here, we briefly mention their effect on the considered platforms.

*State reduction* (which merges duplicate NFA paths) reduces the memory requirements on GPU and AP, and the logic requirements on FPGA. In addition, it reduces the number of states that can be active in parallel, which for GPU is beneficial to the throughput.

*Alphabet compression* (which consolidates the alphabet based on the symbols appearing on the NFA transitions) reduces the wiring and LUT utilization on FPGA. However, because the AP stores a 256-bit mask in each STE, this optimization does not benefit AP unless combined with software striding.

*Software striding* (which allows processing multiple characters in one step) can be beneficial on all platforms if combined with alphabet reduction. This technique is applicable to the AP only if the alphabet generated by combining alphabet reduction and software striding does not exceed 256 symbols. GPUs and FPGAs offer also platform-specific striding schemes [6, 10, 15], which we have included in our toolchain.

### 3.5 Partitioning criteria

An NFA must be partitioned if it exceeds the resources

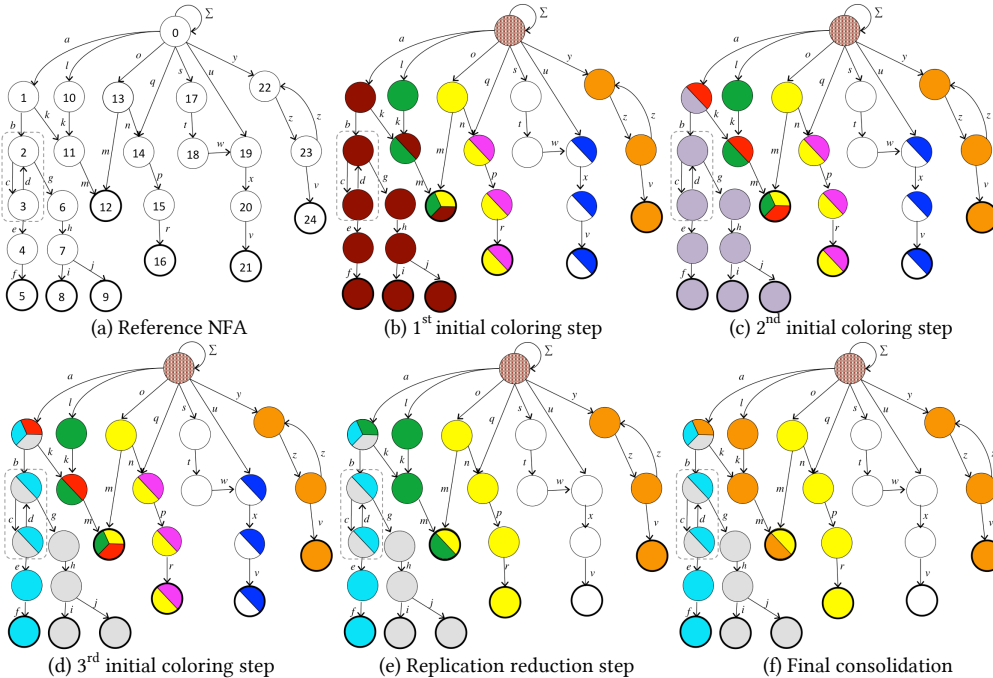


Figure 4: Example of application of our coloring scheme ( $N_{max} = 8$ ).

available on a particular device. Here, we indicate the platform-specific partitioning criteria we use. In Section 4, we describe our proposed partitioning algorithm.

**GPU:** GPU partitioning is required if the shared or global memory capacity is exceeded, or if the state identifier space is exhausted. In this paper, we use 16-bit state identifiers, leading to a maximum of 65536 states per NFA partition. This constraint is more restrictive than those on the global and shared memory capacity (and, due to thread-block concurrency, is not a limiting factor on performance – see Section 5).

**FPGA:** The logic design used stores states in flip-flops and transitions in LUTs. We experimentally found flip-flops to be the bottleneck resource. Thus, we perform NFA partitioning when the number of NFA states exceeds that of available flip-flops.

**AP –** The AP does not allow transitions across half-cores, and has a limited number of STEs, Counter Elements and Boolean Elements per half-core (see Section 2.2). Thus, the AP NFA partitioning criterion is based on these constraints.

## 4 NFA PARTITIONING ALGORITHM

In this section, we describe our NFA partitioning algorithm. For the sake of simplicity, we discuss the algorithm on traditional NFAs: its extension to counters and boolean elements is straightforward. In order to preserve functional equivalence, NFA partitioning requires state replication. For example, let us assume to break the NFA of Figure 4(a) into two partitions to be deployed and operated on two devices: one partition containing states from 0 to 16, and the other containing states from 17 to 24. In order to maintain functional equivalence with the original NFA, the entry state 0, which is shared by the patterns matched in the both partitions, must be replicated into the second partition. In general, very large NFAs may require replication of

sets of states shared by several patterns.

The goal of our partitioning algorithm is to split the NFA into a small number of balanced partitions, while minimizing the required state replications. In particular, given a threshold  $N_{max}$  on the number of states that can be accommodated on a particular device or hardware component, the algorithm must split the NFA into as few partitions as possible, each with size not exceeding  $N_{max}$ . Balanced partitions allow load balancing within (for GPUs and Micron’s AP) and across (for FPGAs) devices, which ultimately has a positive effect on throughput. It is worth noting that existing partitioning schemes for

generic graphs [31] aim to minimize the size of the cut (number of inter-partition transitions), but, when applied to NFAs, they do not necessarily minimize the number of state replications required to preserve functional equivalence with the unpartitioned NFA. Thus, the need for a partitioning scheme tailored to NFAs.

