

Enhanced Hierarchical Multipattern Matching Algorithm for Deep Packet Inspection

D.Sahithi [M.Tech (CSE)],
Dept. of Computer Science & Engineering
St. Ann's College of Engg & Tech
Chirala, A.P, India

Dr.P.Harini [H.O.D of CSE]
Dept. of Computer Science & Engineering
St. Ann's College of Engg & Tech
Chirala, A.P, India

Abstract— Detection engines capable of inspecting packet payloads for application-layer network information are urgently required. The most important technology for fast payload inspection is an efficient multipattern matching algorithm, which performs exact string matching between packets and a large set of predefined patterns. This paper proposes a novel Enhanced Hierarchical Multipattern Matching Algorithm (EHMA) for packet inspection. Based on the occurrence frequency of grams, a small set of the most frequent grams is discovered and used in the EHMA. EHMA is a two-tier and cluster-wise matching algorithm, which significantly reduces the amount of external memory accesses and the capacity of memory. Using a skippable scan strategy, EHMA speeds up the scanning process. Furthermore, independent of parallel and special functions, EHMA is very simple and therefore practical for both software and hardware implementations. Simulation results reveal that EHMA significantly improves the matching performance. The speed of EHMA is about 0.89-1,161 times faster than that of current matching algorithms. Even under real-life intense attack, EHMA still performs well.

Keywords-component; Network-level security and protection, network security, intrusion detection, pattern matching, content inspection.

I. INTRODUCTION

Network services are extremely important since many companies provide services over the Internet. A variety of Internet-based applications have created a strong demand for content-aware services, network policy, and security management. Furthermore, increasing amounts of important information exist in packet payloads. Therefore, low-layer network equipment is inadequate for checking the information, since it only checks specified fields of the packet headers. High-layer network equipment providing in-depth packet inspection, such as intrusion detection systems (IDSs), application firewalls, antivirus appliances, and layer-7 switches, is a prerequisite in a network. Such equipment typically contains a policy or rule database applied to finding certain packets over the network. Every rule in the database consists of several patterns (also called signatures) and a matching action (or a series of actions). These patterns describe the fingerprints of packets. The network equipment applies the predefined patterns to identify and manage the monitored packets over the network.

Different network equipment may have different pattern databases applied, respectively, to attack detection, bandwidth management, load balancing, and virus blocking over the network. However, they have similar features in terms of patterns and matching procedures. The number of patterns is typically a few thousands, and the lengths of the patterns are varied. The patterns may appear anywhere in any packet payload. Consequently, the emerging high-layer network equipment needs a pattern detection engine capable of in-depth packet inspection, which searches the entire packet headers and payloads for pattern matching. Network equipment then employs the detection results to manage network systems intelligently. For instance, Snort is an open-source network-based intrusion detection system (NIDS) and is adopted for detecting anomalous intruder behavior with a set of patterns and generating logs and alerts from predefined actions [1]. One of the patterns of Nimda worm is described as "GET/scripts/root.exe?/c+dir." When the detection engine of Snort finds this pattern existing in a packet, the corresponding alert is generated to warn network administrators. The pattern matching is considered as the most resource-intensive task in the Snort detection engine [2]. Hence, this study focuses on the nascent issues of the payload inspection.

The most important part of a detection engine is a powerful multi pattern matching algorithm, which can efficiently process the pattern matching task to keep up with the growing data volume in the network. However, conventional string-matching algorithms are impractical for packet inspection [3], [4], [5]. Due to the large pattern database, an effective detection engine must be able to search for a set of patterns simultaneously, rather than iteratively performing the single-pattern matching. While considering implementation issues of the network equipment, the performance of processing packets is not only affected by the computation time but also strongly affected by the memory latency. As is well known, the rate of improvement in processor speed exceeds that of improvement in memory speed [6]. The gap has been the largest problem for system builders. Therefore, the vital issue of designing a high-speed detection engine is to reduce the number of external memory accesses[8]. This study proposes a novel Enhanced Hierarchical Multipattern Matching Algorithm (EHMA) for fast packet inspection, which simultaneously searches the packet payload

for a set of patterns. This study contributes modifications to the hierarchical matching algorithm (HMA) [9] and introduces the idea of a sampling window and a Safety Shift Strategy in addition. EHMA is a two-tier and cluster-wise matching algorithm and can perform fast skippable payload scan. Based on the occurrence frequency of grams, this study discovers a small set of signatures from the patterns themselves to narrow the searching domain. A Min-Max strategy is used in the EHMA. The hit rate of the first-tier table in the EHMA is minimized, while the spread of patterns in the second-tier table is maximized. Accordingly, EHMA significantly reduces the number of memory accesses and pattern comparisons. EHMA can skip unnecessary payload scans by applying the proposed Safety Shift Strategy, which is based on a frequency-based bad gram heuristic. The frequency-based bad gram heuristic is a modification of the bad grouped character heuristic of Wu-Manber (WM) algorithm [10]. Therefore, EHMA has the advantages of both HMA and WM.

The memory space and the number of external memory accesses required by the proposed EHMA are much smaller than those required by state-of-the-art multi pattern matching algorithms. EHMA needs less than 40-Kbyte memory space to construct required tables for the Snort patterns and, therefore, enables small-scale and cost-effective hardware implementations. Using only 768-byte on-chip memory, EHMA reduces the average number of external memory accesses to 0.06-0.19 and, thus, significantly improves the matching time of the detection engine. Simulation results reveal that EHMA outperforms the state-of-the-art algorithms. Even under real-life intense attack, EHMA still outperforms others. Because it employs only basic instructions and two small index tables, EHMA is very simple for hardware and software implementations. Consequently, the proposed EHMA is a very cost-effective and efficient mechanism for real-life network detection systems.

The rest of this paper is organized as follows: Section 2 Presents previously proposed pattern matching algorithms and the fundamental definitions. Section 3 then describes the proposed EHMA in detail. Next, Section 4 presents the performance and memory requirements of EHMA. Conclusions are finally drawn in Section 5.

2. RELATED WORK

This section discusses the main concepts and the limitations of the state-of-the-art exact string matching algorithms that have been used or modified for packet inspection. Some fundamental definitions and notations used in this study are presented.

2.1 Notations

An array is used to represent a string of characters from an alphabet set A . Namely, an element representing string T at the position i is given by $T[i]$, where $T[i] \in A$. The absolute value of an object means the size of the object. For instance, $|T|$ denotes the length of the string T , and $|A|$ is the number of elements in the set A . A function $\text{sub}(T, i, B)$ is defined as the substring of T from $T[i]$ to $T[i + B - 1]$. A string can also be denoted as a set of B -grams, where a gram is

defined as a group of characters, and B is the number of characters in a gram. For instance, the string "green" can be converted into a set of 2-grams {"gr", "re", "ee", "en"} when $B = 2$. The i th B -gram of a string T is represented as $TB[i]$. Let $P = \{p_i\}$ be a set of distinct patterns, where p_i denotes a pattern with an identification number (ID) i . The payload of an input packet T and the pattern $p_i \in P$ are both strings drawn over A with finite length $|T|$ and $|p_i|$, respectively. The notation $e:f$ denotes the value of the field (or offset) f at the entry (or address) e . If e is a table, then $e:f$ means all fields named f of the table e .

A single-pattern matching algorithm is used to search a string (or text) T for the first occurrence or all occurrences of one given pattern. A multipattern matching algorithm is applied to search the input T for all occurrences of any pattern $p_i \in P$, or to corroborate that no pattern of P is in T , where the number of patterns is from hundreds to thousands. In other words, the algorithm aims to find all the matched patterns in T , say $P_M \subseteq P$ such that $P_M = \{p_j \mid p_j \in T \text{ and } p_j \in P\}$. P_M can be applied to any high-level detecting rule, such as the high-priority-win, first-matched-win, or other state-concerned rules.

