

Bypassing Space Explosion in Regular Expression Matching for Network Intrusion Detection and Prevention Systems

Jignesh Patel Alex X. Liu Eric Torng
Department of Computer Science and Engineering
Michigan State University
East Lansing, MI 48823, U.S.A.
{patelji1, alexliu, torng}@cse.msu.edu

Abstract

Network intrusion detection and prevention systems commonly use regular expression (RE) signatures to represent individual security threats. While the corresponding DFA for any one RE is typically small, the DFA that corresponds to the entire set of REs is usually too large to be constructed or deployed. To address this issue, a variety of alternative automata implementations that compress the size of the final automaton have been proposed such as XFA and D²FA. The resulting final automata are typically much smaller than the corresponding DFA. However, the previously proposed automata construction algorithms do suffer from some drawbacks. First, most employ a “Union then Minimize” framework where the automata for each RE are first joined before minimization occurs. This leads to an expensive NFA to DFA subset construction on a relatively large NFA. Second, most construct the corresponding large DFA as an intermediate step. In some cases, this DFA is so large that the final automaton cannot be constructed even though the final automaton is small enough to be deployed. In this paper, we propose a “Minimize then Union” framework for constructing compact alternative automata focusing on the D²FA. We show that we can construct an almost optimal final D²FA with small intermediate parsers. The key to our approach is a space and time efficient routine for merging two compact D²FA into a compact D²FA. In our experiments, our algorithm runs up to 302 times faster and uses 1390 times less memory than previous algorithms. For example, we are able to construct a D²FA with over 80,000,000 states using only 1GB of main memory in only 77 minutes.

1 Introduction

1.1 Background and Problem Statement

The core component of today’s network security devices such as Network Intrusion Detection and Prevention Systems is signature based deep packet inspection. The contents of every packet need to be compared against a set of signatures. Application level signature analysis can also be used for detecting peer-to-peer traffic, providing advanced QoS mechanisms. In the past, the signatures were specified as simple strings. Today, most deep packet inspection engines such as Snort [1, 24], Bro [23], TippingPoint X505 and Cisco security appliances use *regular expressions (REs)* to define the signatures. REs are used instead of simple string patterns because REs are fundamentally more expressive and thus are able to describe a wider variety of attack signatures [28]. As a result, there has been a lot of recent work on implementing high speed RE parsers for network applications.

Most RE parsers use some variant of the Deterministic Finite State Automata (DFA) representation of REs. A DFA is defined as a 5-tuple $(Q, \Sigma, \delta, q_0, A)$, where Q is a set of states, Σ is an alphabet, $\delta: Q \times \Sigma \rightarrow Q$ is the transition function, $q_0 \in Q$ is the start, and $A \subseteq Q$ is the set of accepting states. Any set of REs can be converted into an equivalent DFA with the minimum number of states [13, 14]. DFAs have the property of needing constant memory access per input symbol, and hence result in predictable and fast bandwidth. The main problem with DFAs is space explosion: a huge amount of memory is needed to store the transition function which has $|Q| \times |\Sigma|$ entries. Specifically, the number of states can be very large (state explosion), and the number of transitions per state is large ($|\Sigma|$).

To address the DFA space explosion problem, a variety of DFA variants have been proposed that require much less memory than DFAs to store. For example, there is the *De-*

layered Input DFA (D²FA) proposed by Kumar *et al.* [17]. The basic idea of D²FA is that in a typical DFA for real world RE set, given two states u and v , $\delta(u, c) = \delta(v, c)$ for many symbols $c \in \Sigma$. We can remove all the transitions for v from δ for which $\delta(u, c) = \delta(v, c)$ and make a note that v 's transitions were removed based on u 's transitions. When the D²FA is later processing input and is in state v and encounters input symbol x , if $\delta(v, x)$ is missing, the D²FA can use $\delta(u, x)$ to determine the next state. We can do the same thing for most states in the DFA, and it results in tremendous transition compression. Kumar *et al.* observe an average decrease of 97.6% in the amount of memory required to store a D²FA when compared to its corresponding DFA.

In more detail, to build a D²FA from a DFA, just do the following two steps. First, for each state $u \in Q$, pick a *deferred* state, denoted by $F(u)$. (We can have $F(u) = u$.) Second, for each state $u \in Q$ for which $F(u) \neq u$, remove all the transitions for u for which $\delta(u, x) = \delta(F(u), x)$.

When traversing the D²FA, if on current state u and current input symbol x , $\delta(u, x)$ is missing (*i.e.* has been removed), we can use $\delta(F(u), x)$ to get the next state. Of course, $\delta(F(u), x)$ might be missing too, in which case we then use $\delta(F(F(u)), x)$ to get the next state, and so on. The only restriction on selecting deferred states is that the function F cannot create a cycle other than a self-loop on the states; otherwise all states on that cycle might have their transitions on some $x \in \Sigma$ removed and there would no way of finding the next state.

Figure 1(a) shows a DFA for the REs set $\{. * a . * b c b, . * c . * b c b\}$, and Figure 1(c) shows the D²FA built from the DFA. The dashed lines represent deferred states. The DFA has $13 \times 256 = 3328$ transitions, whereas the D²FA only has 1030 actual transitions and 9 deferred transitions.

D²FA are very effective at dealing with the DFA space explosion problem. In particular, D²FA exhibit tremendous transition compression reducing the size of the DFA by a huge factor; this makes D²FA much more practical for a software implementation of RE matching than DFAs. D²FAs are also used as starting point for advanced techniques like those in [18, 20].

This leads us to the fundamental problem we address in this paper. Given as input a set of REs \mathcal{R} , build a compact D²FA as efficiently as possible that also supports frequent updates. Efficiency is important as current methods for constructing D²FA may be so expensive in both time and space that they may not be able to construct the final D²FA even if the D²FA is small enough to be deployed in networking devices that have limited computing resources. Such issues become doubly important when we consider the issue of the frequent updates (typically additions) to \mathcal{R} that occur as new security threats are identified. The limited resource networking device must be able to efficiently compute the

new D²FA. One subtle but important point about this problem is that the resulting D²FA must report which RE (or REs) from \mathcal{R} matched a given input; this applies because each RE typically corresponds to a unique security threat. Finally, while we focus on D²FA in this paper, we believe that our techniques can be generalized to other compact RE matching automata solutions [5, 8, 16, 26, 27].

1.2 Summary of Prior Art

Given the input RE set \mathcal{R} , any solution that builds a D²FA for \mathcal{R} will have to do the following two operations: (a) union the automata corresponding to each RE in \mathcal{R} and (b) minimize the automata, both in terms of the number of states and the number of edges. Previous solutions [6, 17] employ a “Union then Minimize” framework where they first build automata for each RE within \mathcal{R} , then perform union operations on these automata to arrive at one combined automaton for all the REs in \mathcal{R} , and only then minimize the resulting combined automaton. In particular, previous solutions typically perform a computationally expensive NFA to DFA subset construction followed by or composed with DFA minimization (for states) and D²FA minimization (for edges).

Consider the D²FA construction algorithm proposed by Kumar *et al.* [17]. They first apply the Union then Minimize framework to produce a DFA that corresponds to \mathcal{R} and then construct the corresponding minimum state DFA. Next, in order to maximum transition compression, they essentially solve a maximum weight spanning tree problem on the following weighted graph which they call a *Space Reduction Graph (SRG)*. The SRG has DFA states as its vertices. The SRG is a complete graph with the weight of an edge $w(u, v)$ equal to the number of common transitions between DFA states u and v . Once the spanning tree is selected, a root state is picked and all edges are directed towards the root. These directed edges give the deferred state for each state. Figure 1(b) shows the SRG built for the DFA in Figure 1(a). Using the resulting maximum weight spanning tree on the SRG, they then set the deferment state for each state eliminating redundant transitions.

Becchi and Crowley also use the Union then Minimize Framework to arrive at a minimum state DFA [6]. At this point, rather than using an SRG to set deferment states for each state, Becchi and Crowley use state levels where the level of a DFA state u is the length of the shortest string that takes the DFA from the start state to state u . Becchi and Crowley observed that if all states defer to a state that is at a lower level than itself, then the deferment function F can never produce a cycle. Furthermore, when processing any input string of length n , at most $n - 1$ deferred transitions will be processed. Thus, for each state u , among all the states at a lower level than u , Becchi and Crowley set $F(u)$

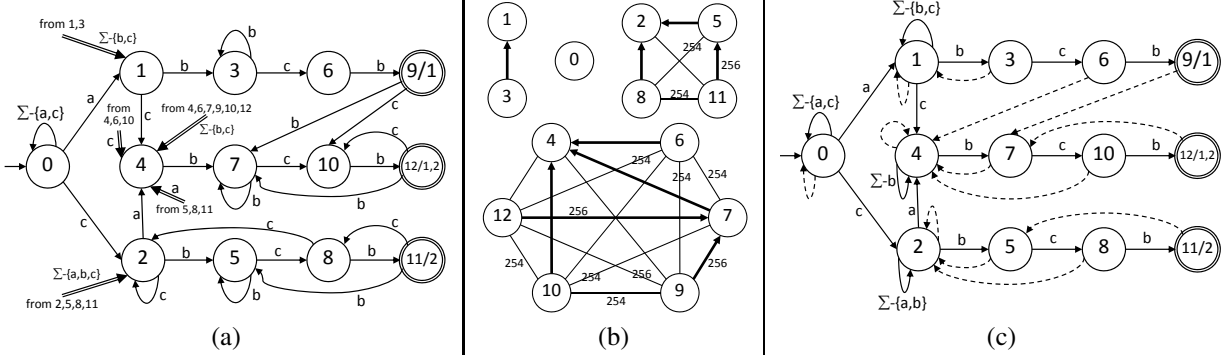


Figure 1. (a) DFA for RE set $\{ . * a . * b c b , . * c . * b c b \}$. (b) SRG for the DFA. Edges with weight ≤ 1 are not shown. Unlabeled edges have weight 255. (c) The D²FA. Dashed edges represent deferment.

to be the state which shares the most transitions with u . The resulting D²FA is typically a bit larger in size than the D²FA built using the SRG which does not have the deferring to lower level state restriction.