We propose an algorithm that colors the NFA so that each color represents a partition, and states assigned multiple colors are shared across partitions and must be replicated in each of these partitions. In order to meet the requirements above, the algorithm must limit the number of colors and of states with multiple colors, while allowing each color to appear in up to  $N_{max}$  states. In the following, we call “color size” the number of states assigned a particular color.

Our algorithm operates in two phases: *initial coloring* and *color consolidation*. In the initial coloring phase, the NFA is traversed from the entry state and recursively colored until the size of each color doesn’t exceed  $N_{max}$ . The color consolidation phase consolidates multiple colors into one while keeping their size below the given threshold.

We note that sets of states connected by cyclic transitions (e.g., states 2 and 3 in Figure 4(a)) cannot be separated into multiple partitions. We recall that, for partitions to be independent, inter-partition activations (that is, cross-partition transitions) must be avoided. As a consequence, a state belonging to multiple partitions must be replicated along with all the states connected to it in a cyclic fashion. Thus, we group states that are cyclically interconnected into “super-states”, and we handle all the states in a super-state together. For example, states 2 and 3 of Figure 4(a) form super-state {2, 3} and are handled as a single state. In order to operate, the algorithm

Table 1: Dataset characteristics and traversal information (ranges correspond to traces with  $p_{\text{forw}}=0.5$  and  $p_{\text{forw}}=0.9$ ).

| Type           | Name            | NFA Characteristics |          |               | # Partitions |      |     | Traversal Information |                 |           |                     |
|----------------|-----------------|---------------------|----------|---------------|--------------|------|-----|-----------------------|-----------------|-----------|---------------------|
|                |                 | # states            | # trans. | # ANML states | GPU          | FPGA | AP  | Avg. active set       | Max. active set | # Matches | Inputs w/ matches % |
| Small NIDS     | l7-filter       | 2794                | 124k     | 3844          | 1            | 1    | 1   | 3.0-4.8               | 5-9             | 4-10      | 0-0                 |
|                | snort534        | 9514                | 79k      | 10k           | 1            | 1    | 1   | 1.7-22.8              | 10-37           | 9-533     | 0.01-0.2            |
| Bioinformatics | 10gene_8k       | 33k                 | 100k     | 55k           | 1            | 1    | 4   | 371.2                 | 617             | 12286     | 99.9                |
|                | 10gene_12k      | 86k                 | 258k     | 243k          | 2            | 1    | 16  | 378.8                 | 713             | 7443      | 14.2                |
|                | 10gene_16k      | 138k                | 415k     | 230k          | 3            | 2    | 14  | 378.0                 | 713             | 493       | 0.1                 |
|                | 10gene_20k      | 190k                | 570k     | 317k          | 4            | 2    | 20  | 377.4                 | 713             | 0         | 0.0                 |
|                | 100gene_8k      | 137k                | 413k     | 229k          | 3            | 2    | 16  | 670.8                 | 826             | 81476     | 100                 |
|                | 100gene_12k     | 621k                | 1863k    | 1035k         | 14           | 9    | 80  | 742.3                 | 1382            | 71140     | 68.6                |
|                | 100gene_16k     | 1124k               | 3372k    | 1873k         | 28           | 16   | 149 | 741.6                 | 1387            | 4997      | 1.04                |
|                | 100gene_20k     | 1619k               | 4858k    | 2699k         | 42           | 22   | 213 | 740.1                 | 1386            | 25        | 0.01                |
| Synth.         | deep-64char     | 800k                | 1801k    | 800k          | 14           | 8    | 76  | 2.1-2.5               | 7-7             | 0-0       | 0-0                 |
|                | deep-256char    | 800k                | 4855k    | 800.3k        | 13           | 8    | 72  | 6.3-6.7               | 11-11           | 0-0       | 0-0                 |
|                | shallow-64char  | 800k                | 1801k    | 800.1k        | 15           | 7    | 74  | 3.1-3.5               | 7-8             | 2-1904    | 0-2.9               |
|                | shallow-256char | 800k                | 4855k    | 800.1k        | 13           | 8    | 78  | 6.1-6.7               | 12-14           | 1-2106    | 0-3.21              |

requires super-states to include fewer than  $N_{max}$  states. If this is not the case, the NFA cannot be split into independent partitions. In the presence of dependent partitions, multiple NFA traversals are required to handle inter-partition activations. Fortunately, NFAs originated from regular expressions datasets tend to have only few super-states of small size. This is because backward-directed transitions in NFAs originate from sub-pattern repetitions within regular expressions (for example, sub-pattern ( $dc$ ) in Figure 4(a), where string  $dc$  can be repeated zero or more times). Sub-pattern repetitions are rare in real-world datasets, and are rarely shared by a large number of patterns.

We now detail the operation of the two phases of the algorithm, and illustrate them in Figure 4. In the example, we assume that the threshold  $N_{max}$  is equal to 8.

**Initial coloring** – The initial coloring procedure starts by assigning distinct colors to the states connected to the entry state (or to the super-state to which it belongs). This is illustrated in Figure 4(b), where the children of the entry state 0 are colored brown, green, yellow, pink, white, blue and orange. The colors are propagated to all the connected states following the transitions. The entry state is then assigned all the colors of its children. As can be seen, this leads to some states (states 11-12, 14-16, 19-21, beside state 0) be assigned multiple colors. This operation must be repeated recursively on all generated NFA partitions until their size does not exceed the threshold  $N_{max}$ . As can be seen, after the first coloring step the brown color has size 12 (states 0-9 and 11-12). Therefore, the coloring procedure is repeated starting from state 1. This causes color brown to be split into colors red and violet, which are again propagated down to the terminal states of the NFA (Figure 4(c)). Since color violet has size 10 (including state 0), the algorithm invokes one additional recursive step on super-state {2,3}, causing color violet to be split into colors cyan and grey (Figure 4(d)). Since the largest color (grey) has now size 8 (equal to  $N_{max}$ ), the initial coloring phase is terminated.