2.2 Previous Work

Single-pattern matching algorithms were originally proposed to perform text searching problem in computer systems. In single-pattern matching, Boyer-Moore (BM)-based algorithms provide the best average-case performance in terms of computation complexity, which is sub linear to the input string [3], [13]. The BM algorithm uses the bad character and good suffix heuristics to build a skip table and a shift table, respectively [13]. The Boyer-Moore-Horspool (BMH) algorithm, which is a variant of BM, slightly modifies the bad character heuristic to construct a single skip table [3]. The tables of BM and BMH are pre computed and used to determine the number of safety shifts of each character for the searching process. Some characters of T can thus be skipped in the matching process on specific conditions. Several approaches apply the BM-based single-pattern matching algorithms iteratively to solve the multi pattern matching problem. However, network equipment usually has a large pattern database. Iteratively performing the single-pattern matching for multi pattern matching in the packet inspection engine is inefficient. Markatos et al.'s approach promotes Snort by using a bitmap filter before BMH but still searches for only one pattern in each iteration [11].

Several modifications to BM-based algorithms have been proposed for the multi pattern matching. Risk and Varghese's (RV) approach groups all patterns to pre calculate the number of safety shifts of each character [5]. The WM approach, which assumes that all patterns are larger than M characters, groups B -grams of the M -character prefixes of all patterns to build a shift table [10]. The WM's shift table contains the valid shifts of each B -gram. Liu et al.'s algorithm [a variant of the WM algorithm using a grouped prefix hash (WM-PH)] groups the B -character prefixes of all patterns to build a large hash table, in which each entry contains valid shifts of the corresponding B -character prefix [12]. However, the maximum shift value of RV and WM must be not larger

than the minimum pattern length in P , in order to avoid missing any pattern. Thus, RV and WM are unfeasible when the pattern set includes single-character patterns. The required memory space of the table in the algorithms WM and WM-PH is in the order of $O(jAB)$.

Generally, $B = 3$, and the table consists of 16 million entries when the alphabet size is 256 as in 1-byte coding. The large tables must be stored in the external memory, which leads to long access delay during the matching process.

It has been pointed out that the Aho-Corasick (AC) algorithm provides the best worst-case computational time complexity [4]. Using a compressed structure, Tuck et al. proposed the AC algorithm with memory compression (AC-C), a modification of AC, and reduced the required memory to about 2 percent of AC [8]. ACM applied a magic number derived from the Chinese Remainder Theorem to AC [14]. ACM reduced the required memory space and computation complexity, thus improving the worst-case performance. However, the required memory of AC-C and ACM is typically too large to be cached in the on-chip memory of embedded systems, field-programmable gate arrays (FPGAs), and network-processor-based platforms. Although the AC-based algorithms have the best worst-case computational time complexity, the latency of external memory accesses dominates the processing performance rather than computational time. Coit et al. proposed a matching algorithm for Snort that combines BM and AC [15]. However, this algorithm requires three times the memory of the standard version and may produce inconsistent matching results.

A Piranha algorithm was proposed based on an idea that if a least popular B -gram of a pattern exists in a packet, then this packet may have a pattern [16]. A least popular gram of a pattern was chosen as an index key of a pattern. However, the Piranha algorithm cannot handle the patterns smaller than B , and the required memory space is very large ($O(jAjB)$). Although the idea of least popular index keys can reduce the collisions of patterns, the hit rate of index table is increased, thus increasing the number of external memory accesses and pattern comparisons.

In the case of hardware solutions, Li et al. presented an FPGA-based detection engine for NIDSs, using the internal content addressable memory (CAM) technology to speed up multi pattern matching [17]. Since an internal CAM of FPGA is not large enough to store all patterns, Li et al.'s approach has to dynamically reload a block of patterns into the CAM, causing long latency. Moreover, the patterns of varied lengths complicate the formulation of a CAM for exact matching, but Li et al.'s approach does not mention the solution for patterns with varied lengths. Dharmapurikar et al. used Bloom Filters (BFs) and Kim and Kim employed mask filters in the FPGA-based packet inspection [18], [19]. However, these two methods only act as pre filters and have to cooperate with another string matching algorithm to verify a match, and furthermore, this BF-based algorithm can be used only in the case that all patterns are longer than a certain length. Lu et al. used several binary CAMs and BFs to implement parallel compressed deterministic finite automata (DFA), and Dharmapurikar et al. combined AC with BFs for packet inspection [20], [21]. These two methods applied parallel BFs and assumed that BFs can execute one query every clock cycle. However, these architectures and assumptions can only be

established in some specific hardware implementations. BFs are inefficient in the software implementations, because one BF consists of several hash functions and the computation time of hash functions is usually expensive in software [6].

3 THE ENHANCED HIERARCHICAL MULTI-PATTERN ALGORITHM

Some network equipment is implemented by network processors, FPGAs, networks-on-chip (NOCs), or systems-on-a-programmable-chip (SOPCs) to improve the performance. The embedded memory of these platforms is typically very small. For instance, the Intel IXP2x00 network processor has only a 4-Kbyte instruction cache and a 2-Kbyte data cache in each micro engine, while the Vitesse IQ2000 network processor has a 4-Kbyte data cache (2 Kbytes for local storage and 2 Kbytes for reserved header buffers) [22], [23]. Although high-end FPGAs providing up to 1-Mbyte embedded memory are available, linking many memory blocks degrades the chip performance. Nevertheless, the required memory of the previous pattern matching algorithms is generally larger than 300 Kbytes for NIDSs. Hence, the patterns and the tables built by matching algorithms need to be stored in external memory.

However, frequently accessing the external memory (to read patterns or tables) significantly decreases the matching efficiency due to the external memory access latency being very long and indeterminable. For example, Intel IXP2x00 needs about one cycle for one microprocessor instruction but about 150 cycles for each access from SRAM (or 250-300 cycles from DRAM) [7]. The memory latency strongly affects the throughput of pattern matching. Therefore, reducing the number of required external memory accesses is more important than reducing the amount of computational time. This study proposes an EHMA based on a hierarchical and cluster-wise architecture. EHMA comprises two small index tables, namely the first-tier table (H_1) and the second-tier table (H_2). These two tables act as filters to avoid unnecessary external memory accesses and pattern comparisons and, thereby, pass the innocuous packets quickly in the online matching process. The second-tier procedure (Tier-2 Matching) activates only after the first-tier procedure (Tier-1 Matching) gets a match. Using H_2 , which indicates a small subset of patterns that are similar to the input packet, EHMA compares only a few selected patterns of P with the suspected substrings of the packet, rather than comparing all patterns with all substrings of the packet. Furthermore, a frequency-based bad gram heuristic is proposed in the EHMA to determine the safety shifts on the input strings during the online matching process. In other words, some characters of the input packets can be safely skipped without any process. External memory accesses are needed only in the Tier-2 Matching state. Consequently, EHMA significantly enhances the matching performance and effectively reduces the number of external memory accesses, string comparisons, and character scans, by utilizing two small index tables.

This study proposes a general frequent-common gram searching (GFGS) algorithm and a cluster balancing strategy (CBS) to lower the size of the tables H_1 and H_2 . The following sections describe the GFGS, CBS, and the Safety

Shift Strategy in detail. The hierarchical online matching using these two index tables, namely Tier-1 Matching and Tier-2 Matching, is then shown.

3.1 The GFGS Algorithm

In the high-layer intrusion detection, patterns may appear anywhere in the packet payload, making the attacking packets difficult to recognize. GFGS assumes that a small set of signatures can be found from the patterns themselves, then the suspicious substrings of T may be easier to distinguish from the innocent parts, and the pattern matching is therefore faster. A set of significant grams is defined as representatives of a pattern set P, given by $\mathcal{G} = \{g_i\}$, where the size of a gram is B_i characters. The set \mathcal{G} is much smaller than \mathcal{P} . Only when at least a significant gram occurs in the payload, a pattern may exist. That is, when at least one B_i -gram of p_i belonging to \mathcal{G} occurs in the payload T, the pattern $p_i \in \mathcal{P}$ may be found in T. Many innocent B_i -grams of T, which do not belong to \mathcal{G} , can be filtered in the Tier-1 Matching when scanning the packet payload. Obviously, smaller \mathcal{G} leads to fewer pattern comparisons and, thus, faster pattern matching. The GFGS is proposed to find the smallest \mathcal{G} from P.