1.3 Limitations of Prior Art

Prior methods have three fundamental limitations. First, they follow the Union then Minimize framework which means they create large automata and only minimize them at the end. This also means they must employ the expensive NFA to DFA subset construction. Second, prior methods build the corresponding minimum state DFA before constructing the final D²FA. This is very costly in both space and time. The D²FA is typically 50 to 100 times smaller than the DFA, so even if the D²FA would fit in available memory, the intermediate DFA might be too large, making it impractical to build the D²FA. This is exacerbated in the case of the Kumar *et al.* algorithm which needs the SRG which ranges from about the size of the DFA itself to over 50 times the size of the DFA. The resulting space and time required to build the DFA and SRG impose serious limits on the D²FA that can be practically constructed. We do observe that the method proposed in [6] does not need to create the SRG. Furthermore, as the authors have noted, there is a way to go from the NFA directly to the D²FA, but implementing such an approach is still very costly in time as many transition tables need to be repeatedly recreated in order to realize these space savings. Third, none of the previous methods provide efficient algorithms for updating the D²FA when a new RE is added to \mathcal{R} .

1.4 Our Approach

To address these limitations, we propose a Minimize then Union framework. Specifically, we first minimize the

small automata corresponding to each RE from \mathcal{R} and then union the minimized automata together. A key property of our method is that our union algorithm automatically produces a minimum state D²FA for the regular expressions involved without explicit state minimization. Likewise, we create deferment states efficiently while performing the union operation using deferment information from the input D²FAs. Together, these optimizations lead to a vastly more efficient D²FA construction algorithm in both time and space.

In more detail, given \mathcal{R} , we first build a DFA and D²FA for each individual RE in \mathcal{R} . The heart of our technique is the D²FA merge algorithm that performs the union. It merges two smaller D²FAs into one larger D²FA such that the merged D²FA is equivalent to the union of REs that the D²FAs being merged were equivalent to. Starting from the the initial D²FAs for each RE, using this D²FA merge subroutine, we merge two D²FAs at a time until we are left with just one final D²FA. The initial D²FAs are each equivalent to their respective REs, so the final D²FA will be equivalent to the union of all the REs in \mathcal{R} . Figures 2(a) and 2(b) show the initial D²FAs for the RE set $\{ . * a . * b c b , . * c . * b c b \}$. The resulting D²FA from merging these two D²FAs using the D²FA merge algorithm is shown in Figure 2(c).

Advantages of our algorithm One of the main advantages of our algorithm is a dramatic increase in time and space efficiency. These efficiency gains are partly due to our use of the Minimize then Union framework instead of the Union then Minimize framework. More specifically, our improved efficiency comes about from the following four factors. First, other than for the initial DFAs that correspond to individual REs in \mathcal{R} , we build D²FA bypassing DFAs. Those initial DFAs are very small (typically < 50 states), so the memory and time required to build the ini-

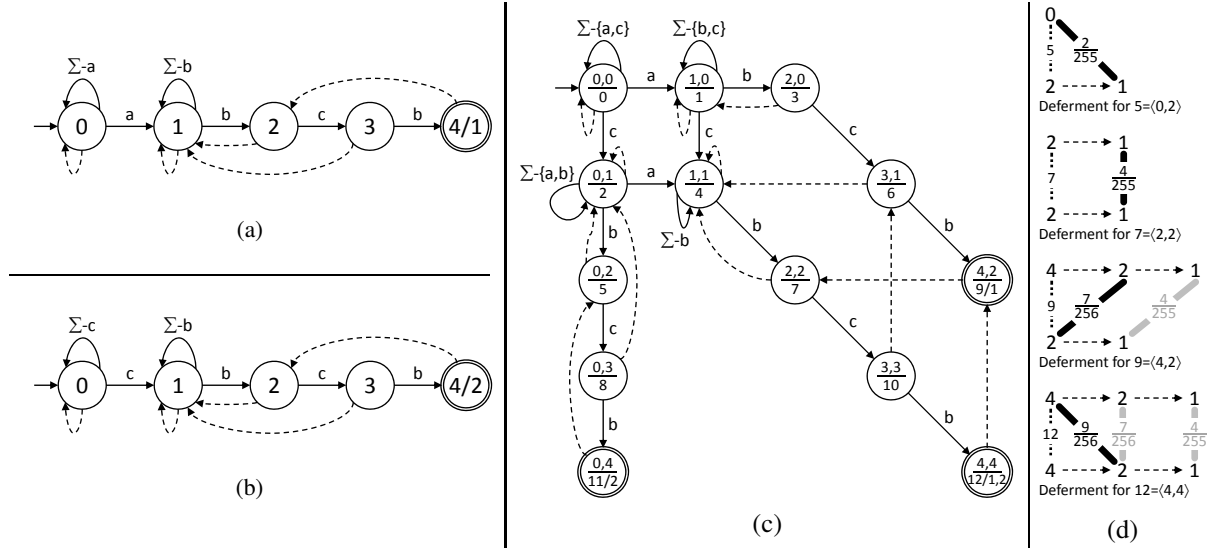


Figure 2. (a) D_1 , the D^2FA for RE $.*a.*bcb$. (b) D_2 , the D^2FA for RE $.*c.*bcb$. (c) D_3 , the merged D^2FA . (d) Illustration of setting deferment for some states in D_3 .

tial DFAs and D^2FA s is negligible. The D^2FA merge algorithm directly merges the two input D^2FA s to get the output D^2FA without creating the DFA first. Second, other than for the initial DFAs, we never have to perform the NFA to DFA subset construction. Third, other than for the initial DFAs, we never have to perform DFA state minimization. Fourth, when setting deferment states in the D^2FA merge algorithm, we use deferment information from the two input D^2FA . This typically involves performing only a constant number of comparisons per state rather than a linear in the number of states comparison per state as is required by previous techniques. All told, our algorithm has a practical time complexity of $O(n|\Sigma|)$ where n is the number of states in the final D^2FA and $|\Sigma|$ is the size of the input alphabet. In contrast, Kumar *et al.*'s algorithm [17] has a time complexity of $O(n^2(\log(n) + |\Sigma|))$ and Becchi and Crowley's algorithm [6] has a time complexity of $O(n^2|\Sigma|)$ just for setting the deferment state for each state and ignoring the cost of the NFA subset construction and DFA state minimization. See Section 5.4 for a more detailed complexity analysis.

These efficiency advantages allow us to build much larger D^2FA s than are possible with previous methods. For the synthetic RE set that we consider in Section 6, given a maximum working memory size of 1GB, we can build a D^2FA with 80, 216, 064 states with our D^2FA merge algorithm whereas the Kumar *et al.* algorithm can only build a D^2FA with 397, 312 states. Also from Section 6, our algorithm is typically 25 to 300 times faster than previous algorithms on our RE sets.

Besides being much more efficient in constructing D^2FA from scratch, our algorithm is very well suited for frequent RE updates. When an RE needs to be added to the current set, we just need to merge the D^2FA for the RE to the current D^2FA using our merge routine which is a very fast operation.

Technical Challenges For our approach to work, the main challenge is to figure out how to efficiently union two minimum state D^2FA s D_1 and D_2 so that the resulting D^2FA D_3 is also a minimum state D^2FA . There are two aspects to this challenge. First, we need to make sure that D^2FA D_3 has the minimum number of states. More specifically, suppose D_1 and D_2 are equivalent to RE sets R_1 and R_2 , respectively. We need to ensure that D_3 has the minimum number of states of any D^2FA equivalent to $R_1 \cup R_2$. We use an existing *union cross product construction* for this and prove that it results in a minimum state D^2FA for our purpose. We emphasize that this is not true in general but holds for applications where D_3 must identify which REs from $R_1 \cup R_2$ match a given input string. Many security applications meet this criteria.

Our second challenge is building the D^2FA D_3 without building the entire DFA equivalent to D_3 while ensuring that D_3 achieves significant transition compression; that is, the number of actual edges stored in D^2FA D_3 must be small. More concretely, as each state in D_3 is created, we need to immediately set a deferred state for it; otherwise, we would be storing the entire DFA. Furthermore, we need to choose a good deferment state that eliminates as many

edges as possible. We address this challenge by efficiently choosing a good deferment state for D_3 by using the deferment information from D_1 and D_2 . Typically, the algorithm only needs to compare a handful of candidate deferment states.

Key Contributions In summary, we make the following contributions: (1) We propose a novel Minimize then Union framework for constructing D^2FA for network security RE sets that we believe will generalize to other RE matching automata. Note, Smith *et al.* used the Minimize then Union Framework when constructing XFA [26, 27], though they did not prove any properties about their union algorithms. (2) To implement this framework, we propose a very efficient D^2FA merge algorithm for performing the union of two D^2FAs . (3) To prove the correctness of this D^2FA merge algorithm, we prove a fundamental property about the standard union cross product construction and minimum state DFAs when applied to network security RE sets that can be applied to other RE matching automata. We implemented our algorithms and conducted experiments on real-world and synthetic RE sets. Our experiments indicate that: (a) Our algorithm generates D^2FA with a fraction of the memory required by existing algorithms making it feasible to build D^2FA with many more states. For the real world RE sets we consider in the experimental section, our algorithm requires an average of 1390 times less memory than the algorithm proposed in [17] and 31 times less memory than the algorithm proposed in [6]. (b) Our algorithm runs much faster than existing algorithms. For the real world RE sets we consider in the experimental section, our algorithm runs an average of 302 times faster than the algorithm proposed in [17] and 28 times faster than the algorithm proposed in [6]. (c) Even with the huge space and time efficiency gains, our algorithm generates D^2FA only slightly larger than existing algorithms in the worst case.

The rest of the paper is organized as follows. We review related work in Section 2. In Section 3 we present our construction for efficiently building DFA for pattern matching RE sets. In Section 4 we present our D^2FA merge algorithm and the incrementally building approach. Section 5 presents some properties of our D^2FA merge algorithm, and some variations of it. We present experimental results and conclusions in Sections 6 and 7 respectively. For space reasons, the proofs of Theorems and Lemmas and algorithm pseudocode are given in the Appendix.

2 Related Work

Initially network intrusion detection and prevention systems used string patterns to specify attack signatures [2, 3, 29, 31–33, 36]. Sommer and Paxson [28] first proposed using REs instead of strings to specify attack sig-

natures. Today network intrusion detection and prevention systems mostly use REs for attack signatures. RE matching solutions are typically software-based or hardware-based (FPGA or ASIC).