**Color consolidation** – While respecting the constraint on the maximum partition size, the partitioning generated by the initial

coloring step has two limitations: it includes small and unbalanced partitions, and it leads to significant state replication. In the example, states 2, 3, 11, 14-16, 19-21 must be replicated once, states 1 and 12 must be replicated twice, and state 0 must be replicated 8 times. The coloring consolidation phase aims to combine different colors into one so as to increase the partition size, decrease the number of partitions and the number of state replications required, and achieve more balanced partitions. This phase is broken down into two steps: *replication reduction* and *final consolidation*. The first step aims to reduce the number of state replications required by merging colors. To determine which colors to consolidate, we sort pairs of colors in descending order according to the number of state replications that their consolidation would save. In the example, cyan/grey, yellow/pink and white/blue would save 3 state replications, green/red would save 2, and red/yellow, green/yellow, red/cyan and red/grey would save only 1. We then consider all pair-wise consolidation opportunities in order, and merge the two colors only if their merging doesn't violate the partition size constraint. In the example, we consolidate yellow+pink into yellow, white+blue into white, and green+red into green. Figure 4(e) shows the result of the replication reduction step. In the final consolidation step, we look for opportunities to consolidate colors according to their size. To this end, we first sort the colors in descending order by size, and then traverse the list and consolidate each color with the next color in the list that doesn't lead to violating the partition size constraint (if such color exists). In the example, colors orange and green are consolidated into orange. Figure 4(f) shows the final coloring, which leads to 5 partitions: two of size 8 (grey and orange) and three of size 6 (cyan, yellow and white).

## 5 EXPERIMENTAL EVALUATION

### 5.1 Hardware platform

We conducted almost all our experiments on a machine equipped with a dual 6-core Intel Xeon CPU @ 2.66GHz and



**Table 2: Resource utilization for GPU (ranges correspond to different numbers of streams) and FPGA (ranges correspond to the minimum and maximum values across partitions)**

| Type           | Name            | GPU       |                        |             |            | FPGA      |                    |       |        |
|----------------|-----------------|-----------|------------------------|-------------|------------|-----------|--------------------|-------|--------|
|                |                 | # devices | GPU Memory Utilization |             |            | # devices | FPGA % utilization |       |        |
|                |                 |           | shared (KB)            | global (MB) |            |           | FF                 | LUT   | Slice  |
|                |                 |           | iNFAnt                 | CSR         |            |           |                    |       |        |
| Small NIDS     | l7-filter       | 1         | 0.6                    | 0.4         | 0.5-0.6    | 1         | 1                  | 1     | 11     |
|                | snort534        | 1         | 2.3                    | 0.3         | 0.4-0.4    | 1         | 5                  | 5     | 16     |
| Bioinformatics | 10gene_8k       | 1         | 8.1                    | 0.4         | 7.6-47.5   | 1         | 20                 | 33    | 86     |
|                | 10gene_12k      | 1         | 13.2                   | 1           | 13.1-78.1  | 1         | 53                 | 52    | 99     |
|                | 10gene_16k      | 1         | 12.3                   | 1.6         | 13.3-74    | 2         | 31-54              | 50-52 | 97-99  |
|                | 10gene_20k      | 1         | 14.8                   | 2.2         | 16.6-90.3  | 2         | 59-59              | 56-57 | 92-99  |
|                | 100gene_8k      | 1         | 14.9                   | 1.6         | 15.4-87.8  | 2         | 18-67              | 32-77 | 75-99  |
|                | 100gene_12k     | 1         | 14.8                   | 7.1         | 27-99.8    | 9         | 38-66              | 62-67 | 96-99  |
|                | 100gene_16k     | 1         | 15.8                   | 15.5        | 42.2-121.3 | 16        | 38-63              | 61-61 | 97-99  |
|                | 100gene_20k     | 1         | 15.7                   | 22.5        | 55-133.4   | 22        | 31-49              | 49-74 | 98-98  |
| Synt.          | deep-64char     | 1         | 15.5                   | 6.9         | 13.7-17.3  | 8         | 32-69              | 22-61 | 99-99  |
|                | deep-256char    | 1         | 15.8                   | 18.5        | 25.1-27.6  | 8         | 15-69              | 12-62 | 34-100 |
|                | shallow-64char  | 1         | 15.8                   | 6.8         | 13.6-17.3  | 7         | 56-68              | 60-76 | 94-96  |
|                | shallow-256char | 1         | 15.8                   | 18.5        | 25.1-27.6  | 8         | 29-69              | 23-75 | 97-99  |

64GB of memory, running CentOS 6.4. Since some of our AP syntheses run out-of-memory on that machine, for our AP experiments we used a server with similar hardware settings but equipped with 256GB memory. For our GPU experiments we used an Nvidia Titan X GPU (Maxwell architecture), equipped with 12GB of global memory and 24 streaming multiprocessors (SMs), each including 128 cores and 96KB of shared memory. We used CUDA 7.0. For our FPGA experiments we used a Xilinx XC6VLX130T device (Virtex-6 family), which includes 20,000 slices (for a total of 160,000 flip-flops and 80,000 LUTs). We used the Xilinx ISE Design suite v13.2 to perform synthesis, mapping and place&route of our HDL designs. This FPGA device was chosen because it is in the same price range (~\$1,200) as our GPU. For AP experiments, we refer to the architecture of a 32-chip AP-D480. Since Micron’s AP hardware is not yet available on the market, we don’t have pricing information for it. For the AP, we used AP SDK v. 1.6.5 to collect resource utilization and preprocessing data, and performed throughput projections using the nominal operating frequency (more details in Section 5.3).