Define \mathcal{P}_g as a subset of P, with $\mathcal{P}_g = \{p_i \mid p_i \text{ has the gram } g, \forall p_i \in \mathcal{P}_g\}$, where g is called the common gram of those patterns in the set \mathcal{P}_g . Notably, if a common gram appears in the distinct patterns more frequently than other grams and it is selected as one of the significant grams, then a smaller \mathcal{G} is found. Based on this inference, the GFGS algorithm is designed to find the frequent-common gram set F, such that F is the minimum set of significant grams to represent a pattern set P. In the GFGS, the common grams are searched only from the sampling window, which is defined as the last W characters of the first m characters of a pattern. The range of m is $M < m < j p_i$, where M denotes the minimum pattern length of all patterns, and $j p_i$ is the current pattern length. Fig.1 illustrates the sampling window, where B_1 is the size of a frequent-common gram, $B_1 < W$, and B_2 is the size of the second pivot in the H2 table, which is explained later.

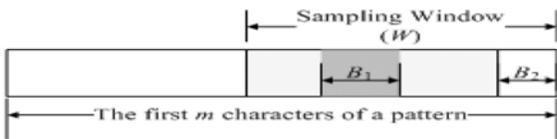


Fig. 1. The sampling window.

The GFGS algorithm is presented in Fig. 2. A bitmap vector $V = (v_i)$ and a matrix $R = (r_{ij})$ are temporary memories, where $0 < i, j < jA_j$. Vector V records the occurrence of each B_i -gram in a pattern; R is used for recording frequency, where $r_{ij}, i = j$, indicates the number of concurrent occurrences of two B_i -grams g_i and g_j in P; and r_{ii} records the frequency of the B_i -gram g_i occurring in distinct patterns. For instance, $r_{ij} = 2$ means there are two patterns, each containing both g_i and g_j . In the GFGS algorithm, each pattern is first transferred into a set of B_i -grams, and the occurrence of each B_i -gram is recorded in the bitmap V, where B_i is predefined and depends on the available on-chip memory space. Matrix R is then derived from V (as shown in line 4 of Fig. 2). Second, the largest occurrence frequency

is found, and its corresponding gram gf is selected as one of F. The elements of R relating to gf are subtracted accordingly to renew R. GFGS is repeated until all elements on the diagonal of R become zero. GFGS uses only a matrix and a vector to discover F from P.

```

GFGS Algorithm;
Input: Given a set of patterns P, the parameters: W, B2, B1, and m.
Output: A set of frequent-common grams F.
1 Initialize: F ← ∅, V and R are set to zero;
2 For each pattern pi of P, 0 ≤ i < |P| do /*build a matrix R*/
3   Transfer the first W- B2 bytes of the sampling window of the pattern pi into
   B1-grams, and set the element of a vector V: vj ← 1 if B1-gram = j; otherwise
   vj ← 0;
4   Read V. For each vj = 1, set the elements of matrix R: rjk ← rjk + vk, ∀ k,
   0 ≤ k < |ΛB1|;
5   While (rii ≠ 0, ∀ 0 ≤ i < |ΛB1|) do
6     Find a frequent-common gram gf, where rff = max{rii | ∀ i, 0 ≤ i < |ΛB1|};
7     Add this gram into F: F ← F ∪ {gf};
8     For 0 ≤ i < |ΛB1| do /* refresh the diagonal of R*/
9       rii ← rii - rff, if rii > rff; otherwise, rii ← 0;
10  Return;
    
```

Fig. 2. The GFGS algorithm.

Fig. 3 plots the pattern spectrum of the Snort patterns with different gram sizes. The pattern spectrum indicates the occurrence frequency of grams of patterns. Fig. 3a shows the distribution of 2-grams of patterns, and Fig. 3b is the distribution of 1-gram of patterns. As shown in the figures, they are not normally distributed and have several peaks, which mean that some grams obviously occur more frequent than others. Hence, GFGS can easily discover the most frequent grams from patterns and obtain a small F as the signatures of pattern set. Since both 1-gram and 2-gram spectrums have peaks, the gram size of F can be one or two, depending on the available size of on-chip cache.

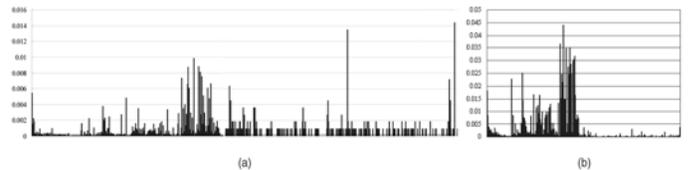


Fig 3. The pattern spectrum when JP J = i, 200. (a) Spectrum of 2-grams. (b) Spectrum of 1-gram.

3.2 Cluster Balancing Strategy (CBS)

Most packets are innocent in general situations. Even a harmful packet may contain only few patterns. Therefore, comparing all of the patterns in the large P with each input packet is time consuming. If the patterns in P can be distributed into different small clusters based on their similarity, then only the pattern in each cluster that is most similar to the suspected packet needs to be compared, thus improving the efficiency of the matching process. This section presents strategies to attain this goal. First, the method of clustering a set P based on the similarity of patterns is described. Then, a CBS is adopted to balance the cluster size. A second-tier table (H2) for online matching can be constructed based on the clusters.

The clustering pivots are the keys used to distribute patterns, where each clustering pivot is a common gram of patterns defined previously. Two common grams are employed as a pair of clustering pivots, called a pivot pair, say (a, b), where the first pivot is a frequent-common gram,

and the second pivot is the substring following the frequent-common gram. Let $P_{a,b}$ represent a cluster of selected patterns (a subset of patterns) with the pivot pair (a, b), which means that $P_{a,b} = \{p_i \mid p_i \text{ contains } 'ab' \text{ as a substring}\}$, where 'ab' is the combination of two strings a and b and is a substring of p_i ; F is the result of GFGS, and B_2 is the length of the second pivot. Notably, a pattern is assigned to only one cluster in the clustering strategy, although a pattern may have more than one pivot pair. That is, the clusters have the following properties: for any cluster $P_{a,b}$, $\forall p \in P_{a,b}, [a, b] \text{ is a substring of } p$ and $\forall p \in P_{a,b}, p \text{ is not a substring of any other } p' \in P_{a,b}$. Since a pattern may have several opportunities to select a cluster, a better assignment can lower the maximum cluster size and, thereby, improve the worst-case performance of EHMA.

The CBS is given as follows:

1. According to GFGS, for any given p_i , there exists a B_1 -gram $g \in F$, where B_1 is the length of a frequent-common gram. To balance the cluster size, CBS finds the smallest $n_{a,b}$, given by $n_{x,y}$, among all available pivot pairs (a, b)'s of p_i , for all $a \in F$ and 'ab' c p_i .

2. After grouping p_i into the smallest cluster $P_{x,y}$, the corresponding $n_{x,y}$ is also incremented.

All patterns are distributed sequentially into the designate clusters. Accordingly, GFGS and CBS divide the large P into smaller subsets.

3.3 Safety Shift Strategy

This section presents a safety shift strategy to derive the values of the shift fields of H_1 and H_2 . H_1 and H_2 can use the same strategy to derive their safety shifts, respectively. As mentioned previously, as long as no frequent-common gram is matched in input strings, then no pattern exists. Therefore, if no frequent-common gram is missed, then no pattern will be missed. The safety shift strategy is based on a modified bad grouped character heuristic [7], named frequency-based bad gram heuristic in this study. The safety shift strategy ensures that no frequent-common gram is missed during a skippable scanning process. The proposed strategy helps EHMA to speed up the online matching process, since certain characters can be skipped unhesitatingly.