Software-based approaches are cheap and deployable on general purpose processors, but their throughput may not be high. To achieve higher throughput, software solutions can be deployed on customized ASIC chips at the cost of low versatility and high deployment cost. To achieve deterministic throughput, software-based solutions must use DFAs, which face a space explosion problem. Specifically, there can be state explosion where the number of states increases exponentially in the number of REs, and the number of transitions per state is extremely high. To address the space explosion problem, transition compression and state minimization software-based solutions have been developed.

Transition compression schemes that minimize the number of transitions per state have mostly used one of two techniques. One is alphabet re-encoding, which exploits redundancy within a state, [6, 7, 10, 15]. The second is default transitions or deferment states which exploits redundancy among states [4, 6, 17, 18]. Kumar *et al.* [17] originally proposed the use of default transitions. Becchi and Crowley [6] proposed a more efficient way of using default transitions. Our work falls into the category of transition compression via default transitions. Our algorithms are much more efficient than those of [6, 17] and thus can be applied to much larger RE sets. For example, if we are limited to 1GB of memory to work with, we show that Kumar *et al.*'s original algorithm can only build a D^2FA with less than 400,000 states whereas our algorithm can build a D^2FA with over 80,000,000 states.

Two basic approaches have been proposed for state minimization. One is to partition the given RE set and build a DFA for each partition [35]. When inspecting packet payload, each input symbol needs to be scanned against each partition's DFA. Our work is orthogonal to this technique and can be used in combination with this technique. The second approach is to modify the automata structure and/or use extra memory to remember history and thus avoid state duplication [5, 8, 16, 26, 27]. We believe our merge technique can be adopted to work with some of these approaches. For example, Smith *et al.* also use the Minimize then Union framework when constructing XFA [26, 27]. One potential drawback with XFA is that there is no fully automated procedure to construct XFAs from a set of regular expressions. Paraphrasing Yang, Karim, Ganapathy, and Smith [34], constructing an XFA from a set of REs requires manual analysis of the REs to identify and eliminate ambiguity.

FPGA-based solutions typically exploit the parallel processing capabilities of FPGAs to implement a Nondeterministic Finite Automata (NFA) [5, 7, 11, 12, 21, 25, 30] or to implement multiple parallel DFAs [22]. TCAM based solu-

tions have been proposed for string matching in [3,9,31,36] and for REs in [20]. Our work can potentially be applied to these solutions as well.

Recently and independently, Liu *et al.* proposed to construct DFA by hierarchical merging [19]. That is, they essentially propose the Minimize then Union framework for DFA construction. They consider merging multiple DFAs at a time rather than just two. However, they do not consider D²FA, and they do not prove any properties about their merge algorithm including that it results in minimum state DFAs.

3 Pattern Matching DFAs

3.1 Pattern Matching DFA Definition

In a standard DFA, defined as a 5-tuple $(Q, \Sigma, \delta, q_0, A)$, each accepting state is equivalent to any other accepting state. However, in many pattern matching applications where we are given a set of REs \mathcal{R} , we must keep track of which REs have been matched. For example, each RE may correspond to a unique security threat that requires its own processing routine. This leads us to define Pattern Matching Deterministic Finite State Automata (PMDFA). The key difference between a PMDFA and a DFA is that for each state q in a PMDFA, we cannot simply mark it as accepting or rejecting; instead, we must record which REs from \mathcal{R} are matched when we reach q . More formally, given as input a set of REs \mathcal{R} , a PMDFA is a 5-tuple $(Q, \Sigma, \delta, q_0, M)$ where the last term M is now defined as $M: Q \rightarrow 2^{\mathcal{R}}$.

3.2 Minimum State PMDFA construction

Given a set of REs \mathcal{R} , we can build the corresponding minimum state PMDFA using the standard Union then Minimize framework: first build an NFA for the RE that corresponds to an OR of all the REs $r \in \mathcal{R}$, then convert the NFA to a DFA, and finally minimize the DFA treating accepting states as equivalent if and only if they correspond to the same set of regular expressions. This method can be very slow, mainly due to the NFA to DFA conversion, which often results in an exponential growth in the number of states. Instead, we propose a more efficient Minimize then Union framework.

Let R_1 and R_2 denote any two disjoint subsets of \mathcal{R} , and let D_1 and D_2 be their corresponding minimum state PMDFAs. We use the standard *union cross product* construction to construct a minimum state PMDFA D_3 that corresponds to $R_3 = R_1 \cup R_2$. Specifically, suppose we are given the two PMDFAs $D_1 = (Q_1, \Sigma, \delta_1, q_{01}, M_1)$ and $D_2 = (Q_2, \Sigma, \delta_2, q_{02}, M_2)$. The union cross product PMDFA of D_1 and D_2 , denoted as $UCP(D_1, D_2)$, is given by $D_3 = UCP(D_1, D_2) = (Q_3, \Sigma, \delta_3, q_{03}, M_3)$ where

$$Q_3 = Q_1 \times Q_2, \delta_3(\langle q_i, q_j \rangle, x) = \langle \delta_1(q_i, x), \delta_2(q_j, x) \rangle, q_{03} = \langle q_{01}, q_{02} \rangle, \text{ and } M_3(\langle q_i, q_j \rangle) = M_1(q_i) \cup M_2(q_j).$$

Each state in D_3 corresponds to a pair of states, one from D_1 and one from D_2 . For notational clarity, we use $\langle \text{ and } \rangle$ to enclose an ordered pair of states. Transition function δ_3 just simulates both δ_1 and δ_2 in parallel. Many states in Q_3 might not be reachable from the start state q_{03} . Thus, while constructing D_3 , we only create states that are reachable from q_{03} .

We now argue that this construction is correct. This is a standard construction, so the fact that D_3 is a PMDFA for $R_3 = R_1 \cup R_2$ is straightforward and covered in standard automata theory textbooks (*e.g.* [14]). We now show that D_3 is also a minimum state PMDFA for R_3 assuming $R_1 \cap R_2 = \emptyset$, a result that does not follow for standard DFAs.

Theorem 3.1. *Given two RE sets, R_1 and R_2 , and equivalent minimum state PMDFAs, D_1 and D_2 , the union cross product DFA $D_3 = UCP(D_1, D_2)$, with only reachable states constructed, is the minimum state PMDFA equivalent to $R_3 = R_1 \cup R_2$ if $R_1 \cap R_2 = \emptyset$.*

Proof. First since only reachable states are constructed, D_3 cannot be trivially reduced. Now assume D_3 is not minimum. That would mean there are two states in D_3 , say $\langle p_1, p_2 \rangle$ and $\langle q_1, q_2 \rangle$, that are indistinguishable. This implies that

$$\forall x \in \Sigma^*, M_3(\delta_3(\langle p_1, p_2 \rangle, x)) = M_3(\delta_3(\langle q_1, q_2 \rangle, x)).$$

Working on both sides of this equality, we get $\forall x \in \Sigma^*$,

$$\begin{aligned} M_3(\delta_3(\langle p_1, p_2 \rangle, x)) &= M_3(\langle \delta_1(p_1, x), \delta_2(p_2, x) \rangle) \\ &= M_1(\delta_1(p_1, x)) \cup M_2(\delta_2(p_2, x)) \end{aligned}$$

as well as $\forall x \in \Sigma^*$,

$$\begin{aligned} M_3(\delta_3(\langle q_1, q_2 \rangle, x)) &= M_3(\langle \delta_1(q_1, x), \delta_2(q_2, x) \rangle) \\ &= M_1(\delta_1(q_1, x)) \cup M_2(\delta_2(q_2, x)) \end{aligned}$$

This implies that

$$\begin{aligned} \forall x \in \Sigma^* M_1(\delta_1(p_1, x)) \cup M_2(\delta_2(p_2, x)) &= \\ M_1(\delta_1(q_1, x)) \cup M_2(\delta_2(q_2, x)). & \end{aligned}$$

Now since $R_1 \cap R_2 = \emptyset$, this gives us

$$\begin{aligned} \forall x \in \Sigma^*, M_1(\delta_1(p_1, x)) &= M_1(\delta_1(q_1, x)) \quad \text{and} \\ \forall x \in \Sigma^*, M_2(\delta_2(p_2, x)) &= M_2(\delta_2(q_2, x)) \end{aligned}$$

This implies that p_1 and q_1 are indistinguishable in D_1 and p_2 and q_2 are indistinguishable in D_2 , implying that both D_1 and D_2 are not minimum state PMDFAs, which is a contradiction and the result follows. \square

Our efficient construction algorithm works as follows. First, for each RE $r \in \mathcal{R}$, we build an equivalent minimum state PMDFA D for r using the standard method, resulting in a set of PMDFAs \mathcal{D} . Then we merge two PMDFAs from \mathcal{D} at a time using the above UCP construction until there is just one PMDFA left in \mathcal{D} . The merging is done in a greedy manner: in each step, the two PMDFAs with the fewest states are merged together. Note the condition $R_1 \cap R_2 \neq \emptyset$ is always satisfied in all the merges.

In our experiments, our Minimize then Union technique runs exponentially faster than the standard Union then Minimize technique because we only apply the NFA to DFA step to the NFAs that correspond to each individual regular expression rather than the composite regular expression. This makes a significant difference even when we have a relatively small number of regular expressions. For example, for our C7 RE set which contains 7 REs, the standard technique requires 385.5 seconds to build the PMDFA, but our technique builds the PMDFA in only 0.66 seconds. For the remainder of this paper, we use DFA to stand for minimum state PMDFA.

4 Efficient D²FA Construction

In this section, we first formally define what a D²FA is and then describe how we can extend the Minimize then Union technique to D²FA bypassing DFA construction.

4.1 D²FA Definition

Let $D = (Q, \Sigma, \delta, q_0, M)$ be a DFA. A corresponding D²FA D' is defined as a 6-tuple $(Q, \Sigma, \rho, q_0, M, F)$. Together, function $F: Q \rightarrow Q$ and partial function $\rho: Q \times \Sigma \rightarrow Q$ are equivalent to DFA transition function δ . Specifically, F defines a unique deferred state for each state in Q , and ρ is a partially defined transition function. We use $dom(\rho)$ to denote the domain of ρ , i.e. the values for which ρ is defined. The key property of a D²FA D' that corresponds to DFA D is that $\forall \langle q, c \rangle \in Q \times \Sigma, \langle q, c \rangle \in dom(\rho) \iff (F(q) = q \vee \delta(q, c) \neq \delta(F(q), c))$; that is for each state, ρ only has those transitions that are different from that of its deferred state in the underlying DFA. When defined, $\rho(q, c) = \delta(q, c)$. States that defer to themselves must have all their transitions defined. We only consider D²FA that correspond to minimum state DFA, though the definition applies to all DFA.