## 5.2 Datasets

We selected datasets allowing to compare the three platforms on different application domains and on NFAs with varying characteristics in terms of number of states and transitions, alphabet size, connectivity and depth. To this end, we used three types on datasets: *small NIDS*, *bioinformatics*, and *synthetic*. Recently proposed benchmark suites for automata processing [32] are not meant for large scale analysis (they include NFAs with up to about 100k states). Table 1 (columns 3-5) summarizes the characteristics of the NFAs for the considered datasets.

*Small NIDS* (Snort538 and l7-filter) are small network intrusion detection datasets that include 538 and 116 regular expressions, respectively (see [10] for more details).

*Bioinformatics* datasets (*ngene\_kk*) consist of a set of Hamming distance automata used to address a motif-finding problem [33]. The problem requires identifying all the substrings

of length  $k$  that appear on multiple genes within hamming distance  $d$ , and can be found in a region of the gene of length  $l$ . Due to space limitation, here we show only the results for  $n$  genes from a yeast genome of about 5000 genes, with  $n=\{10, 100\}$ ,  $k=\{8, 12, 16, 20\}$ ,  $l=500$ ,  $d=2$  and a 4-symbol alphabet (A, C, G, T). A Hamming distance NFA has  $(k+1)(d+1)-d(d+1)/2$  states, and each gene region of length  $l$  leads to  $(l-k+1)$  of these NFAs. The NFAs in Table 1 (used on all three platforms) have been state-reduced. However, previous

work [22, 30] has shown that, on the AP, preprocessing time can be significantly reduced if NFAs with known structure are precompiled. Thus, on the AP we also use a non state-reduced variant of these bioinformatics datasets (see Table 5), leading to networks of  $n(l-k+1)$  small NFAs (each with  $(2d+1)k-d^2$  STEs) with fixed topology.

*Synthetic* automata exhibit the structure of NFAs accepting sets of regular expressions with shared prefixes: large state outdegrees in the proximity of the entry state, and low state outdegrees as we move deeper in the NFA. Our synthetic NFAs have configurable number of states, alphabet size, entry state outdegree, outdegree decrease factor  $\gamma$  (the outdegree decreases as  $\gamma^{\text{depth}}$ ), and frequency of wildcards, character sets and their repetitions. We set these parameters so as to generate 800k-state NFAs with two alphabet sizes (64 and 256), and two structures (*deep* and *shallow*, about 180- and 16-level deep, respectively).

Whenever required, we partition these NFAs with the algorithm described in Section 4. We recall that the partitioning threshold is platform-specific (Section 3.5). This leads to the number of partitions shown in Table 1 (columns 6-8). In order to simulate the NFA traversal, we use two kinds of input streams. For bioinformatics datasets, we generate traces of length 500,000 (for 1,000 genes) by randomly selecting symbols from the {A, C, G, T} alphabet. For NIDS and synthetic datasets, we generate 256k character traces through our trace generator [34], setting the probability to move deeper in the NFA ( $p_{\text{forw}}$ ) to 0.5 and 0.9. The traversal characteristics (average and maximum number of active states per input character, and number and frequency of matches) are reported in Table 1 (columns 9-12).

## 5.3 Results

**Resource utilization** is reported in Table 2 for GPU and FPGA, and in Table 5 (columns 7-12) for the AP.

For GPU, we recall that the NFA is stored in global memory, while the active state information (encoded in a bit vector) is stored in shared memory. In the CSR case, the active state

**Table 3: Traversal throughput (ranges correspond to different numbers of streams for GPU and to different partitions for FPGA)**

| Type           | Name            | GPU          |                                |                                   |                                | FPGA<br>Throughput<br>per device<br>(Gbps) |
|----------------|-----------------|--------------|--------------------------------|-----------------------------------|--------------------------------|--|
|                |                 | #<br>streams | iNFAnt<br>Throughput<br>(Mbps) | CSR-state<br>Throughput<br>(Mbps) | CSR-tx<br>Throughput<br>(Mbps) |  |
| Small<br>NIDS  | l7-filter       | 48-4800      | 73.0-675.7                     | 1.4-9.47                          | 2.2-19.7                       | 12.16                                      |
|                | snort534        | 48-4800      | 60.1-466.5                     | 24.2-39.9                         | 40.9-30.1                      | 19.82                                      |
| Bioinformatics | 10gene_8k       | 72-480       | 21.9-32.5                      | 3.1-3.2                           | 0.6-0.8                        | 2.51                                       |
|                | 10gene_12k      | 36-240       | 7.1-11.9                       | 4.7-4.8                           | 0.6-0.8                        | 1.68                                       |
|                | 10gene_16k      | 24-160       | 5.3-7.9                        | 8.9-14.0                          | 0.6-0.9                        | 1.8-3.2                                    |
|                | 10gene_20k      | 18-120       | 3.2-5.5                        | 7.3-12.5                          | 0.6-0.9                        | 1.2-1.6                                    |
|                | 100gene_8k      | 24-160       | 4.2-7.5                        | 0.5-0.6                           | 0.2-0.3                        | 1.1-3.0                                    |
|                | 100gene_12k     | 6-35         | 0.8-1.4                        | 0.6-0.6                           | 0.2-0.2                        | 1.5-2.4                                    |
|                | 100gene_16k     | 3-18         | 0.3-0.5                        | 1.2-2.2                           | 0.2-0.4                        | 0.7-2.2                                    |
|                | 100gene_20k     | 2-12         | 0.2-0.3                        | 1.0-1.7                           | 0.2-0.3                        | 1.7-3.5                                    |
| Synth.         | deep-64char     | 6-37         | 1.5-2.2                        | 1.8-2.7                           | 1.8-2.8                        | 2.5-3.6                                    |
|                | deep-256char    | 6-37         | 1.7-2.6                        | 0.5-1.7                           | 1.4-2.7                        | 1.7-2.9                                    |
|                | shallow-64char  | 5-32         | 1.2-1.9                        | 0.4-1.8                           | 0.5-2.2                        | 1.5-1.7                                    |
|                | shallow-256char | 6-35         | 1.7-2.7                        | 0.2-1.0                           | 0.3-2.6                        | 1.1-1.9                                    |