Assume that x identifies all possible index keys and that the length of x is B . Because the index keys of H_1 and H_2 are different, the parameters used to determine the shift fields of these two tables are different. For H_1 , as the length of a frequent-common gram is B_1 , thus $x \in AB_1$ and $B = B_1$. For H_2 , since x is all the possible of the pivot pairs (a, b), $x \in F \times AB_2$ and $B = B_1 + B_2$. The basic concept of the safety shift strategy is that: if x is not a gram of any pattern, and any suffix of x is not any prefix of any pattern in P , then it is safe to shift m when x is scanned; otherwise, the number of safety shifts is the offset between the rightmost occurrence position of x and the position of the frequent-common gram nearest to x . Two parameters are needed to derive the safety shifts, namely W and m , as shown in Fig. 1. Assume that $B < W < m$, and define the safety shifts of each entry $(H(x)):\text{shif } t$ as follows:

1. Initially, all shift fields of the table H are set as

If $m > W$, then $H(x):\text{shif } t = m - W + q$,

Where $q = \min \{q \mid \exists \text{ sub}(x, q+1, B-q) = \text{sub}(p, 1, B-q), \forall p \in P \text{ and } 1 < q < B \text{ when } B > 1 \text{ and } q \text{ exists; otherwise, } q = B. \text{ Else } H(x):\text{shif } t = r, \text{ where } r = \min \{r \mid \exists \text{ sub}(x, r+1, B-r) = \text{sub}(f, 1, B-r), \forall f \in F, 1 < r < B, \text{ and } r+B < W \text{ when } B > 1 \text{ and } r \text{ exists; otherwise, } r = B. \}$

2. If the current $H(x):\text{shif } t > m - W - i + 1$, then update the entry, so that $H(x):\text{shif } t = m - W - i + 1$.

3. For each i th B -gram of each pattern $pB[i]$, where $m - W < i < m - B + 1$, set $x \in pB[i]$ if the entry $H(x)$ exists:

If $x \in F$, then $H(x):\text{shif } t = 0$;

Else If the current $H(x):\text{shif } t > r$, then update the entry: $H(x):\text{shif } t = r$, Where $r = \min \{r \mid \exists \text{ sub}(x, r+1, B-r) = \text{sub}(f, 1, B-r), \forall f \in F, 1 < r < B, \text{ and } r+B < W \text{ when } B > 1 \text{ and } r \text{ exists; otherwise, } r = B. \}$

Sub (f, 1, B-r), $\forall f \in F, 1 < r < B$, and $r+B < W$ g

When $B > 1$ and r exists; otherwise, $r = B$.

Notably, the maximum shift of EHMA is m while $W = B$. The frequent-common grams and the sampling window are introduced in the proposed frequency-based bad gram heuristic to improve the flexibility and the efficiency. Additionally, comparing EHMA with WM, the maximum safety shift is raised from $m - B + 1$ to m . The shift value of the proposed EHMA is similar to but larger than the shift value of WM, when the given parameters are $m = M$ and $W = B$.

3.4 THE ONLINE HIERARCHICAL AND CLUSTER-WISE MATCHING

The previous sections presented the offline stage of EHMA, which builds two index tables H_1 and H_2 , holding the indexing and pattern information in the cache memory and external memory, respectively. These two tables are regarded as the two-tier filters and indices for the online matching. This section presents the online matching procedure in detail.

3.4.1 Tier-1 Matching

In online matching, the payload T is scanned from left to right, and each B_i -gram of T is the key to fetch the entry $H_i(t_i)$, where $t_i = T B_i[i]$. The H_i acts as the first-tier filter of EHMA, by checking whether T may likely contain patterns belonging the pattern set P . Because H_i is small enough to be stored in the on-chip memory during the online matching procedure, the latency of accessing H_i is very small.

In the Tier-1 Matching, first the shift field is checked. If $H_i(t_i):\text{shif } t = 0$, i.e., $t_i \in F$, then no external memory is necessary. The obtained $H_i(t_i):\text{shif } t$ also determines the number of grams that can be skipped without further process. The next gram to check is then $T B_i[i + H_i(t_i):\text{shif } t]$. After reading the next gram, the matching process repeats as in the previous steps and remains in the Tier-1 Matching. Because $\prod_{j \in J} p_j$, the probability of $t_i \in F$ is small and most grams of T gain the shifts, thus avoiding the Tier-2 Matching. Consequently, both the number of string comparisons and the costly memory accesses can be significantly reduced.

Otherwise, if $t_i \in F$, then T may contain a malicious pattern $p \in P$, where $t_i \in p$. Simply stated, if $H_i(t_i):\text{shif } t = 0$, then T may have a pattern that belongs to the cluster of pivot pair (t_1, t_2) , where $t_2 = T B_2[i + B_1]$.

Therefore, the matching procedure activates Tier-2 Matching to identify the pattern. If $H_i(t_i):pid$ is not NULL, then the current gram t_i itself is a pattern, and this matched pattern is also added into P^M .

3.4.2 Tier-2 Matching

After the Tier-1 Matching, if $H_i(t_i):shift = 0$, then the matching procedure proceeds to the Tier-2 Matching. The function $H_2(t_i, t_2)$ indicates the location of the corresponding cluster according to input T . Since EHMA is a cluster-wise matching algorithm, only the patterns in the small cluster of pivot pair (t_i, t_2) , which are similar to T , are loaded to the processing unit for further checks.

Tier-2 Matching first checks the pid field of H_2 . If $H_2(t_i, t_2):pid$ is NULL, then the cluster (t_i, t_2) contains no pattern, and no pattern comparison is necessary. Otherwise, if $H_2(t_i, t_2):pid$ is not NULL, then this cluster contains patterns. The pattern content in the $H_2(t_i, t_2):data$ is then compared with the corresponding substring of T : $sub(T, i - H_2(t_i, t_2):offset, H_2(t_i, t_2):size)$. If $H_2(t_i, t_2):next$ is valid and points to the next entry, here given by $H_2(a, b)$, then the cluster contains other patterns. Similarly, the pattern in $H_2(a, b):data$ is also fetched and compared with the substring of T starting at $T[i - H_2(a, b):offset]$ of length $H_2(a, b):size$. Every matched pattern is added to the matched pattern set P^M and its corresponding matched pid set PID^M in order. Until all patterns in this cluster are checked, the next gram $T_{Bi}[i + H_2(t_i, t_2):shift]$ is then read, and the online matching procedure returns to the Tier-1 Matching. $H_2(t_i, t_2):shift$ also indicates the number of characters of T that can be skipped, since the next possible frequent-common gram may only appear far away than $H_2(t_i, t_2):shift$.