The function F defines a directed graph on the states of Q . A D²FA is well defined if and only if there are no cycles of length > 1 in this directed graph which we call a deferral forest. We use $p \rightarrow q$ to denote $F(p) = q$, i.e. p directly defers to q . We use $p \rightsquigarrow q$ to denote that there is a path from p to q in the deferral forest defined by F . We use $p \sqcap q$ to

denote the number of transitions in common between states p and q ; i.e. $p \sqcap q = |\{c \mid c \in \Sigma \wedge \delta(p, c) = \delta(q, c)\}|$.

The total transition function for a D²FA is defined as

$$\delta'(u, c) = \begin{cases} \rho(u, c) & \text{if } \langle u, c \rangle \in dom(\rho) \\ \delta'(F(u), c) & \text{else} \end{cases}$$

It is easy to see that δ' is well defined and equal to δ if the D²FA is well defined.

4.2 D²FA Merge Algorithm

The UCP construction merges two DFAs together. We extend the UCP construction to merge two D²FAs together as follows. During the UCP construction, as each new state u is created, we define $F(u)$ at that time. We then define ρ to only include transitions for u that differ from $F(u)$.

To help explain our algorithm, Figure 2 shows an example execution of the D²FA merge algorithm. Figures 2(a) and 2(b) show the D²FA for the REs $. * a . * b c b$ and $. * c . * b c b$, respectively. Figure 2(c) shows the merged D²FA for the D²FAs in figures 2(a) and 2(b). We use the following conventions when depicting a D²FA. The dashed lines correspond to the deferred state for a given state. For each state in the merged D²FA, the pair of numbers above the line refer to the states in the original D²FAs that correspond to the state in the merged D²FA. The number below the line is the state in the merged D²FA. The number(s) after the '/' in accepting states give the id(s) of the pattern(s) matched. Figure 2(d) shows how the deferred state is set for a few states in the merged D²FAs D_3 . We explain the notation in this figure as we give our algorithm description.

For each state $u \in D_3$, we set the deferred state $F(u)$ as follows. While merging D²FAs D_1 and D_2 , let state $u = \langle p_0, q_0 \rangle$ be the new state currently being added to the merged D²FA D_3 . Let $p_0 \rightarrow p_1 \rightarrow \dots \rightarrow p_l$ be the maximal deferral chain DC_1 (i.e. p_l defers to itself) in D_1 starting at p_0 , and $q_0 \rightarrow q_1 \rightarrow \dots \rightarrow q_m$ be the maximal deferral chain DC_2 in D_2 starting at q_0 . For example, in Figure 2 (d), we see the maximal deferral chains for $u = 5 = \langle 0, 2 \rangle$, $u = 7 = \langle 2, 2 \rangle$, $u = 9 = \langle 4, 2 \rangle$, and $u = 12 = \langle 4, 4 \rangle$. For $u = 9 = \langle 4, 2 \rangle$, the top row is the deferral chain of state 4 in D_1 and the bottom row is the deferral chain of state 2 in D_2 . We will choose some state $\langle p_i, q_j \rangle$ where $0 \leq i \leq l$ and $0 \leq j \leq m$ to be $F(u)$. In Figure 2(d), we represent these candidate $F(u)$ pairs with edges between the nodes of the deferral chains. For each candidate pair, the number on the top is the corresponding state number in D_3 and the number on the bottom is the number of common transitions in D_3 between that pair and state u . For example, for $u = 9 = \langle 4, 2 \rangle$, the two candidate pairs represented are state 7 ($\langle 2, 2 \rangle$) which shares 256 transitions in common with state 9 and state 4 ($\langle 1, 1 \rangle$) which shares 255 transitions in common with state 9. Note that a candidate pair

is only considered if it is reachable in D_3 . In Figure 2(d) with $u = 9 = \langle 4, 2 \rangle$, three of the candidate pairs corresponding to $\langle 4, 1 \rangle$, $\langle 2, 1 \rangle$, and $\langle 1, 2 \rangle$ are not reachable, so no edge is included for these candidate pairs. Ideally, we want i and j to be as small as possible though not both 0. For example, our best choices are typically $\langle p_0, q_1 \rangle$ or $\langle p_1, q_0 \rangle$. In the first case, $p_0 \sqcap p_1 = \langle p_0, q_0 \rangle \sqcap \langle p_1, q_0 \rangle$, and we already have $p_0 \rightarrow p_1$ in D_1 . In the second case, $q_0 \sqcap q_1 = \langle p_0, q_0 \rangle \sqcap \langle p_0, q_1 \rangle$, and we already have $q_0 \rightarrow q_1$ in D_2 . In Figure 2 (d), we set $F(u)$ to be $\langle p_0, q_1 \rangle$ for $u = 5 = \langle 0, 2 \rangle$ and $u = 12 = \langle 4, 4 \rangle$, and we use $\langle p_1, q_0 \rangle$ for $u = 9 = \langle 4, 2 \rangle$. However, it is possible that both states are not reachable from the start state in D_3 . This leads us to consider other possible $\langle p_i, q_j \rangle$. For example, in Figure 2 (d), both $\langle 2, 1 \rangle$ and $\langle 1, 2 \rangle$ are not reachable in D_3 , so we use reachable state $\langle 1, 1 \rangle$ as $F(u)$ for $u = 7 = \langle 2, 2 \rangle$.

We consider a few different algorithms for choosing $\langle p_i, q_j \rangle$. The first algorithm which we call the *first match method* is to find a pair of states $\langle p_i, q_j \rangle$ for which $\langle p_i, q_j \rangle \in Q_3$ and $i + j$ is minimum. Stated another way, we find the minimum $z \geq 1$ such that the set of states $Z = \{ \langle p_i, q_{z-i} \rangle \mid (\max(0, z-m) \leq i \leq \min(l, z)) \wedge (\langle p_i, q_{z-i} \rangle \in Q_3) \} \neq \emptyset$. From the set of states Z , of which there are at most two choices, we choose the state that has the most transitions in common with $\langle p_0, q_0 \rangle$ breaking ties arbitrarily. If Z is empty for all $z > 1$, then we just pick $\langle p_0, q_0 \rangle$, *i.e.* the deferrer pointer is not set (or the state defers to itself). The idea behind the first match method is that $\langle p_0, q_0 \rangle \sqcap \langle p_i, q_j \rangle$ decreases as $i + j$ increases. In Figure 2(d), all the selected $F(u)$ correspond to the first match method.

A second more complete algorithm for setting $F(u)$ is the *best match method* where we always consider all $(l + 1) \times (m + 1) - 1$ pairs and pick the pair that is in Q_3 and has the most transitions in common with $\langle p_0, q_0 \rangle$. The idea behind the best match method is that it is not always true that $\langle p_0, q_0 \rangle \sqcap \langle p_x, q_y \rangle \geq \langle p_0, q_0 \rangle \sqcap \langle p_{x+i}, q_{y+j} \rangle$ for $i + j > 0$. For instance we can have $p_0 \sqcap p_2 < p_0 \sqcap p_3$, which would mean $\langle p_0, q_0 \rangle \sqcap \langle p_2, q_0 \rangle < \langle p_0, q_0 \rangle \sqcap \langle p_3, q_0 \rangle$. In such cases, the first match method will not find the pair along the deferrer chains with the most transitions in common with $\langle p_0, q_0 \rangle$. In Figure 2(d), all the selected $F(u)$ also correspond to the best match method. It is difficult to create a small example where first match and best match differ.

When adding the new state u to D_3 , it is possible that some state pairs along the deferrer chains that were not in Q_3 while finding the deferred state for u will later on be added to Q_3 . This means that after all the states have been added to Q_3 , the deferrer for u can potentially be improved. Thus, after all the states have been added, for each state we again find a deferred state. If the new deferred state is better than the old one, we reset the deferrer to the new deferred state. Algorithm 1 shows the pseudocode for the D²FA merge algorithm with the first match method for

choosing a deferred state. Note that we use u and $\langle u_1, u_2 \rangle$ interchangeably to indicate a state in the merged D²FA D_3 where u is a state in Q_3 , and u_1 and u_2 are the states in Q_1 and Q_2 , respectively, that state u corresponds to.

Algorithm 1: D2FAMerge(D_1, D_2)

Input: A pair of D²FAs, $D_1 = (Q_1, \Sigma, \rho_1, q_{01}, M_1, F_1)$ and $D_2 = (Q_2, \Sigma, \rho_2, q_{02}, M_2, F_2)$, corresponding to RE sets, say R_1 and R_2 , with $R_1 \cap R_2 = \emptyset$.
Output: A D²FA corresponding to the RE set $R_1 \cup R_2$

```

1 Initialize  $D_3$  to an empty D2FA;
2 Initialize queue as an empty queue;
3 queue.push( $\langle q_{01}, q_{02} \rangle$ );
4 while queue not empty do
5    $u = \langle u_1, u_2 \rangle :=$  queue.pop();
6    $Q_3 := Q_3 \cup \{u\}$ ;
7   foreach  $c \in \Sigma$  do
8      $nxt := \langle \delta'_1(u_1, c), \delta'_2(u_2, c) \rangle$ ;
9     if  $nxt \notin Q_3 \wedge nxt \notin$  queue then queue.push( $nxt$ );
10    Add  $(u, c) \rightarrow nxt$  transition to  $\rho_3$ ;
11     $M_3(u) := M_1(u_1) \cup M_2(u_2)$ ;
12     $F_3(u) :=$  FindDefState( $u$ );
13    Remove transitions for  $u$  from  $\rho_3$  that are in common with  $F_3(u)$ ;
14 foreach  $u \in Q_3$  do
15   newDptr := FindDefState( $u$ );
16   if  $(newDptr \neq F_3(u)) \wedge (newDptr \sqcap u > F_3(u) \sqcap u)$  then
17      $F_3(u) := newDptr$ ;
18     Reset all transitions for  $u$  in  $\rho_3$  and then remove ones that are in common with  $F_3(u)$ ;
19 return  $D_3$ ;

20 FindDefState( $\langle v_1, v_2 \rangle$ )
21   Let  $\langle p_0 = v_1, p_1, \dots, p_l \rangle$  be the list of states on the deferrer chain from  $v_1$  to the root in  $D_1$ ;
22   Let  $\langle q_0 = v_2, q_1, \dots, q_m \rangle$  be the list of states on the deferrer chain from  $v_2$  to the root in  $D_2$ ;
23   for  $z = 1$  to  $(l + m)$  do
24      $S := \{ \langle p_i, q_{z-i} \rangle \mid (\max(0, z-m) \leq i \leq \min(l, z)) \wedge (\langle p_i, q_{z-i} \rangle \in Q_3) \}$ ;
25     if  $S \neq \emptyset$  then return  $\operatorname{argmax}_{v \in S} (\langle v_1, v_2 \rangle \sqcap v)$ ;
26   return  $\langle v_1, v_2 \rangle$ ;

```

4.3 Direct D²FA construction for RE set

Similar to efficient DFA construction, we first build the D²FA for each RE in \mathcal{R} using the method described in [20]. We then merge the D²FAs together using a balanced binary tree structure to minimize the worst-case number of merges that any RE experiences.