information is also stored (in queue format) in global memory, and therefore the global memory requirement increases with the number of thread-blocks run. However, as can be seen in Table 2, the global memory utilization is very limited even for the CSR format, and even the largest dataset occupies only up to 133MB of the 12GB global memory. We recall that NFA partitioning is driven by the use of 16-bit state identifiers, and shared memory stores two bitmaps indicating the states active at the beginning and the end of each traversal step. Therefore, the use of partitions with at most 64k states limits the per-block shared memory utilization to 16KB in the worst case, allowing at least 6 blocks to reside on a SM and hide each other’s memory latencies.

For FPGA, to facilitate the place&route process, we sized the partitions so as to use up to 70% of the flip-flop capacity. Since the considered device has twice as many flip-flops as LUTs, on most experiments this setting leads to near full slice utilization.

For the AP (Table 5), we report both the ideal utilization (the number of blocks and AP cores that a dataset would require based on the number of its STEs and reporting elements), and the utilization numbers reported by the AP’s SDK (real utilization). The utilization efficiency in column 11 is the ratio between the ideal and real block utilization. As can be seen, due to the place&route constraints on the routing matrix, the real utilization is significantly higher than the ideal one. Note that shallow synthetic datasets have significantly lower utilization efficiency (~20%) than deep ones (>80%): this is because the node out-degree of non-terminal states is large for shallow and low for deep datasets, making the former much harder to route. Due to the generally low utilization efficiency, we partitioned all the state-reduced NFAs so that each partition would require 50% (rather than the whole) half-core capacity. This led to the number of AP partitions shown

in Table 1 (column 8). In addition, we experienced that the AP SDK tools run out-of-memory when processing large datasets. To avoid this, for state-reduced NFAs we grouped partitions into batches of 32-64 (depending on the transition density of the dataset), and we run the AP SDK on one batch at a time. In the table, for each state-reduced dataset we report the cumulative results over all batches. In case of large fixed topology datasets (100gene\*), which consist of many small NFAs with the same topology, we sized each batch so as to use all 32 cores on the AP. Since the place&route algorithm used by the SDK is proprietary, this was a trial-and-error process. For these datasets, we report the number of batches (which corresponds to the number of AP boards required), and the per-batch data. As can be seen, for small  $k$  (i.e., small hamming distance NFAs) the place&route is easier and the number of STEs/batch and utilization efficiency are higher. Since larger hamming distance NFAs are harder to place, the utilization efficiency decreases as  $k$  increases.

Traversal throughput is computed using the following formulas, which assume 8-bit inputs.

$$throughput_{GPU} = \frac{N_{inputs} \times N_{streams} \times 8}{time_{kernel+transfer}}$$

$$throughput_{FPGA/AP} = \frac{N_{inputs} \times N_{streams} \times 8}{N_{inputs} + N_{output\_processing\_cycle}} \times f_{clk}$$

We assume that matches are reported every 64K inputs (maximum IP packet length) for NIDS datasets, every 500 inputs (length of relevant portion of a gene) for bioinformatics datasets, and every 1000 inputs for synthetic datasets ( $N_{inputs}$ ). For FPGA, we used the worst-case, post-place&route operating frequency reported by the Xilinx tools. The number of cycles required to report the matches ( $N_{output\_processing\_cycles}$ ) is equal to the ratio between the number of matching states in the NFA and the number of output pins on the FPGA device. For the AP, we performed estimates based on the 133 MHz nominal operating frequency and the 291 clock cycle output processing time.

**Table 4: Preprocessing overhead (in case of large datasets, we show minimum and maximum per-partition data)**

| Type           | Name            | GPU              |                       |                      | FPGA             |                       |                     |
|----------------|-----------------|------------------|-----------------------|----------------------|------------------|-----------------------|---------------------|
|                |                 | Parsing<br>(sec) | Mem. L.<br>gen. (sec) | Loading<br>mem. (ms) | Parsing<br>(sec) | Verilog gen.<br>(sec) | Synt.+<br>p&r (min) |
| Small<br>NIDS  | l7-filter       | 5                | 0.01                  | 0.004                | 0.4              | 0.6                   | 6                   |
|                | snort534        | 5.1              | 0.01                  | 0.003                | 0.6              | 0.3                   | 6                   |
| Bioinformatics | 10gene_8k       | 3.2              | 0.05                  | 3-6                  | 3.3              | 9.2                   | 37                  |
|                | 10gene_12k      | 5.8              | 0.3                   | 8-12                 | 7.2              | 58.28                 | 76                  |
|                | 10gene_16k      | 12               | 0.5                   | 12-18                | 4.3-7.7          | 19.6-60.2             | 69-78               |
|                | 10gene_20k      | 12.4             | 0.9                   | 16-26                | 7.8-8.2          | 69.8-71.7             | 89-129              |
|                | 100gene_8k      | 12.3             | 0.4                   | 12-18                | 3.1-11.2         | 7.8-99.9              | 34-165              |
|                | 100gene_12k     | 70               | 8.5                   | 47-76                | 5.5-9.2          | 29.2-90.4             | 54-153              |
|                | 100gene_16k     | 114.8            | 36.7                  | 99-167               | 5-8.4            | 27.9-78.1             | 54-554              |
|                | 100gene_20k     | 193.2            | 80.8                  | 145-249              | 4-8.9            | 19.9-93.3             | 70-92               |
| Synth.         | deep-64char     | 61.6             | 7.0                   | 40-80                | 2.6-5.3          | 0.4-0.9               | 73-178              |
|                | deep-256char    | 267.8            | 8.4                   | 110-150              | 1.5-6.6          | 0.5-2.6               | 15-189              |
|                | shallow-64char  | 76.5             | 6.1                   | 40-80                | 6.6-7.9          | 0.8-1.0               | 103-133             |
|                | shallow-256char | 279.5            | 6.8                   | 120-150              | 3.6-8.5          | 3.9-9.9               | 38-192              |