Notably, if a pattern pk exists in T , then all grams of pk appear in T . The clustering pivot pair of pattern pk ($pk[j], pk[j + B_i]$) is certainly scanned, say at t_i and t_2 , so that $t_i = pB_i[j] \in F$ and $t_2 = pB_2[j + B_i]$. Pattern pk is then recognized when T is compared with the patterns in the cluster (t_i, t_2) during the online matching procedure. Based on the Safety Shift Strategy, EHMA never skips any frequent-common gram. Consequently, no patterns in the payload T are missed.

```

Procedure Tier-1Matching( $T, H^1, M, W, B_1$ )
Input: Packet payload  $T$ , a first-tier hash table:  $H^1$ , the minimum pattern length  $M$ , the length of the frequent-common gram  $B$  and the length of the sampling window  $W$ .
Output: The output of Tier-2Matching.
1  $i \leftarrow M - W + 1$ ;
2 While  $i \leq |T| - B_1$  do
3   Read the  $i$ -th  $B_1$ -gram of  $T$ :  $gram \leftarrow T^{B_1}[i]$ ;
4   If  $H^1(gram).shift > 0$ , then  $shift \leftarrow H^1(gram).shift$ ;
5   Else
6     If  $H^1(gram).fid \neq \text{NULL}$ , then  $shift \leftarrow \text{Tier-2Matching}(T, H^2, B_2, i)$ ;
7     If  $H^1(gram).pid \neq \text{NULL}$ , then
8        $P^M \leftarrow P^M \cup \{gram\}$ ;
9       If  $shift = 0$ , then  $shift \leftarrow 1$ ;
10    Jump over the string:  $i \leftarrow i + shift$ ; /*shift and read the next*/
11 End While
12 Return;

Procedure Tier-2Matching( $T, H^2, B_2, i$ )
Input: Packet payload  $T$ , a preprocessed indexing table:  $H^2$ , the length of the second pivot  $B_2$ , and the current pointer  $i$ 
Output: A safety shift number for Tier-1 Matching:  $shift$ , the matched pattern set of  $T$ :  $P^M$ , and its corresponding pid  $PID^M$ 
1 Load data from the external RAM at entry  $H^2(T^{B_2}[i], T^{B_2}[i + B_1])$  to a local buffer  $LB$ ;
2  $shift \leftarrow LB.shift$ ;
3 While  $(k \leftarrow LB.pid) \neq \text{NULL}$  do
4   Compare the substring of  $T$ :  $sub(T, i - LB.offset, LB.size)$  with the pattern  $LB.data$ ; /*Assume no fragmentation here*/
5   If it is matched then  $P^M \leftarrow P^M \cup \{pk\}$  and  $PID^M \leftarrow PID^M \cup \{k\}$ ;
6   If  $LB.next \neq \text{NULL}$  then
7     Load data from the external RAM at entry  $LB.next$  to the local buffer  $LB$ ;
8   Else
9     Jump to Line 10;
10 End While
11 Return  $shift$ ;
    
```

Fig.4. the online matching procedure, including Tier-1 Matching and Tier-2 Matching.

The online matching procedure of EHMA is described in Fig. 4, including Tier-1 Matching and Tier-2 Matching. Since EHMA introduces H_i and H_2 as filters, and CBS is employed, only a few suspected patterns are loaded from external memory and compared with T . Because generally most of the packets are innocent over the network, and the frequent-common grams (F) narrow the searching field, EHMA performs a fast scan over the packets. The returned result P^M includes all matched patterns for a given T and is applied to make the final decision and analyze the impending attacks. The final decision depends on decision-making rules.

An example is provided to demonstrate the online matching of EHMA. Assume that the H_i and H_2 tables have been built as Fig. 4, where $W = 3$ and $M = 6$. Assume that the input T is “kangaroo” as given in Fig. 5. The scan runs from left to right. The scan starts at “g” ($(M - W + i)$ th gram), obtaining H_i (“g”): $shift = 4$. Therefore, Tier-1 Matching shifts four characters. Because the pointer goes beyond $jT_j - B_i$ after the shift, EHMA completes scanning the input T . This example only requires one on-cache table lookup and no external memory access. By only checking T with the embedded table H_i , EHMA can know that T contains no pattern.



Fig. 5. An example of matching process with input “kangaroo.”

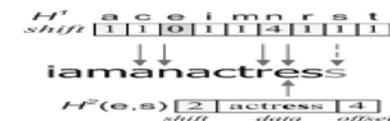


Fig.6 An example of matching process with input “iamanactress.”

Considering another example where $T = \text{'iamanactress'}$ as shown in Fig. 6, the first scanned B_i -gram is “a,” yielding H_i

(‘a’):shif t = i. Thus, the matching process stays in the Tier-1 Matching, and the next Bi -gram “n” is read after shifting one character, yielding Hi (‘n’):shif t = 4. Similarly, staying in the Tier-1 Matching, and the next Bi -gram “n” is read after shifting one character, yielding Hi (‘n’):shif t = 4. Similarly, staying in the Tier-1 Matching, the matching process obtains Hi (‘r’):shif t = i and Hi (‘e’):shif t = 0 in order after shifting. While Hi (‘e’):shif t = 0, the Tier-2 Matching is activated. After checking the field H2 (‘e’, ‘s’):pid and finding that it is not NULL, EHMA knows a suspected pattern may exist. The Tier-2 Matching then compares input T with the pattern in the cluster P e.s , where H 2 (‘e’, ‘s’):data = ‘actress’, and gets a match. Because this cluster contains no other patterns, the matching process returns to Tier-1 Matching with H 2 (‘e’, ‘s’):shif t = 2. Since the pointer goes beyond jT j— Bi after shifting two characters, the matching process for the input T is finished. In this case, Hi is checked four times, and H 2 is fetched only once for the string T of 12 characters. EHMA thus significantly reduces the latency caused by memory accesses.

4 RESULTS

As the number of network security threats rises, the NIDS has become one of the most important applications of packet inspection .Therefore, this study demonstrates the feasibility of integrating the proposed EHMA with the promising NIDS. This section presents the simulation results of EHMA deployed in the NIDS, compared with the original HMA [9], BMH algorithm [3], WM algorithm [10], WM-PH [12], and AC-C [8]. In the simulations, the assembly-like micro programs were emulated for EHMA, BMH, WM, WM-PH, and AC-C using RISC instructions of general network processors (such as ADD, XOR, MOV), and the number of instructions and the number of memory accesses needed to process a packet were calculated. To simplify the evaluation, the simulation assumed that one microprocessor was employed.

Items	Value
Time of one RISC instruction or one local memory access (w_l)	1 cycle
Latency for each external memory access (w_e)	10, 100 cycles
Packet payload length for Model I	512 bytes
Number of patterns in P ($ P $)	200, 400, ..., 5000
Simulation time for Model I	10 million packets

TABLE 1 the Simulation Parameters

4.1 Measurements

Define I as the average number of RISC instructions (including comparisons and calculations) and L as the average number of local memory accesses (including reading data from the cache to the registers for further processes), for each payload character in the pattern matching. E represents the average number of external memory accesses per input character, which includes loading the input packets, querying the entries of tables in the external memory, and fetching the patterns. w_l indicates the time needed by one instruction or one local memory/register access, and w_e indicates the time for one external memory access. The following measurements are given: the average computation cycles $I = I \times w_l$; the average memory latency $M = E \times w_e + L \times w_l$; and the total average matching time $T = I + M$, which is regarded as the overall performance.

In the simulations, the skip table of BMH was assumed to be small enough to be loaded into the cache memory, and therefore, only one external memory access was counted during the matching process of BMH for each pattern. One external memory access was assumed for AC-C, although it typically needs two memory references to fetch the transition matrices, and the fail table or the matched patterns. Table 1 lists the simulation parameters.

4.2 Traffic Models

The simulations used two free and real pattern sets, R1 and R2 , from Snort in August 2004 and May 2008, respectively [1], although the pattern set can be self-defined or any commercial pattern set. The number of distinct patterns is about 1,250 in the R1 , where the average length of a pattern is about 11.2 bytes (the statistics of the pattern set listed in Table 2); while the number of distinct patterns becomes up to about 5,000 in the R2 . Since Snort patterns are written in mixed plain text and hex formatted byte codes, the alphabet size (jA_j) was set to 256 in the simulations. In the simulation traffic models, Models I and II use R1, and Model III uses R2 as the matching pattern sets.

TABLE 2

The Pattern Size Distribution of Snort Rule Set R1

Pattern Size	=1	≤ 4	≤ 8	≤ 12	≤ 16	>16
Ratio	0.028	0.245	0.482	0.653	0.813	0.187

Table 3 shows the relationships between the number of patterns jP j and the number of frequent-common grams jF j of the EHMA, where the lengths of patterns are in the range from 1 to 122, $m = jpi j$, and the patterns are randomly selected from R1 . The results in Table 3 reveal that the growth rate of jF j is much slower than that of jP j.