5 D²FA Merge Algorithm Properties

5.1 Proof of Correctness

The D²FA merge algorithm exactly follows the UCP construction to create the states. So the correctness of the underlying DFA follows from the the correctness of the UCP construction.

Theorem 5.1 shows that the merged D²FA is also well defined (no cycles in deferment forest).

Lemma 5.1. *In the D²FA $D_3 = \text{D2FAMerge}(D_1, D_2)$, $\langle u_1, u_2 \rangle \twoheadrightarrow \langle v_1, v_2 \rangle \Rightarrow u_1 \twoheadrightarrow v_1 \wedge u_2 \twoheadrightarrow v_2$.*

Proof. If $\langle u_1, u_2 \rangle = \langle v_1, v_2 \rangle$ then the lemma is trivially true. Otherwise, let $\langle u_1, u_2 \rangle \rightarrow \langle w_1, w_2 \rangle \twoheadrightarrow \langle v_1, v_2 \rangle$ be the deferment chain in D_3 . When selecting the deferred state for $\langle u_1, u_2 \rangle$, D2FA Merge always choose a state that corresponds to a pair of states along deferment chains for u_1 and u_2 in D_1 and D_2 , respectively. Therefore, we have that $\langle u_1, u_2 \rangle \rightarrow \langle w_1, w_2 \rangle \Rightarrow u_1 \twoheadrightarrow w_1 \wedge u_2 \twoheadrightarrow w_2$. By induction on the length of the deferment chain and the fact that the \twoheadrightarrow relation is transitive, we get our result. \square

Theorem 5.1. *If D²FAs D_1 and D_2 are well defined, then the D²FA $D_3 = \text{D2FAMerge}(D_1, D_2)$ is also well defined.*

Proof. Since D_1 and D_2 are well defined, there are no cycles in their deferment forests. Now assume that D_3 is not well defined, *i.e.* there is a cycle in its deferment forest. Let $\langle u_1, u_2 \rangle$ and $\langle v_1, v_2 \rangle$ be two distinct states on the cycle. Then, we have that

$$\langle u_1, u_2 \rangle \twoheadrightarrow \langle v_1, v_2 \rangle \wedge \langle v_1, v_2 \rangle \twoheadrightarrow \langle u_1, u_2 \rangle$$

Using Lemma 5.1 we get

$$(u_1 \twoheadrightarrow v_1 \wedge u_2 \twoheadrightarrow v_2) \wedge (v_1 \twoheadrightarrow u_1 \wedge v_2 \twoheadrightarrow u_2)$$

i.e. $(u_1 \twoheadrightarrow v_1 \wedge v_1 \twoheadrightarrow u_1) \wedge (u_2 \twoheadrightarrow v_2 \wedge v_2 \twoheadrightarrow u_2)$

Since $\langle u_1, u_2 \rangle \neq \langle v_1, v_2 \rangle$, we have $u_1 \neq v_1 \vee u_2 \neq v_2$ which implies that at least one of D_1 or D_2 has a cycle in their deferment forest which is a contradiction. \square

5.2 Limiting Deferment Depth

Since no input is consumed while traversing a deferred transition, in the worst case, the number of lookups needed to process one input character is given by the depth of the deferment forest. As previously proposed, we can guarantee a worst case performance by limiting the depth of the deferment forest.

For a state u_1 of a D²FA D_1 , the *deferment depth* of u_1 , denoted as $\psi(u_1)$, is the length of the maximal deferment chain in D_1 from u_1 to the root. $\Psi(D_1) = \max_{v \in Q_1} \psi(v)$ denotes the deferment depth of D_1 (*i.e.* the depth of the deferment forest in D_1).

Lemma 5.2. *In the D²FA $D_3 = \text{D2FAMerge}(D_1, D_2)$, $\forall \langle u_1, u_2 \rangle \in Q_3$, $\psi(\langle u_1, u_2 \rangle) \leq \psi(u_1) + \psi(u_2)$.*

Proof. Let $\psi(\langle u_1, u_2 \rangle) = d$. If $\psi(\langle u_1, u_2 \rangle) = 0$, then $\langle u_1, u_2 \rangle$ is a root and the lemma is trivially true. So, we consider $d \geq 1$ and assume the lemma is true for all states with $\psi < d$. Let $\langle u_1, u_2 \rangle \rightarrow \langle w_1, w_2 \rangle \twoheadrightarrow \langle v_1, v_2 \rangle$ be the deferment chain in D_3 . Using the inductive hypothesis, we have

$$\psi(\langle w_1, w_2 \rangle) \leq \psi(w_1) + \psi(w_2)$$

Given $\langle u_1, u_2 \rangle \neq \langle w_1, w_2 \rangle$, we assume without loss of generality that $u_1 \neq w_1$. Using Lemma 5.1 we get that $u_1 \twoheadrightarrow w_1$. Therefore $\psi(w_1) \leq \psi(u_1) - 1$. Combining the above, we get $\psi(\langle u_1, u_2 \rangle) = \psi(\langle w_1, w_2 \rangle) + 1 \leq \psi(w_1) + \psi(w_2) + 1 \leq (\psi(u_1) - 1) + \psi(u_2) + 1 \leq \psi(u_1) + \psi(u_2)$. \square

Lemma 5.2 directly gives us the following Theorem.

Theorem 5.2. *If $D_3 = \text{D2FAMerge}(D_1, D_2)$, then $\Psi(D_3) \leq \Psi(D_1) + \Psi(D_2)$.*

For an RE set \mathcal{R} , if the initial D²FAs have $\Psi = d$, in the worst case, the final merged D²FA corresponding to \mathcal{R} can have $\Psi = d \times |\mathcal{R}|$. Although Theorem 5.2 gives the value of Ψ in the worst case, in practical cases, $\Psi(D_3)$ is very close to $\max(\Psi(D_1), \Psi(D_2))$. Thus the deferment depth of the final merged D²FA is usually not much higher than d .

Let Ω denote the desired upper bound on Ψ . To guarantee $\Psi(D_3) \leq \Omega$, we modify the `FindDefState` subroutine in Algorithm 1 as follows: When selecting candidate pairs for the deferred state, we only consider states with $\psi < \Omega$. Specifically, we replace line 24 with the following

$$S := \{ \langle p_i, q_{z-i} \rangle \mid (\max(0, z - m) \leq i \leq \min(l, z)) \wedge \langle p_i, q_{z-i} \rangle \in Q_3 \wedge (\psi(\langle p_i, q_{z-i} \rangle) < \Omega) \}$$

When we do the second pass (lines 14-19), we may increase the deferment depth of nodes that defer to nodes that we readjust. We record the affected nodes and then do a third pass to reset their deferment states so that the maximum depth bound is satisfied. In practice, this happens very rarely.

When constructing a D²FA with a given bound Ω , we first build D²FAs without this bound. We only apply the bound Ω when performing the final merge of two D²FAs to create the final D²FA.

5.3 Deferment to a Lower Level

In [6], the authors propose a technique to guarantee an amortized cost of 2 lookups per input character without limiting the depth of the deferment tree. They achieve this by having states only defer to lower level states where the level of any state u in a DFA (or D²FA), denoted $level(u)$, is defined as the length of the shortest string that ends in that

state (from the start state). More formally, they ensure that for all states u , $level(u) > level(F(u))$ if $u \neq F(u)$. We call this property the *back-pointer* property. If the back-pointer property holds, then every deferred transition taken decreases the level of the current state by at least 1. Since a regular transition on an input character can only increase the level of the current state by at most 1, there have to be fewer deferred transitions taken on the entire input string than regular transitions. This gives an amortized cost of at most 2 transitions taken per input character.

In order to guarantee the D²FA D_3 has the back-pointer property, we perform a similar modification to the `FindDefState` subroutine in Algorithm 1 as we performed when we wanted to limit the maximum deferment depth. When selecting candidate pairs for the deferred state, we only consider states with a lower level. Specifically, we replace line 24 with the following:

$$S := \{ \langle p_i, q_{z-i} \rangle \mid (\max(0, z-m) \leq i \leq \min(l, z)) \wedge (\langle p_i, q_{z-i} \rangle \in Q_3) \wedge (level(\langle v_1, v_2 \rangle) > level(\langle p_i, q_{z-i} \rangle)) \}$$

For states for which no candidate pairs are found, we just search through all states in Q_3 that are at a lower level for the deferred state. In practice, this search through all the states needs to be done for very few states because if D²FAs D_1 and D_2 have the back-pointer property, then almost all the states in D²FAs D_3 have the back-pointer property. As with limiting maximum deferment depth, we only apply this restriction when performing the final merge of two D²FAs to create the final D²FA.