Table 5: AP Results

| Type           | Name           | # batches   | ANML-NFA characteristics |                |                 | Ideal utilization |          | Resource utilization from SDK profiling |          |                      | # AP boards | SDK preprocessing time |           |             | Throughput per device (Mbps) |      |      |
|----------------|----------------|-------------|--------------------------|----------------|-----------------|-------------------|----------|---|----------|----------------------|-------------|------------------------|-----------|-------------|------------------------------|------|------|
|                |                |             | # states/ STEs           | # start states | # report states | # cores           | # blocks | # cores                                 | # blocks | % utiliz. efficiency |             | Comp (sec)             | p&r (sec) | total (min) |                              |      |      |
| Small NIDS     | L7-Filter      | 1           | 4k                       | 78             | 489             | 1                 | 16       | 1                                       | 99       | 16.1                 | 1           | 1                      | 732       | 16          | 16946                        |      |      |
|                | Snort534       | 1           | 11k                      | 24             | 582             | 1                 | 43       | 1                                       | 135      | 31.8                 | 1           | 2                      | 2261      | 38          | 16946                        |      |      |
| Bioinformatics | fixed-topology | 10gene_8k   | 1                        | 177k           | 10k             | 25k               | 4        | 694                                     | 6        | 1089                 | 63.7        | 1                      | 77        | 714         | 14                           | 2690 |      |
|                |                | 10gene_20k  | 1                        | 462k           | 10k             | 24k               | 10       | 1804                                    | 37       | 7055                 | 25.5        | 2                      | 1587      | 2412        | 68                           | 672  |      |
|                |                | 100gene_8k  | 2                        | 1300k          | 72k             | 181k              | 27       | 5079                                    | 32       | 6051                 | 83.9        | 2                      | 2109      | 5022        | 122                          | 672  |      |
|                |                | 100gene_12k | 5                        | 597k           | 21k             | 53k               | 13       | 2333                                    | 32       | 6123                 | 38.1        | 5                      | 1383      | 1861        | 55                           | 672  |      |
|                |                | 100gene_16k | 9                        | 436k           | 12k             | 29k               | 9        | 1704                                    | 32       | 6141                 | 27.7        | 9                      | 1226      | 2161        | 57                           | 672  |      |
|                |                | 100gene_20k | 12                       | 406k           | 9k              | 21k               | 9        | 1587                                    | 32       | 6118                 | 25.9        | 12                     | 1156      | 2086        | 55                           | 672  |      |
|                | state-reduced  | 10gene_8k   | 1                        | 56k            | 8               | 21k               | 2        | 218                                     | 6        | 971                  | 22.4        | 1                      | 68        | 674         | 13                           | 2690 |      |
|                |                | 10gene_20k  | 1                        | 317k           | 8               | 22k               | 7        | 1239                                    | 35       | 6625                 | 18.7        | 2                      | 1030      | 2390        | 58                           | 672  |      |
|                |                | 100gene_8k  | 1                        | 230k           | 32              | 136k              | 5        | 899                                     | 30       | 5643                 | 15.9        | 1                      | 1696      | 6098        | 131                          | 672  |      |
|                |                | 100gene_12k | 2                        | 985k           | 162             | 201k              | 21       | 3850                                    | 77       | 15k                  | 25.6        | 3                      | 4537      | 12k         | 192                          | 672  |      |
|                |                | 100gene_16k | 5                        | 1746k          | 300             | 197k              | 38       | 6822                                    | 153      | 29k                  | 23.5        | 7                      | 5273      | 13k         | 331                          | 672  |      |
|                |                | 100gene_20k | 7                        | 2592k          | 428             | 201k              | 56       | 10k                                     | 256      | 49k                  | 20.4        | 13                     | 9132      | 21k         | 502                          | 672  |      |
|                | Synth.         | deep        | 64char                   | 3              | 800k            | 84                | 4447     | 18                                      | 3128     | 20                   | 3623        | 86.3                   | 3         | 324         | 300                          | 12   | 1648 |
|                |                |             | 256char                  | 3              | 752k            | 130               | 5096     | 17                                      | 2940     | 19                   | 3300        | 89.0                   | 3         | 404         | 249                          | 12   | 1648 |
| shallow        |                | 64char      | 3                        | 779k           | 101             | 382k              | 18       | 3045                                    | 79       | 15k                  | 20.3        | 5                      | 2725      | 11k         | 224                          | 824  |      |
|                |                | 256char     | 3                        | 801k           | 193             | 317k              | 17       | 3128                                    | 66       | 12k                  | 26.0        | 3                      | 2258      | 43k         | 757                          | 824  |      |

Throughput data are shown in Table 3 for GPU and FPGA and in Table 5 for the AP. As can be seen, while able to fit even large datasets on a single device, GPU reports the lowest throughput data. In the GPU experiments, we configured the thread-block size to 256 and 32 for bioinformatics and NIDS/synthetic datasets, respectively. This is because we expected bioinformatics datasets to have larger active sets (as confirmed in Table 1). We recall that the number of thread-blocks run is equal to the product between the number of partitions and the number of input streams processed. To avoid idle SMs and ensure processing all partitions, we set the number of blocks to be at least equal to the number of SMs and of partitions. We then increased the number of blocks (and, as a consequence, of streams) until noticeable throughput improvements could no longer be observed. We make two observations. First, GPU resources are better utilized when processing a large number of input streams, leading to better throughput. Second, while the iNFAnT kernel greatly outperforms the CSR kernels on small datasets, the CSR-state kernel reports better performance on bioinformatics datasets with large  $k$ . On datasets with a large number of partitions, iNFAnT is penalized by looping through a large number of transitions that originate from inactive states.