TABLE 3

The Number of Frequent-Common Grams versus the Pattern Set Size

P	100	200	300	400	500	600	700	800	900	1000	1100	1200
F	11	28	32	37	45	49	52	58	66	74	75	77

TABLE 4

The Memory Requirements

	EHMA	HMA	WM	WM-PH	AC-C	BMH	BMH-O
Cache Memory	$O(A)$	$O(A)$	$O(1)$	$O(1)$	$O(1)$	$O(A)$	$O(1)$
External Memory	$O(P \times A)$	$O(P \times A)$	$O(A ^2)$	$O(A ^2)$	$O(S)$	$O(P \times A)$	$O(P \times A)$

4.2.1 Model I

In Model I, the synthetic malicious packets are generated by randomly choosing patterns from the pattern set P and spreading over the packet payloads. The attack load A is defined to represent the expected number of malicious patterns existing in one packet. For instance, if A = 2, then each packet contains two harmful patterns on average. Except for the injected patterns parameterized by A, the background characters of a packet were randomly drawn from

A to imitate the normal packet content. Hence the random background may unconsciously contain patterns.

4.2.2 Model II

To evaluate the performance of algorithms in a real intense attack, a trace from the Capture-the-Flag contest held at Defcon9 was adopted as the input traffic in Model II. The Defcon Capture-the-Flag contest is the largest security hacking game, in which competitors try to break into the servers of others while protecting their own servers, each hiding several security holes [26].

4.2.3 Model III

Model III uses a real 2-hour trace as the input traffic, and the more recent Snort rules R2 as the pattern set jP_j . This real trace recorded all IP packets in a laboratory of Providence University for 2 hours. The laboratory has an FTP server, a Web server, and three PCs running several network application clients.

Table 5 lists the statistics of the traffic traces used in Model II and Model III, where the values are measured by traffic analysis tools: tcpstat and tcptrace.

TABLE 5

The Statistics of the Traffic Traces

Statistics	Model II	Model III
Average Packet Size (Byte)	467.71	896.1
The Standard Deviation of the Size of each Packet (Byte)	651.06	690.99
Data Transmission Rate (Kbps)	254.13	280.03
Number of Packets per second	69.55	40

4.3 Memory Requirements

For fast lookup and matching, the lookup information and patterns are usually saved in the memory using a tabular structure. Therefore, the memory requirements are estimated according to the number of entries. Since all algorithms need to keep the pattern content in the (external) memory, this section only discusses the extra memory requirement for the tables of each algorithm. In the simulations, the numbers of characters in the clustering pivots (B1 and B2) were both assumed to be 1. Because the H1 of EHMA is a direct index table, the cache memory

Space (MI) of EHMA comprises jAj entries. Based on GFGS and CBS, the number of entries in H2 is the total number of possible clusters (plus a small memory pool). Since the domain of possible pivot pairs is $F \times A$, the external memory space for H2 (ME) of EHMA is $O(jF_j \times jAj)$. HMA has the same memory requirement as EHMA. The shift table of WM is also a direct hash table. The gram size of WM (block size B) was 3 in the simulations, so the shift table of WM had jAj entries. The grouped skip table of WM-PH used in the simulations was a direct prefix hash table with a prefix length of three characters. Therefore, the skip table of WM-PH comprises jAj entries. Every pattern in the BMH has its own skip table of jAj entries, so that the table of BMH has $jP_j \times jAj$ entries. Because each skip table of BMH (for one pattern) is small enough to be loaded into the local memory, for fairness, a cache memory space was allocated to lower the

number of external memory accesses. The BMH-O is the original BMH with no local cache and assesses the latency penalty. Notably, WM-PH, AC-C, and BMH-O also require cache memory to store the skip value or one state during the matching process. Table 4 lists the memory requirements of EHMA, HMA, WM, WM-PH, BMH, and AC-C. The scale relation of the parameters is $jF_j < jAj$ $jP_j < S$ jAj .

In the simulations using Model I, when jP_j is 1,200, the H1 and H2 of EHMA needs 256 and 19,712 entries, respectively (about 768 bytes on-chip memory and 38.5-Kbyte external memory, including the shared memory pool); HMA has the same number of entries as EHMA but needs smaller entry size as HMA has no shift field; the table of WM needs more than 16 million entries (16-Mbyte external memory, in the case without using an additional prefix table); the table size of WM-PH is the same as that of WM; BMH and BMH-O need more than 300,000 entries (300-Kbyte external memory); and AC-C needs 10,731 states (461 Kbytes with each node of 44 bytes). The memory size of all algorithms listed previously excludes pattern content. Obviously, the required memory space of EHMA is quite small.

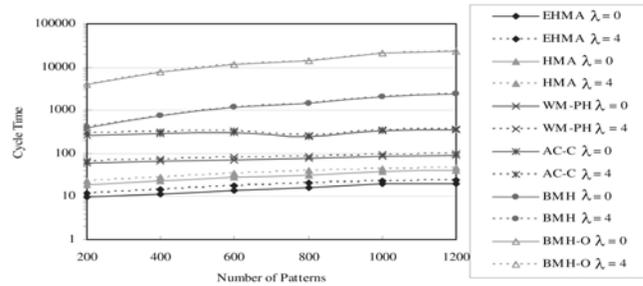


Fig. 7. The average matching time () versus the number of patterns (jP_j), using Model I with $A = 0$ and $A = 4$, where $wE = 100$.

4.4 Results and Discussion

The minimum pattern length of the feeding patterns in Figs. 7, 8, 9, and 10 is only one character, i.e., $M = 1$. Because the minimum pattern length of WM is restricted to be larger than the gram size, in this case three characters, WM is not compared in these figures. In Figs. 7,8, 9, and 10, the results labeling EHMA in the following simulations use the sampling window with parameters $W = m = jP_j$, which means that each pattern is sampled in its entirety.

Fig. 7 compares the average matching time of EHMA, HMA, WM-PH, AC-C, BMH, and BMH-O using Model I with different attack loads $A = 0$ and $A = 4$, respectively. It also shows the impact of the number of patterns (jP_j) on the matching time. Simulation results reveal that EHMA outperforms others even when jP_j and A increase. EHMA has slightly higher growth rate than WM-PH, because it has a much smaller table size. WM-PH gains performance by having a large direct index table. Notably, the matching time of the original AC using basic structure is independent from jP_j and A . The curves of AC-C increase with jP_j and A owing to the popsum used in the AC-C algorithm. The increasing jP_j makes the matching time of BMH (BMH-O) rise steeply, because the BMH is originally a single-pattern matching algorithm that simply executes iteratively for multipattern matching.

The case $A = 0$ means that the traffic has no malicious packets. In this case, the proposed EHMA needs only 9.5- 19.9 cycles per character on average, which is about 0.9,

3.3-5.3, 16.3-26.8, 40-117, and 408-1,161 times less than the matching time of HMA, WM-PH, AC-C, BMH, and BMH-O, respectively, under various pattern set sizes. We can say that EHMA is very appropriate for network equipment, because generally most packets are innocent ($A = 0$). The time available for the detection engine to process the malicious packets rises as the innocent packets are processed more quickly.

When $A = 4$, then the systems are under heavy attack, and the traffic contains many monitored patterns. In this situation, the matching time of EHMA is about 0.89-0.94, 3.1-4.5, 14.1-24.9, 33.2-96.4, and 335-957 times less than that of HMA, WM-PH, AC-C, BMH, and BMH-O, respectively. Additionally, the performance of EHMA is quite stable, since rises only slightly as A or jP_j rises.

Fig. 8 displays the proportion of I to M and M to I , respectively, for all approaches using Model I with $jP_j = 1, 200$, where Fig. 8a shows the results under $A = 0$, and Fig. 9b shows the results under $A = 4$. In Fig. 8, the upper part of the bar is I and the lower part of the bar is M . The results show that the I of EHMA is close to that of HMA and WM-PH, but M of EHMA is much less than others. The proportion of M to I of BMH seems smaller than others, because the whole skip table of a pattern is idealistically assumed to be loaded within one external memory access and kept in the cache during the matching process for each pattern. Because AC-C compresses the data structure of the state machine, it requires more time to derive the next state pointer. Therefore, AC-C does not have the smallest I . Simulation results show that the I does not significantly rise with A in any of the experiments, because each algorithm has already tried to reduce the computation load (I). However, M dominates the overall matching cost. This reveals that the number of external memory accesses is the bottleneck of almost all algorithms. This result also reflects our opinion mentioned previously that the essential issue in designing a high-speed detection engine is to reduce the number of required external memory accesses.