5.4 Algorithmic Complexity

The time complexity of the original D²FA algorithm proposed in [17] is $O(n^2(\log(n) + |\Sigma|))$. The SRG has $O(n^2)$ edges, and $O(|\Sigma|)$ time is required to add each edge to the SRG and $O(\log(n))$ time is required to process each edge in the SRG during the maximum spanning tree routine. The time complexity of the D²FA algorithm proposed in [6] is $O(n^2|\Sigma|)$. Each state is compared with $O(n)$ other states, and each comparison requires $O(|\Sigma|)$ time.

The time complexity of our new D2FAMerge algorithm to merge two D²FAs is $O(n\Psi_1\Psi_2|\Sigma|)$, where n is the number of states in the merged D²FA, and Ψ_1 and Ψ_2 are the maximum deferment depths of the two input D²FAs. When setting the deferment for any state $u = \langle u_1, u_2 \rangle$, in the worst case the algorithm compares $\langle u_1, u_2 \rangle$ with all the pairs along the deferment chains of u_1 and u_2 , which are at most Ψ_1 and Ψ_2 in length, respectively. Each comparison requires $O(|\Sigma|)$ time. In practice, the time complexity is $O(n|\Sigma|)$ as each state needs to be compared with very few states for the following three reasons. First, the maximum deferment depth Ψ is usually very small. The largest value of Ψ among our 8 primary RE sets in Section 6 is 7.

Second, the length of the deferment chains for most states is much smaller than Ψ . The largest value of average deferment depth $\bar{\psi}$ among our 8 RE sets is 2.54. Finally, many of the state pairs along the deferment chains are not reachable in the merged D²FA. Among our 8 RE sets, the largest value of the average number of comparisons needed is 1.47.

When merging all the D²FAs together for an RE set \mathcal{R} , the total time required in the worst case would be $O(n\Psi_1\Psi_2|\Sigma| \log(|\mathcal{R}|))$. The worst case would happen when the RE set contains strings and there is no state explosion. In this case, each merged D²FA would have a number of states roughly equal to the sum of the sizes of the D²FAs being merged. When there is state explosion, the last D²FA merge would be the dominating factor, and the total time would just be $O(n\Psi_1\Psi_2|\Sigma|)$.

When modifying the D2FAMerge algorithm to maintain back-pointers, the worst case time would be $O(n^2|\Sigma|)$ because we would have to compare each state with $O(n)$ other states if none of the candidate pairs are found at a lower level than the state. In practice, this search needs to be done for very few states, typically less than 1%.

6 Experimental Results

In this section, we evaluate the effectiveness of our algorithm (D²FAMERGE) on real-world and synthetic RE sets. We compare D²FAMERGE with the original D²FA algorithm proposed in [17] (ORIGINAL) that optimizes transition compression and the D²FA algorithm proposed in [6] (BACKPTR) that enforces the back-pointer property described in Section 5.3.

6.1 Methodology

6.1.1 Data Sets

Our main results are based on eight real RE sets, four proprietary RE sets C7, C8, C10, and C613 from a large networking vendor and four public RE sets Bro217, Snort 24, Snort31, and Snort 34, that we partition into three groups, STRING, WILDCARD, and SNORT, based upon their RE composition. For each RE set, the number indicates the number of REs in the RE set. The STRING RE sets, C613 and Bro217, contain mostly string matching REs. The WILDCARD RE sets, C7, C8 and C10, contain mostly REs with multiple wildcard closures ‘.*’. The SNORT RE sets, Snort24, Snort31, and Snort34, contain a more diverse set of REs, roughly 40% of which have wildcard closures. To test scalability, we use Scale, a synthetic RE set consisting of 26 REs of the form $/.*c_u0123456.*c_l789!#\%&/$, where c_u and c_l are the 26 uppercase and lowercase alphabet letters. Even though all the REs are nearly identical differing only in the character after the two ‘.*’s, we still get

the full multiplicative effect where the number of states in the corresponding minimum state DFA roughly doubles for every RE added.

6.1.2 Metrics

We use the following metrics to evaluate the algorithms. First, we measure the resulting D²FA size (# transitions) to assess transition compression performance. Our D²FAMERGE algorithm typically performs almost as well as the other algorithms even though it builds up the D²FA incrementally rather than compressing the final minimum state DFA. Second, we measure the the maximum deferment depth (Ψ) and average deferment depth ($\bar{\psi}$) in the D²FA to assess how quickly the resulting D²FA can be used to perform regular expression matching. Smaller Ψ and $\bar{\psi}$ mean that fewer deferment transitions that process no input characters need to be traversed when processing an input string. Our D²FAMERGE significantly outperforms the other algorithms. Finally, we measure the space and time required by the algorithm to build the final automaton. Again, our D²FAMERGE significantly outperforms the other algorithms. When comparing the performance of D²FAMERGE with another algorithm A on a given RE or RE set, we define the following quantities to compare them: transition increase is (D²FAMERGE D²FA size - A D²FA size) divided by A D²FA size, transition decrease is (A D²FA size - D²FAMERGE D²FA size) divided by A D²FA size, average (maximum) deferment depth ratio is A average (maximum) deferment depth divided by D²FAMERGE average (maximum) deferment depth, space ratio is A space divided by D²FAMERGE space, and time ratio is A build time divided by D²FAMERGE build time.

6.1.3 Measuring Space

When measuring the required space for an algorithm, we measure the maximum amount of memory required at any point in time during the construction and then final storage of the automaton. This is a difficult quantity to measure exactly; we approximate this required space for each of the algorithms as follows. For D²FAMERGE, the dominant data structure is the D²FA. For a D²FA, the transitions for each state can be stored as pairs of input character and next state id, so the memory required to store a D²FA is calculated as $= (\#transitions) \times 5$ bytes. However, the maximum amount of memory required while running D²FAMERGE may be higher than the final D²FA size because of the following two reasons. First, when merging two D²FAs, we need to maintain the two input D²FAs as well as the output D²FA. Second, we may create an intermediate output D²FA that has more transitions than needed; these extra transitions will be eliminated once all D²FA states are added. We keep

track of the worst case required space for our algorithm during D²FA construction. This typically occurs when merging the final two intermediate D²FA to form the final D²FA.

For ORIGINAL, we measure the space required by the minimized DFA and the SRG. For the DFA, the transitions for each state can be stored as an array of size Σ with each array entry requiring four bytes to hold the next state id. For the SRG, each edge requires 17 bytes as observed in [6]. This leads to a required memory for building the D²FA of $= |Q| \times |\Sigma| \times 4 + (\#edges\ in\ SRG) \times 17$ bytes.

For BACKPTR, we consider two variants. The first variant builds the minimized DFA directly from the NFA and then sets the deferment for each state. For this variant, no SRG is needed, so the space required is the space needed for the minimized DFA which is $|Q| \times |\Sigma| \times 4$ bytes. The second variant goes directly from the NFA to the final D²FA; this variant uses less space but is much slower as it stores incomplete transition tables for most states. Thus, when computing the deferment state for a new state, the algorithm must recreate the complete transition tables for each state to determine which has the most common transitions with the new state. For this variant, we assume the only space required is the space to store the final D²FA which is $= (\#transitions) \times 5$ bytes even though more memory is definitely needed at various points during the computation. We also note that both implementations must perform the NFA to DFA subset construction on a large NFA which means even the faster variant runs much more slowly than D²FAMERGE.

6.1.4 Correctness

We tested correctness of our algorithms by verifying the final D²FA is equivalent to the corresponding DFA. Note, we can only do this check for our RE sets where we were able to compute the corresponding DFA. Thus, we only verified correctness of the final D²FA for our eight real RE sets and the smaller Scale RE sets.

6.2 D²FAMERGE versus ORIGINAL

We first compare D²FAMERGE with ORIGINAL that optimizes transition compression when both algorithms have unlimited maximum deferment depth. These results are shown in Table 1 for our 8 primary RE sets. Table 2 summarizes these results by RE group. We make the following observations.

(1) *D²FAMERGE uses much less space than ORIGINAL.* On average, D²FAMERGE uses 1390 times less memory than ORIGINAL to build the resulting D²FA. This difference is most extreme when the SRG is large, which is true for the two STRING RE sets and Snort24 and Snort34. For these RE sets, D²FAMERGE uses between 347 and 3613

RE set	# States	ORIGINAL					D ² FAMERGE				
		# Trans	Def. depth		RAM (MB)	Time (s)	# Trans	Def. depth		RAM (MB)	Time (s)
			Avg.	Max.				Avg.	Max.		
Bro217	6533	9816	3.90	8	179.3	542.1	12325	2.16	5	0.10	6.1
C613	11308	21633	4.38	11	1042.7	1892.5	34991	2.54	7	0.29	17.1
C7	24750	205633	16.38	27	47.4	1274.7	208564	1.14	3	1.07	2.1
C8	3108	23209	8.60	14	4.9	36.3	24604	1.14	2	0.14	0.5
C10	14868	96793	16.39	26	25.5	505.6	99124	1.17	3	0.53	1.6
Snort24	13886	38485	9.67	18	861.2	1856.2	44883	1.56	4	0.35	0.5
Snort31	20068	70701	9.17	16	298.5	1086.8	94339	1.97	6	0.86	5.3
Snort34	13825	40199	10.95	18	795.2	1911.3	45642	1.38	5	0.28	3.7

Table 1. The D²FA size, D²FA average $\bar{\psi}$ and maximum Ψ deferment depths, space estimate and time required to build the D²FA for ORIGINAL and D²FAMERGE on our eight primary RE sets.

RE set group	Trans increase	Def. depth ratio		Space ratio	Time ratio
		Avg.	Max.		
All	20.1%	7.3	4.8	1390.0	301.6
STRING	44.0%	1.8	1.6	2672.8	99.5
WILDCARD	3.0%	12.0	8.2	42.7	338.2
SNORT	21.3%	6.3	3.6	1882.1	399.7

Table 2. Average values of transition increase, deferment depth ratios, space ratios, and time ratios for D²FAMERGE compared with ORIGINAL for our RE set groups.

times less memory than ORIGINAL. For the RE sets with relatively small SRGs such as those in the WILDCARD and Snort31, D²FAMERGE uses between 35 and 49 times less space than ORIGINAL.

(2) *D²FAMERGE is much faster than ORIGINAL.* On average, D²FAMERGE builds the D²FA 300 times faster than ORIGINAL. This time difference is maximized when the deferment chains are shortest. For example, D²FAMERGE only requires an average of 0.12 msec and 0.19 msec per state for the WILDCARD and SNORT RE sets, respectively, so D²FAMERGE is, on average, 338 and 399 times faster than ORIGINAL for these RE sets, respectively. For the STRING RE sets, the deferment chains are longer, so D²FAMERGE requires an average of 1.23 msec per state, and is, on average, 100 times faster than ORIGINAL.