Since large datasets require multiple FPGAs and AP boards (or multiple iterations through the same board), for FPGAs and the AP we report the traversal throughput per device. Since for most partitions the slice capacity is fully utilized, the number of FPGA devices required is equal to the number of FPGA partitions (Table 2/column 7), while the number of AP boards required is reported in Table 5/column 12. For small datasets requiring only a small portion of the device, both platforms can run multiple streams by replicating the NFA. In case of FPGA to utilize ~70% of slice capacity, we run 6 and 4 streams for l7-filter and snort534 respectively. In case of the AP, we consider that

chips can be grouped into logical cores processing streams in parallel (Section 2.2). As can be seen, on large datasets (100ups\* and synthetic) FPGAs outperform the AP up to a factor ~2.6x, while requiring 2-3x more devices than the AP.

**Preprocessing cost:** In this section, we focus on the platform-specific preprocessing time. The NFA optimization and partitioning steps, common to all platforms, take from 3 to 249 sec (smallest to largest dataset). After these two steps, we save the NFA into file. As can be seen from Table 4, the GPU preprocessing is mostly related to the parsing of the NFA partition files, and varies from 5 sec to about 4.5 min. For FPGA, synthesis and place&route account for most of the preprocessing time, and preprocessing a large partition may require up to 165 minutes (leading to several hours for the full datasets). Similar preprocessing times are observed on the AP (for example, the preprocessing time for the shallow-256-char dataset is about 12 hours). In addition, the preprocessing time increases with the transition density (deep datasets are preprocessed must faster than shallow ones), whereas the alphabet size has a lesser effect (since on the AP transition symbols are associated to STEs and stored in memory). As mentioned, the AP preprocessing time can be reduced in case of datasets with known topology (i.e., fixed topology datasets) by pre-compilation. However, finding a configuration that fully uses the AP is a trial-and-error process.

**Overall Comparison:** Figure 5 summarizes the results (note that throughput and preprocessing time are in logarithmic scale). As can be seen, FPGAs provide the best traversal throughputs (up to ~2.6x those of the AP) at the cost of significant preprocessing times (~hours); GPUs deliver modest traversal throughputs (~Mbps) but incur limited preprocessing time (~seconds-minutes) and can accommodate large datasets on a single device; Micron's AP is an intermediate choice between FPGAs and GPUs, and is most suited for applications that use datasets consisting of many small NFAs with a fixed topology.

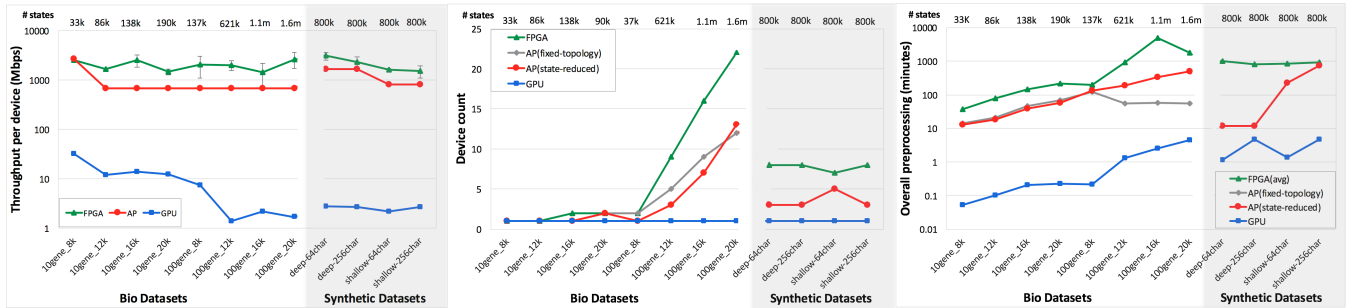


Figure 5: Traversal throughput, device utilization (in terms of number of devices) and overall preprocessing time

**Power consumption:** While the AP is not yet on the market, its design aims to a worst-case power consumption of 4W per chip [23]. Due to lack of space, here we report power data only on a medium-sized dataset (100genes\_12k). Xilinx’s Power Analyzer estimates the FPGA power consumption to be between 2.09W and 2.36W on different partitions. In contrast, GPU experiments on Texas State’s Marcher system report an average GPU power consumption of 185.11W and 105.25W on the best and worst implementations/kernel configurations, respectively.

## 6 CONCLUSION

To summarize, large datasets with more than 100-200 thousand states must be partitioned in order to be deployed on GPUs, FPGAs and Micron’s AP. While for GPUs, partitioning is required only to effectively use the GPU resources (e.g., on-chip memory), FPGAs and the AP require splitting large NFAs onto multiple devices. On these large datasets, logic-based FPGA designs can outperform the AP by a factor  $\sim 2x$ , while requiring 2-3x more devices to accommodate the dataset; GPUs underperform FPGAs by up to a factor 900x. GPUs in general deliver low performance on a single input stream, but their cumulative throughput scales up to thousands of input streams. GPUs offer the advantage of limited preprocessing time (up to a few minutes on million-state NFAs), while FPGAs and AP can take several hours to preprocess the same datasets. Precompiling the NFA can hide the AP’s preprocessing time, but this is possible only if the topology of the NFA is known a priori (e.g., Hamming or Levenshtein distance NFAs). Finding an NFA configuration that uses all 32 AP cores is a trial-and-error process that can require about an hour per experiment. Finally, due to routing constraints, AP’s SDK can keep utilization efficiency as low as 20%, while the FPGA utilization is more predictable given the NFA size.