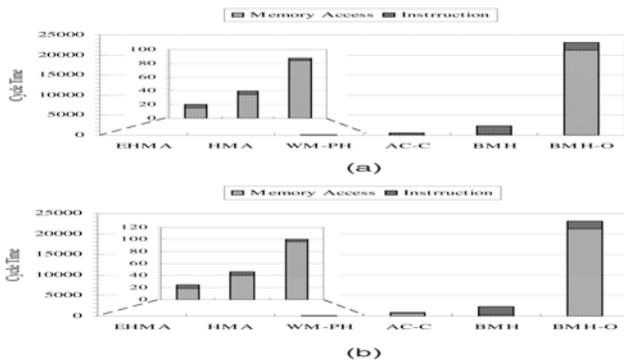


Fig. 8. The proportion of I to M and M to I using Model I with $jP_j = 1,200$ and $wE = 100$. (a) $A = 0$ and $jP_j = 1,200$. (b) $A = 4$ and $jP_j = 1,200$.

Fig. 9 compares the average number of external memory accesses per character (E) of the state-of-the-art pattern matching algorithms. The figure shows that the E of EHMA is

only 0.06-0.19, which is much smaller than others. In other words, EHMA can successfully filter out about 94 percent payloads when $jP_j = 200$ and 81 percent when $jP_j = 1, 200$, requiring no external memory accesses and string comparisons. The E of EHMA rises only slightly with rising A . The increasing rate of E is slightly higher in EHMA than in WM-PH when jP_j rises, because EMHA has much smaller table size than WM-PH. Since BMH is based on the single-pattern matching algorithm, its E is proportional to jP_j . Consequently, the hierarchical matching along with the safety shift strategy is highly effective in reducing the memory latency.

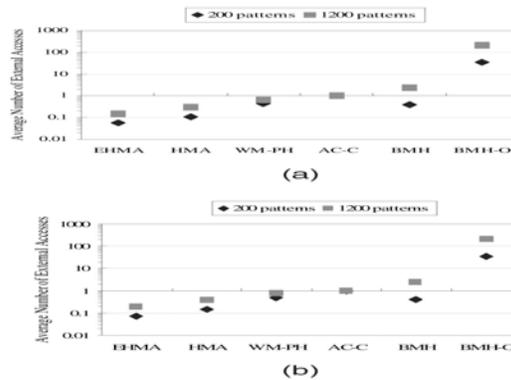


Fig.9. The comparisons of average number of external memory accesses (E) using Model I with $wE = 100$. (a) $A = 0$. (b) $A = 4$.

Figs. 9 and 10 adopted Model II as a real-life network environment under intense attack to evaluate the performance of the state-of-the-art algorithms. Since different implementation systems may have different external memory costs (wE), Fig. 9 illustrates two results with $wE = 100$ and $wE = 10$, respectively. To lower the impact of wE on an algorithm, a very small value of wE is adopted in Fig. 10b. The results in Fig. 9 indicate that EHMA significantly outperforms others in both cases of small and large pattern set sizes even in the intense attack. EHMA still performs better than others even when the penalty on the external memory access (wE) is reduced (as shown in Fig. 9b). Comparing EHMA with HMA in Figs. 6,7,8, and 9 reveals that the proposed safety shift strategy significantly reduces the number of external memory accesses and thus improves the matching performance.

The minimum length of Snort patterns is one character. However, some detection systems, such as virus detection systems, have larger minimum pattern lengths. The performance of matching algorithms with long minimum pattern lengths was examined using Model II, including only the patterns with lengths greater than 10 ($M = 10$) from Snort patterns, as drawn in Fig. 10. Since the number of patterns whose length is larger than 10 characters in R1 is around 500, Fig. 10 shows the cases of $jP_j = 200$ and $jP_j = 500$, respectively. Fig.10a shows the average processing time (T); Fig. 10b shows the memory requirement of the fast index/hash tables, excluding the memory for pattern contents. Since here M is larger than the gram size of WM, which is three as mentioned before, the performance of WM is compared here. The result labeling EHMA ($W = 5$) is the case using EHMA algorithm with $m = M = 10$ and $W = 5$. Recall that the sampling window of EHMA is the entire pattern

content, that is, $m = M = jpi j$. To observe the performance of WM and WM-PH with smaller hash tables, Fig. 10 also displays two additional cases with block size of two characters, WM(B = 2) and WM-PH(B = 2).

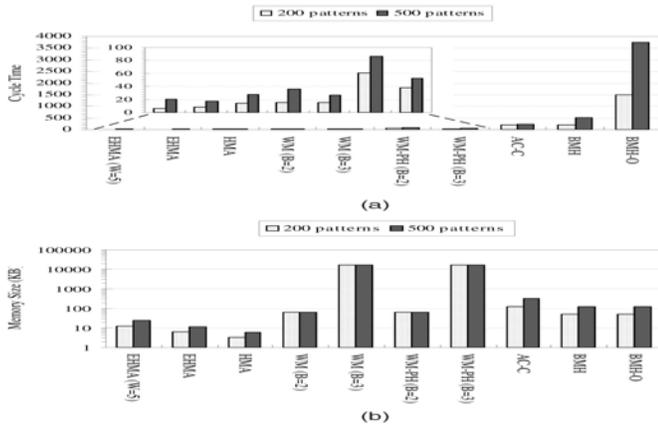


Fig. 10. The costs versus the number of patterns ($jP j$) using Model II, $wE = 100$ and $M = 10$. (a) Average matching time. (b) Extra memory requirement.

Before discussing the simulation results of Fig. 10, Table 6 presents the effect of the size of sampling window (W) on the performance of EHMA in terms of the average shift values of H^1 and H^2 , the size of the set of frequent-common grams ($jF j$) derived from GFGS, the average number of actual shifts, and the average number of external memory accesses, using the same traffic model as in Fig. 10. Table 6 shows that the number of candidate common grams increases with increasing W , resulting in smaller $jF j$. The average number of H^1 : shift and H^2 : shift increases when W decreases. Since the traffic spectrum is not normally distributed, the actual average number of shifts during matching process is not the same as the average of H^1 : shift and H^2 : shift. However, the trend is the same. E is effected by both $jF j$ and the actual average shift.

Fig. 10a shows that EHMA($W = 5$) outperforms EHMA and others when $jP j = 200$; while EHMA performs better than EHMA($W = 5$) and others when $jP j = 500$. Therefore, reducing $jF j$ becomes more important than increasing the average number of shift values when $jP j$ is large. Since all algorithms need a copy of the pattern contents, Fig. 10b only displays the extra memory requirement of every algorithm for the index/hash tables. Fig. 10 b shows that the required memory of EHMA is only slightly larger than that of HMA but much smaller than that of others. The required memory of EHMA grows moderately with $jP j$. The memory of EHMA($W = 5$) is greater than that of EHMA due to the larger $jF j$. As shown in Fig. 10, EHMA is highly effective in reducing the required external memory, providing efficient performance even in the virus-detection-like model.

Fig. 11 uses Model III as real-life normal traffic to show the performance of the algorithms. Meanwhile, to demonstrate the effect of the rising number of patterns on the matching performance, a more recent Snort rule set R2 of about 5,000 patterns are used in Model III. Fig. 12 shows that EHMA performs better than others even when the pattern set is very large. The matching time of EHMA only moderately increases with the rising $jP j$.