(3) *D²FAMERGE produces D²FA with much smaller average and maximum deferment depths than ORIGINAL.* On average, D²FAMERGE produces D²FA that have average deferment depths that are 7.3 times smaller than ORIGINAL and maximum deferment depths that are 4.8 times smaller than ORIGINAL. In particular, the average deferment depth for D²FAMERGE is less than 2 for all but the two STRING RE sets, where the average deferment depths are 2.16 and 2.54. Thus, the expected number of deferment transitions to be traversed when processing a length n string is less than n . One reason D²FAMERGE works so well is that it eliminates low weight edges from the SRG so that the deferment forest has many shallow deferment trees instead of one deep tree. This is particularly effective for the

WILDCARD RE sets and, to a lesser extent, the SNORT RE sets. For the STRING RE sets, the SRG is fairly dense, so D²FAMERGE has a smaller advantage relative to ORIGINAL.

(4) *D²FAMERGE produces D²FA with only slightly more transitions than ORIGINAL, particularly on the RE sets that need transition compression the most.* On average, D²FAMERGE produces D²FA with roughly 20% more transitions than ORIGINAL does. D²FAMERGE works best when state explosion from wildcard closures creates DFA composed of many similar repeating substructures. This is precisely when transition compression is most needed. For example, for the WILDCARD RE sets that experience the greatest state explosion, D²FAMERGE only has 3% more transitions than ORIGINAL. On the other hand, for the STRING RE sets, D²FAMERGE has, on average, 44% more transitions. For this group, ORIGINAL needed to build a very large SRG and thus used much more space and time to achieve the improved transition compression. Furthermore, transition compression is typically not needed for such RE sets as all string matching REs can be placed into a single group and the resulting DFA can be built.

In summary, D²FAMERGE achieves its best performance relative to ORIGINAL on the WILDCARD RE sets (except for space used for construction of the D²FA) and its worst performance relative to ORIGINAL on the STRING RE sets (except for space used to construct the D²FA). This is desirable as the space and time efficient D²FAMERGE is most needed on RE sets like those in the WILDCARD because those RE sets experience the greatest state explosion.

6.3 D²FAMERGE versus ORIGINAL with Bounded Maximum Deferment Depth

We now compare D²FAMERGE and ORIGINAL when they impose a maximum deferment depth bound Ω of 1, 2, and 4. Because time and space do not change significantly, we focus only on number of transitions and average defer-

RE set	ORIGINAL						D ² FAMERGE					
	# Trans			Avg. def. depth			# Trans			Avg. def. depth		
	$\Omega = 1$	$\Omega = 2$	$\Omega = 4$	$\Omega=1$	$\Omega=2$	$\Omega=4$	$\Omega = 1$	$\Omega = 2$	$\Omega = 4$	$\Omega=1$	$\Omega=2$	$\Omega=4$
B217	690999	260076	38898	0.59	1.17	2.21	50281	15633	12345	1.00	1.82	2.16
C613	1197312	473541	85797	0.59	1.14	2.19	155308	50891	35962	1.00	1.90	2.50
C7	2013157	529448	206682	0.69	1.27	2.28	216964	209068	208564	0.97	1.13	1.14
C8	205036	31786	23261	0.75	1.32	2.48	25360	24604	24604	0.98	1.14	1.14
C10	1103830	278839	97074	0.71	1.27	2.34	102826	99130	99124	0.98	1.17	1.17
S24	1361481	514277	79911	0.62	1.18	2.15	72054	47558	44883	1.00	1.47	1.56
S31	2179545	1000435	194294	0.58	1.08	2.21	197043	119087	97039	1.00	1.52	1.95
S34	1351430	530188	70070	0.62	1.15	2.19	59936	47106	45732	1.00	1.34	1.38

Table 3. The D²FA size and D²FA average $\bar{\psi}$ deferment depth for ORIGINAL and D²FAMERGE on our eight primary RE sets given maximum deferment depth bounds of 1, 2 and 4.

ment depth. These results are shown in Table 3. Note that for these data sets, the resulting maximum depth Ψ typically is identical to the maximum depth bound Ω ; the only exception is for D²FAMERGE and $\Omega = 4$; thus we omit the maximum deferment depth from Table 3. Table 4 summarizes the results by RE group highlighting how much better or worse D²FAMERGE does than ORIGINAL on the two metrics of number of transitions and average deferment depth $\bar{\psi}$.

Overall, D²FAMERGE performs very well when presented a bound Ω . In particular, the average increase in the number of transitions for D²FAMERGE with Ω equal to 1, 2 and 4, is only 108%, 14% and 0.7% respectively, compared to D²FAMERGE with unbounded maximum deferment depth. Stated another way, when D²FAMERGE is required to have a maximum deferment depth of 1, this only results in slightly more than twice the number of transitions in the resulting D²FA. The corresponding values for ORIGINAL are 3095%, 1099% and 119%.

RE set group	$\Omega = 1$		$\Omega = 2$		$\Omega = 4$	
	Trans decr.	Avg. def. depth ratio	Trans decr.	Avg. def. depth ratio	Trans decr.	Avg. dptr len ratio
All	91.3%	0.7	75.1%	0.9	30.8%	1.5
STRING	90.0%	0.6	91.5%	0.6	63.0%	1.0
WILDCARD	89.3%	0.7	49.3%	1.1	-3.0%	2.1
SNORT	94.0%	0.6	90.0%	0.8	43.0%	1.4

Table 4. Average values of transition decrease and average deferment depth ratios for D²FAMERGE compared with ORIGINAL for our RE set groups given maximum deferment depth bounds of 1, 2 and 4.

These results can be partially explained by examining the average deferment depth data. Unlike in the unbounded maximum deferment depth scenario, here we see that D²FAMERGE has a larger average deferment depth $\bar{\psi}$ than ORIGINAL except for the WILDCARD when Ω is 1 or 2. What this means is that D²FAMERGE has more states that defer to at least one other state than ORIGINAL does. This leads to the lower number of transitions in the final D²FA. Overall, for $\Omega = 1$, D²FAMERGE produces D²FA with roughly 90% fewer transitions than ORIGINAL for all

RE set groups. For $\Omega = 2$, D²FAMERGE produces D²FA with roughly 50% fewer transitions than ORIGINAL for the WILDCARD RE sets and roughly 90% fewer transitions than ORIGINAL for the other RE sets.

6.4 D²FAMERGE versus BACKPTR

We now compare D²FAMERGE with BACKPTR which enforces the back-pointer property described in Section 5.3. We adapt D²FAMERGE to also enforce this back-pointer property. The results for all our metrics are shown in Table 5 for our 8 primary RE sets. We consider the two variants of BACKPTR described in Section 6.1.3, one which constructs the minimum state DFA corresponding to the given NFA and one which bypasses the minimum state DFA and goes directly to the D²FA from the given NFA. We note the second variant appears to use less space than D²FAMERGE. This is partially true since BACKPTR creates a smaller D²FA than D²FAMERGE. However, we underestimate the actual space used by this BACKPTR variant by simply assuming its required space is the final D²FA size. We ignore, for instance, the space required to store intermediate complete tables or to perform the NFA to DFA subset construction. Table 6 summarizes these results by RE group displaying ratios for many of our metrics that highlight how much better or worse D²FAMERGE does than BACKPTR.

Similar to D²FAMERGE versus ORIGINAL, we find that D²FAMERGE with the backpointer property performs well when compared with both variants of BACKPTR. Specifically, with an average increase in the number of transitions of roughly 23%, D²FAMERGE runs on average 27 times faster than the fast variant of BACKPTR and 120 times faster than the slow variant of BACKPTR. For space, D²FAMERGE uses on average almost 31 times less space than the first variant of BACKPTR and on average roughly 42% more space than the second variant of BACKPTR. Furthermore, D²FAMERGE creates D²FA with average deferment depth 2.9 times smaller than BACKPTR and maximum deferment depth 1.9 times smaller than BACKPTR. As was the case with ORIGINAL, D²FAMERGE achieves its best performance relative to BACKPTR on the WILD-

RE set	BACKPTR							D ² FAMERGE with back-pointer				
	# Trans	Def. depth		RAM (MB)	Time (s)	RAM2 (MB)	Time2 (s)	# Trans	Def. depth		RAM (MB)	Time (s)
		Avg.	Max.						Avg.	Max.		
B217	11247	2.61	6	6.4	238.83	0.05	404.54	13998	2.31	6	0.11	9.68
C613	26222	2.50	5	11.0	118.25	0.13	1087.96	43617	2.13	5	0.32	17.74
C7	217812	5.94	13	24.2	601.2	1.04	2857.76	219684	1.15	4	1.12	6.51
C8	34636	2.44	8	3.0	25.63	0.17	41.43	35476	1.20	4	0.19	1.02
C10	157139	2.13	7	14.5	206.28	0.75	664.34	158232	1.19	4	0.81	14.42
S24	46005	8.74	17	13.6	230.58	0.22	1713.83	58273	1.62	8	0.41	47.77
S31	82809	2.87	8	19.6	269.7	0.39	1472.86	126508	1.66	6	0.90	10.29
S34	46046	7.05	14	13.5	228.44	0.22	1381.64	52057	1.41	5	0.31	8.66

Table 5. The D²FA size, D²FA average $\bar{\psi}$ and maximum Ψ deferment depths, space estimate and time required to build the D²FA for both variants of BACKPTR and D²FAMERGE with the back-pointer property on our eight primary RE sets.

CARD RE sets and its worst performance relative to BACKPTR on the STRING RE sets. This is desirable as the space and time efficient D²FAMERGE is most needed on RE sets like those in the WILDCARD because those RE sets experience the greatest state explosion.

RE set group	Trans increase	Def. depth ratio		Space ratio	Time ratio	Space2 ratio	Time2 ratio
		Avg.	Max.				
All	23.4%	2.9	1.9	30.9	27.6	0.7	120.9
STRING	45.0%	1.2	1.0	46.6	15.7	0.4	51.6
WILDCARD	1.3%	3.0	2.3	18.4	43.9	0.9	175.2
SNORT	31.0%	4.0	2.1	32.8	19.1	0.6	112.9

Table 6. Average values of transition increase, deferment depth ratios, space ratios, and time ratios for D²FAMERGE compared with both variants of BACKPTR for RE set groups.