## ACKNOWLEDGMENTS

This work has been supported by NSF awards CNS-1724934 and CCF-1421765, and by the Institute for Critical Technology and Applied Science (ICTAS: <http://www.ictas.vt.edu>).

## REFERENCES

- [1] S. Kumar *et al.*, “Algorithms to accelerate multiple regular expressions matching for deep packet inspection,” in Proc. of SIGCOMM 2006.
- [2] S. Kumar *et al.*, “Curing regular expressions matching algorithms from insomnia, amnesia, and acalculia,” in Proc. of ANCS 2007.
- [3] S. Kumar *et al.*, “Advanced algorithms for fast and scalable deep packet inspection,” in Proc. of ANCS 2016.

- [4] M. Becchi, and P. Crowley, “An improved algorithm to accelerate regular expression evaluation,” in Proc. of ANCS 2007.
- [5] M. Becchi, and P. Crowley, “A hybrid finite automaton for practical deep packet inspection,” in Proc. of CoNEXT 2007.
- [6] M. Becchi, and P. Crowley, “Extending finite automata to efficiently match Perl-compatible regular expressions,” in Proc. of CoNEXT 2008.
- [7] R. Smith *et al.*, “Deflating the big bang: fast and scalable deep packet inspection with extended finite automata,” in Proc. of SIGCOMM 2008.
- [8] A. X. Liu, and E. Torng, “An overlay automata approach to regular expression matching,” in Proc. of INFOCOM 2014.
- [9] X. Yu *et al.*, “Revisiting State Blow-up: Automatically Building Augmented-FA while Preserving Functional Equivalence,” *JSAC* 2014.
- [10] N. Cascarano *et al.*, “iNFAnt: NFA pattern matching on GPGPU devices,” *SIGCOMM Comput. Commun. Rev.*, vol. 40, no. 5, pp. 20-26, 2010.
- [11] Y. Zu *et al.*, “GPU-based NFA implementation for memory efficient high speed regular expression matching,” in Proc. of PPOPP 2012.
- [12] X. Yu, and M. Becchi, “GPU acceleration of regular expression matching for large datasets: exploring the implementation space,” in Proc. of CF 2013.
- [13] R. Sidhu, and V. K. Prasanna, “Fast Regular Expression Matching Using FPGAs,” in Proc. of FCCM 2001.
- [14] M. Becchi, and P. Crowley, “Efficient regular expression evaluation: theory to practice,” in Proc. of ANCS 2008.
- [15] Y.-H. E. Yang *et al.*, “Compact architecture for high-throughput regular expression matching on FPGA,” in Proc. of ANCS 2008.
- [16] A. Mitra *et al.*, “Compiling PCRE to FPGA for accelerating SNORT IDS,” in Proc. of ANCS 2007.
- [17] B. C. Brodie *et al.*, “A Scalable Architecture For High-Throughput Regular-Expression Pattern Matching,” in Proc. of ISCA 2006.
- [18] J. Van Lunteren *et al.*, “Designing a Programmable Wire-Speed Regular-Expression Matching Accelerator,” in Proc. of MICRO 2012.
- [19] Y. Fang *et al.*, “Fast support for unstructured data processing: the unified automata processor,” in Proc. of MICRO 2015.
- [20] M. Becchi *et al.*, “Evaluating regular expression matching engines on network and general purpose processors,” in Proc. of ANCS 2009.
- [21] P. Dlugosch *et al.*, “An Efficient and Scalable Semiconductor Architecture for Parallel Automata Processing,” *TPDS*, vol. PP, no. 99, pp. 1-1, 2014.
- [22] I. Roy, and S. Aluru, “Finding Motifs in Biological Sequences Using the Micron Automata Processor,” in Proc. of IPDPS 2014.
- [23] K. Wang *et al.*, “Association Rule Mining with the Micron Automata Processor,” in Proc. of IPDPS 2015.
- [24] K. Zhou *et al.*, “Regular expression acceleration on the micron automata processor: Brill tagging as a case study,” Proc. of Big Data 2015.
- [25] I. Roy *et al.*, “High Performance Pattern Matching Using the Automata Processor,” in Proc. of IPDPS 2016.
- [26] I. Roy *et al.*, “Algorithmic Techniques for Solving Graph Problems on the Automata Processor,” in Proc. of IPDPS 2016.
- [27] K. Wang *et al.*, “Sequential pattern mining with the Micron automata processor,” in Proc. of CF 2016.
- [28] J. E. Hopcroft, and J. Ullman, *Introduction to automata theory, languages, and computation*: Addison-Wesley, Reading, Massachusetts, 1979.
- [29] F. Yu *et al.*, “Fast and memory-efficient regular expression matching for deep packet inspection,” in Proc. of ANCS 2006.
- [30] K. Angstadt *et al.*, “RAPID Programming of Pattern-Recognition Processors,” in Proc. of ASPLOS 2016.
- [31] G. Karypis, and V. Kumar, “A Fast and High Quality Multilevel Scheme for Partitioning Irregular Graphs,” *SIAM J. Sci. Comp.*, v. 20, n. 1, pp. 359-392, 1998.
- [32] J. Wadden *et al.*, “ANMLzoo: a benchmark suite for exploring bottlenecks in automata processing engines and architectures,” in Proc. of IISWC 2016.
- [33] A. Todd *et al.*, “Parallel Gene Upstream Comparison via Multi-Level Hash Tables on GPU,” in Proc. of ICPADS 2016.
- [34] M. Becchi *et al.*, “A workload for evaluating deep packet inspection architectures,” in Proc. of IISWC 2008.