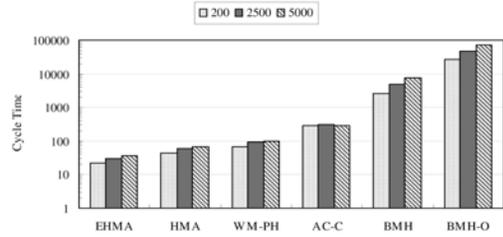


Fig. 11. The average matching time () versus the number of patterns ($jP j$) using Model III, $wE = 100$.

$ P $	200				500			
	EHMA	EHMA (W=7)	EHMA (W=5)	EHMA (W=3)	EHMA (W=7)	EHMA (W=5)	EHMA (W=3)	
H^1 shift	0.94	2.71	3.66	4.74	0.91	1.86	2.02	2.49
H^2 shift	1.99	4.89	6.79	8.71	1.99	4.84	6.72	8.65
$ F $	13	20	25	39	23	33	47	65
Average Shift	1.5	1.74	1.79	1.84	1.49	1.68	1.74	1.8
E	0.0377	0.0441	0.0431	0.0434	0.1243	0.16	0.1635	0.2512

TABLE 6

The Impact of the Size of Sampling Window (W) in the Shift Values of Tables, $jF j$, Actual Matching Shifts, and E Using Model II

5 CONCLUSIONS

The increasing variety of network applications and stakes held by various users are creating a strong demand for fast in-depth packet inspection. The most important component of in-depth packet inspection is an efficient multipattern matching algorithm. This study proposes a novel EHMA for packet inspection. EHMA applies the frequent-common grams obtained by the proposed GFGS to narrow the searching scope and to quickly filter out the innocent packets. The matching process then focuses only on the most suspected packets. EHMA concentrates the patterns into a small on-chip table and performs simple and fast checks. Additionally, EHMA uses the frequency-based bad gram heuristic to speed up the scanning process. The hierarchical matching significantly reduces the average number of external memory accesses to only 6 percent to 19 percent, thus improving the matching performance. The required memory of EHMA is only about 40 Kbytes in addition to the pattern contents of Snort rules. Particularly, EHMA is very simple and can be easily implemented in both software-based and hardware-based platforms. This study also discusses and evaluates current multi pattern matching algorithms for NIDSs. Simulation results show that EHMA performs about 0.89-1,161 times better than others. Even under real-life intense attack, EHMA significantly outperforms others. EHMA also works well for the systems with larger minimum pattern size, such as virus detection systems. In conclusion, EHMA facilitates the creation of efficient and cost-effective pattern detection engines for packet inspection.

REFERENCES

[1] Snort, <http://www.snort.org>, 2008.
 [2] S. Antonatos, K.G. Anagnostakis, and E.P. Markatos, "Generating Realistic Workloads for Network Intrusion Detection Systems," Proc. Fourth Int'l ACM Workshop Software and Performance (WOSP),

2004.

[3] R.N. Horspool, "Practical Fast Searching in Strings," Software

Practice and Experience, vol. 10, no. 6, pp. 501-506, 1980.

[4] A.V. Aho and M.J. Corasick, "Efficient String Matching: An Aid to Bibliographic Search," *Comm. ACM*, vol. 18, no. 6, pp. 330-340, June 1975.

[5] M. Fisk and G. Varghese, "Fast Content-Based Packet Handling for Intrusion Detection," UCSD Technical Report CS2001-0670, May 2001.

[6] O. Erdogan and P. Cao, "Hash-AV: Fast Virus Signature Scanning by Cache-Resident Filters," *Proc. IEEE Global Telecomm. Conf. (GLOBECOM '05)*, Nov. 2005.

[7] S. Lakshmanamurthy, K.-Y. Liu, Y. Pun, L. Huston, and U. Naik, "Network Processor Performance Analysis Methodology," *Intel Technology J.*, vol. 6, Aug. 2002.

[8] N. Tuck, T. Sherwood, B. Calder, and G. Varghese, "Deterministic Memory-Efficient String Matching Algorithms for Intrusion Detection," *Proc. IEEE INFOCOM '04*, Mar. 2004.

[9] T.-F. Sheu, N.-F. Huang, and H.-P. Lee, "A Novel Hierarchical Matching Algorithm for Intrusion Detection Systems," *Proc. IEEE Global Telecomm. Conf. (GLOBECOM '05)*, Nov. 2005.

[10] S. Wu and U. Manber, "A Fast Algorithm for Multi-Pattern Searching," Technical Report TR94-17, Dept. Computer Science, Univ. of Arizona, May 1994.

[11] E. Markatos, S. Antonatos, M. Polychronakis, and K. Anagnostakis, "Exclusion-Based Signature Matching for Intrusion Detection," *Proc. IASTED Int'l Conf. Comm. and Computer Networks (CCN '02)*, Oct. 2002.

[12] R.-T. Liu, N.-F. Huang, C.-H. Chen, and C.-N. Kao, "A Fast String Matching Algorithm for Network Processor-Based Intrusion Detection System," *ACM Trans. Embedded Computing Systems*, vol. 3, no. 3, Aug. 2004.

[13] R.S. Boyer and J.S. Moor, "A Fast String Searching Algorithm," *Comm. ACM*, vol. 20, no. 10, pp. 762-772, Oct. 1977.

[14] T.-F. Sheu, N.-F. Huang, and H.-P. Lee, "A Time- and Memory- Efficient String Matching Algorithm for Intrusion Detection Systems," *Proc. IEEE Global Telecomm. Conf. (GLOBECOM '06)*, Nov. 2006.

[15] C.J. Coit, S. Staniford, and J. McAlerney, "Towards Faster String Matching for Intrusion Detection or Exceeding the Speed of Snort," *Proc. Second DARPA Information Survivability Conf. and Exposition (DISCEX)*, 2001.

[16] S. Antonatos, M. Polychronakis, P. Akritidis, K.G. Anagnostakis, and E.P. Markatos, "Piranha: Fast and Memory-Efficient Pattern Matching for Intrusion Detection," *Proc. 20th IFIP Int'l Information Security Conf. (SEC '05)*, May 2005.

[17] S. Li, J. Torresen, and O. Soraasen, "Exploiting Reconfigurable Hardware for Network Security," *Proc. 11th Ann. IEEE Symp. Field-Programmable Custom Computing Machines (FCCM)*, 2003.

[18] S. Kim and Y. Kim, "A Fast Multiple String-Pattern Matching Algorithm," *Proc. 17th AoM/IAoM Int'l Conf. Computer Science*, Aug. 1999.

[19] S. Dharmapurikar, P. Krishnamurthy, T. Sproull, and J. Lockwood, "Deep Packet Inspection Using Parallel Bloom Filters," *Proc. 11th Symp. High Performance Interconnects*, Aug. 2003.

[20] H. Lu, K. Zheng, B. Liu, X. Zhang, and Y. Liu, "A Memory-Efficient Parallel String Matching Architecture for High-Speed Intrusion Detection," *IEEE J. Selected Area in Comm.*, vol. 24, no. 10, Oct. 2006.

[21] S. Dharmapurikar and J. Lockwood, "Fast and Scalable Pattern Matching for Network Intrusion Detection Systems," *IEEE J. Selected Area in Comm.*, vol. 24, no. 10, Oct. 2006.

[22] Vitesse Network Processors, <http://www.vitesse.com>, 2008.

[23] Intel Network Processors, <http://www.intel.com/design/network/products/npfamily/index.htm>, 2008.

[24] C. Kruegel, F. Valeur, G. Vigna, and R. Kemmerer, "Stateful Intrusion Detection for High-Speed Networks," *Proc. IEEE Symp. Security and Privacy (SP '02)*, May 2002.

[25] M. Handley, V. Paxson, and C. Kreibich, "Network Intrusion Detection: Evasion, Traffic Normalization, and End-to-End Protocol Semantics," *Proc. Ninth USENIX Security Symp.*, 2000.

[26] C. Cowan, "Defcon Capture the Flag: Defending Vulnerable Code from Intense Attack," *Proc. DARPA Information Survivability Conf. and Exposition (DISCEX III '03)*, Apr. 2003.