6.5 Scalability results

Finally, we assess the improved scalability of D²FAMERGE relative to ORIGINAL using the Scale RE set assuming that we have a maximum memory size of 1 GB. For both ORIGINAL and D²FAMERGE, we add one RE at a time from Scale until the space estimate to build the D²FA goes over the 1GB limit. For ORIGINAL, we are able only able to add 12 REs; the final D²FA has 397,312 states and requires over 71 hours to compute. As explained earlier, we include the SRG edges in the RAM size estimate. If we exclude the SRG edges and only include the DFA size in the RAM size estimate, we would only be able to add one more RE before we reach the 1GB limit. For D²FAMERGE, we are able to add 19 REs; the final D²FA has 80,216,064 states and requires only 77 minutes to compute. This data set highlights the quadratic versus linear running time of ORIGINAL and D²FAMERGE, respectively. Figure 3 shows how the space and time requirements grow for ORIGINAL and D²FAMERGE as RE's from Scale are added one by one until 19 have been added.

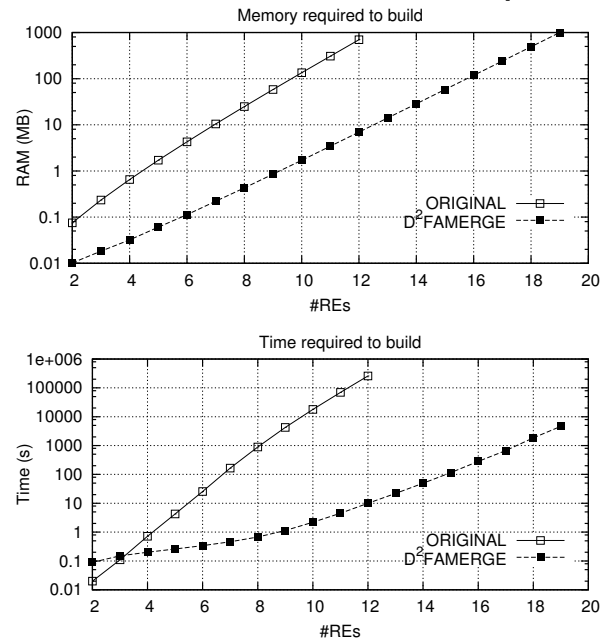


Figure 3. Memory and time required to build D²FA versus number of Scale REs used for ORIGINAL's D²FA and D²FAMERGE's D²FA.

7 Conclusions

In this paper, we propose a novel Minimize then Union framework for constructing D²FAs using D²FA merging. Our approach requires a fraction of memory and time compared to current algorithms. This allows us to build much larger D²FAs than was possible with previous techniques. Our algorithm naturally supports frequent RE set updates. We conducted experiments on real-world and synthetic RE sets that verify our claims. For example, our algorithm requires an average of 1400 times less memory and 300 times less time than the original D²FA construction algorithm of Kumar *et al.* We believe our Minimize then Union framework can be incorporated with other alternative automata for RE matching.

References

- [1] Snort. <http://www.snort.org/>.
- [2] A. V. Aho and M. J. Corasick. Efficient string matching: an aid to bibliographic search. *Communications of the ACM*, 18(6):333–340, 1975.
- [3] M. Alicherry, M. Muthuprasanna, and V. Kumar. High speed pattern matching for network ids/ips. In *Proc. 2006 IEEE International Conference on Network Protocols*, pages 187–196. Ieee, 2006.
- [4] M. Becchi and S. Cadambi. Memory-efficient regular expression search using state merging. In *Proc. INFOCOM*. IEEE, 2007.
- [5] M. Becchi and P. Crowley. A hybrid finite automaton for practical deep packet inspection. In *Proc. CoNext*, 2007.
- [6] M. Becchi and P. Crowley. An improved algorithm to accelerate regular expression evaluation. In *Proc. ACM/IEEE ANCS*, 2007.
- [7] M. Becchi and P. Crowley. Efficient regular expression evaluation: Theory to practice. In *Proc. ACM/IEEE ANCS*, 2008.
- [8] M. Becchi and P. Crowley. Extending finite automata to efficiently match perl-compatible regular expressions. In *Proc. ACM CoNEXT*, 2008. Article Number 25.
- [9] A. Bremler-Barr, D. Hay, and Y. Koral. Compactdfa: Generic state machine compression for scalable pattern matching. In *Proc. IEEE INFOCOM*, pages 1–9. Ieee, 2010.
- [10] B. C. Brodie, D. E. Taylor, and R. K. Cytron. A scalable architecture for high-throughput regular-expression pattern matching. *SIGARCH Computer Architecture News*, 2006.
- [11] C. R. Clark and D. E. Schimmel. Efficient reconfigurable logic circuits for matching complex network intrusion detection patterns. In *Proc. Field-Programmable Logic and Applications*, pages 956–959, 2003.
- [12] C. R. Clark and D. E. Schimmel. Scalable pattern matching for high speed networks. In *Proc. 12th Annual IEEE Symposium on Field-Programmable Custom Computing Machines (FCCM)*, Washington, DC, 2004.
- [13] J. E. Hopcroft. *The Theory of Machines and Computations*, chapter An nlogn algorithm for minimizing the states in a finite automaton, pages 189–196. Academic Press, 1971.
- [14] J. E. Hopcroft, R. Motwani, and J. D. Ullman. *Introduction to Automata Theory, Languages, and Computation*. Addison-Wesley, 2000.
- [15] S. Kong, R. Smith, and C. Estan. Efficient signature matching with multiple alphabet compression tables. In *Proc. 4th Int. Conf. on Security and privacy in communication networks (SecureComm)*, page 1. ACM Press, 2008.
- [16] S. Kumar, B. Chandrasekaran, J. Turner, and G. Varghese. Curing regular expressions matching algorithms from insomnia, amnesia, and acalculia. In *Proc. ACM/IEEE ANCS*, pages 155–164, 2007.
- [17] S. Kumar, S. Dharmapurikar, F. Yu, P. Crowley, and J. Turner. Algorithms to accelerate multiple regular expressions matching for deep packet inspection. In *Proc. SIGCOMM*, pages 339–350, 2006.
- [18] S. Kumar, J. Turner, and J. Williams. Advanced algorithms for fast and scalable deep packet inspection. In *Proc. IEEE/ACM ANCS*, pages 81–92, 2006.
- [19] Y. Liu, L. Guo, M. Guo, and P. Liu. Accelerating DFA construction by hierarchical merging. In *Proc. IEEE Ninth International Symposium on Parallel and Distributed Processing with Applications*, 2011.
- [20] C. R. Meiners, J. Patel, E. Norige, E. Torng, and A. X. Liu. Fast regular expression matching using small teams for network intrusion detection and prevention systems. In *Proc. 19th USENIX Security Symposium*, Washington, DC, August 2010.
- [21] A. Mitra, W. Najjar, and L. Bhuyan. Compiling PCRE to FPGA for accelerating SNORT IDS. In *Proc. ACM/IEEE ANCS*, 2007.
- [22] J. Moscola, J. Lockwood, R. P. Loui, and M. Pachos. Implementation of a content-scanning module for an internet firewall. In *Proc. IEEE Field Programmable Custom Computing Machines*, 2003.
- [23] V. Paxson. Bro: a system for detecting network intruders in real-time. *Computer Networks*, 31(23-24):2435–2463, 1999.
- [24] M. Roesch. Snort: Lightweight intrusion detection for networks. In *Proc. 13th Systems Administration Conference (LISA)*, USENIX Association, pages 229–238, 1999.
- [25] R. Sidhu and V. K. Prasanna. Fast regular expression matching using fpgas. In *Proc. FCCM*, pages 227–238, 2001.
- [26] R. Smith, C. Estan, and S. Jha. Xfa: Faster signature matching with extended automata. In *Proc. IEEE Symposium on Security and Privacy*, pages 187–201, 2008.
- [27] R. Smith, C. Estan, S. Jha, and S. Kong. Deflating the big bang: fast and scalable deep packet inspection with extended finite automata. In *Proc. SIGCOMM*, pages 207–218, 2008.
- [28] R. Sommer and V. Paxson. Enhancing bytelevel network intrusion detection signatures with context. In *Proc. ACM CCS*, pages 262–271, 2003.
- [29] I. Sourdis and D. Pnevmatikatos. Pnevmatikatos: Fast, large-scale string match for a 10gbps fpga-based network intrusion detection system. In *Proc. Int. on Field Programmable Logic and Applications*, pages 880–889, 2003.
- [30] I. Sourdis and D. Pnevmatikatos. Pre-decoded cams for efficient and high-speed nids pattern matching. In *Proc. Field-Programmable Custom Computing Machines*, 2004.
- [31] J.-S. Sung, S.-M. Kang, Y. Lee, T.-G. Kwon, and B.-T. Kim. A multi-gigabit rate deep packet inspection algorithm using team. In *Proc. IEEE GLOBECOM*, 2005.
- [32] L. Tan and T. Sherwood. A high throughput string matching architecture for intrusion detection and prevention. In *Proc. 32nd Annual Int. Symposium on Computer Architecture (ISCA)*, pages 112–122, 2005.
- [33] N. Tuck, T. Sherwood, B. Calder, and G. Varghese. Deterministic memory-efficient string matching algorithms for intrusion detection. In *Proc. Infocom*, pages 333–340, 2004.
- [34] L. Yang, R. Karim, V. Ganapathy, and R. Smith. Fast, memory-efficient regular expression matching with NFA-OBDDs. *Computer Networks*, 55(55):3376–3393, 2011.
- [35] F. Yu, Z. Chen, Y. Diao, T. V. Lakshman, and R. H. Katz. Fast and memory-efficient regular expression matching for deep packet inspection. In *Proc. ANCS*, pages 93–102, 2006.
- [36] F. Yu, R. H. Katz, and T. V. Lakshman. Gigabit rate packet pattern-matching using TCAM. In *Proc. 12th IEEE Int. Conf. on Network Protocols (ICNP)*, pages 174–183, 2004.