

# Exscind: A Faster Pattern Matching For Intrusion Detection Using Exclusion and Inclusion Filters

1

**Monther Aldwairi and Duaa Alansari**

**Seventh International Conference on  
Next Generation Web Services  
Practices (NWeSP'11)**

**19-21, October 2011  
Salamanca, Spain**

# Introduction to IDS

2

- **IDS inspects ingress/egress traffic for attack signatures/behaviour**
  1. Drops the packet
  2. Alerts the system security administrator
  3. Logs malicious activity
- **Anomaly-based IDS**
  - Employs machine learning to profile, then classify traffic into either normal or abnormal behavior
  - + Detects new attacks
    - High false positives
    - Longer detection time
- **Signature-based IDS (Misuse)**
  - Employs exact pattern matching to look for specific signatures within packets
  - + Faster than anomaly but not enough for inline deployment
  - + Low false positives
    - But unable to detect new attacks
    - Manual signatures generation
  - + - Accuracy depends on signatures accuracy

# The Problem

3

```
alert tcp $EXTERNAL_NET any -> $HTTP_SERVERS 80
(msg:"WEB-IIS CodeRed v2 root.exe access";
content:"/root.exe"; sid:1256)
```

- Snort rules or attack signatures

Year	2003 <sup>[2]</sup>	2011
Rules	3000	11,416
Signatures	<b>1542</b>	<b>9945</b>

- **30%-60%** of Snort processing time → pattern matching <sup>[1]</sup>
- Moore's law! Internet speeds double every eighteen months
- Internet traffic is doubling every six months <sup>[4]</sup>
- Better Pattern Matching → Speed, Memory, Configurability

# Contributions

4

1. Add a **modified** BF in front of modified WM to skip as much clean traffic as possible
  - Program ONLY the prefix (first n-gram) and store SIDs
  - Produce a probable matches list
  - Provide an index for the first match position
2. **Modify** the WM
  - Search ONLY the probable matches list
  - Start from the first index of probable match
3. Support all ASCII codes
  - Extended, special characters and the Null character
4. Vary the prefix 4-to -10

# Wu-Manber Algorithm

5

## • Preprocessing Stage

- Find the minimum pattern length ( $m$ )
- Choose Block size  $B$  (2 or 3)
- Builds three hash tables: SHIFT table, HASH table, and PREFIX table
- For any Block  $x$ , does it appear in any pattern?

$x \in patterns?$	yes	No
$SHIFT [i]$	$m-q$	$m-B+1$

- ✦  $q$  is the right most occurrence of  $x$  in a pattern
- HASH and PREFIX tables contain linked lists of signatures with zero shift

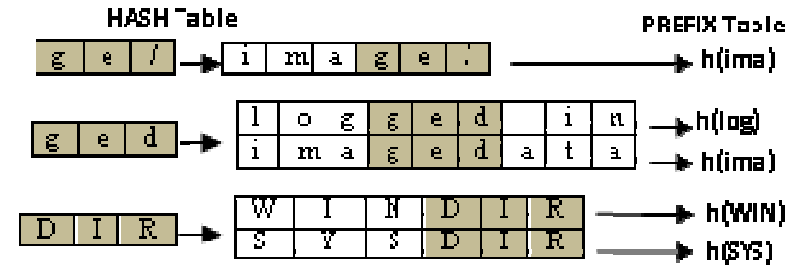
## • Search Stage

- Slide a window,  $w$ , over the packet  $\in$  SHIFT Table
- Hash last  $B$  chars of sliding window and access SHIFT table
  - ✦  $SHIFT[i] > \text{zero}$  slide the window and repeat
  - ✦ If  $SHIFT[i] = \text{zero}$ ,  $\rightarrow$  search HASH and PREFIX tables to verify matches

# Wu-Manber Example

6

Signature	SID	Block	Shift	Block	Shift
"image/"	2706	mai	3	WIN	3
"logged in"	162	mag	2	IND	2
"imagedata"	12280	age	1	NDI	1
"WINDIR"	3010	ge/	0	DIR	0
"SYSDIR"	3011	log	3	SYS	3
		ogg	2	YSD	2
		gge	1	SDI	1
		ged	0	others	4



Step	z t i m a g e / l k S Y S D I R o	shift	output
1	z t i <b>m a g</b>	2	
2	i m a <b>g e /</b>	0	image/
3	m a g <b>e / l</b>	4	
4	/ l k <b>S Y S</b>	3	
5	S Y S <b>D I R</b>	0	SYSDIR
6	Y S D <b>I R o</b>	4	

# Related Work

7

- **Mike Fisk et al.** proposed a new hybrid system for pattern matching for IDS. This system provide a combination of Boyer-Moore algorithm and Aho-Corasick algorithm as a hybrid system
- **Kostas Anagnostakis et al.** proposed a new exclusion filter for IDS called  $E^2xB$ . This stage marks the corresponding cell on the (256-cell) map for each character appearing in the text. Then, the decision that the text contains a character will be only if the corresponding cell of the map is marked with the num of the packet
- **Sarang Dharmapurikar et al.** proposed a new hardware-based technique for Bloom filters. They divided the signatures into groups according to their length and stored each group in a unique Bloom filter. Each Bloom filter used an analyzer for probing detected strings to determine whether it is a false positive or not
- **Haoyu Song et al.** proposed Extended Bloom Filter (EBF) to implement multi-pattern signature matching. They introduced a new technique to support long signature matching to minimize the maintained signature lengths. They embedded multi-port memories in FPGAs

# Outline

8

- Introduction to Intrusion Detection
- The Problem
- Related Work
- **Exscind: The Big Picture**
- Bloom Filters
- Exscind Example
- Experimental Evaluation
- Conclusions



# Exscind the Big picture

9

- We propose Exscind: meaning to exclude from the union
  - Exscind is an exclusion-inclusion intrusion detection
  - It can exclude as many clean packets as possible
  - Modified Exclusion-Inclusion Bloom filter using only 4-gram
    - ✦ Produces a probable matches list
    - ✦ Provides an index for the first match position
  - Modifies the WM
    - ✦ Search ONLY the probable matches list
    - ✦ Start from the first index of probable match
    - ✦ Support all ASCII codes
    - ✦ Extended, special characters and the Null character

# Excind the Big picture

10

Incoming packets



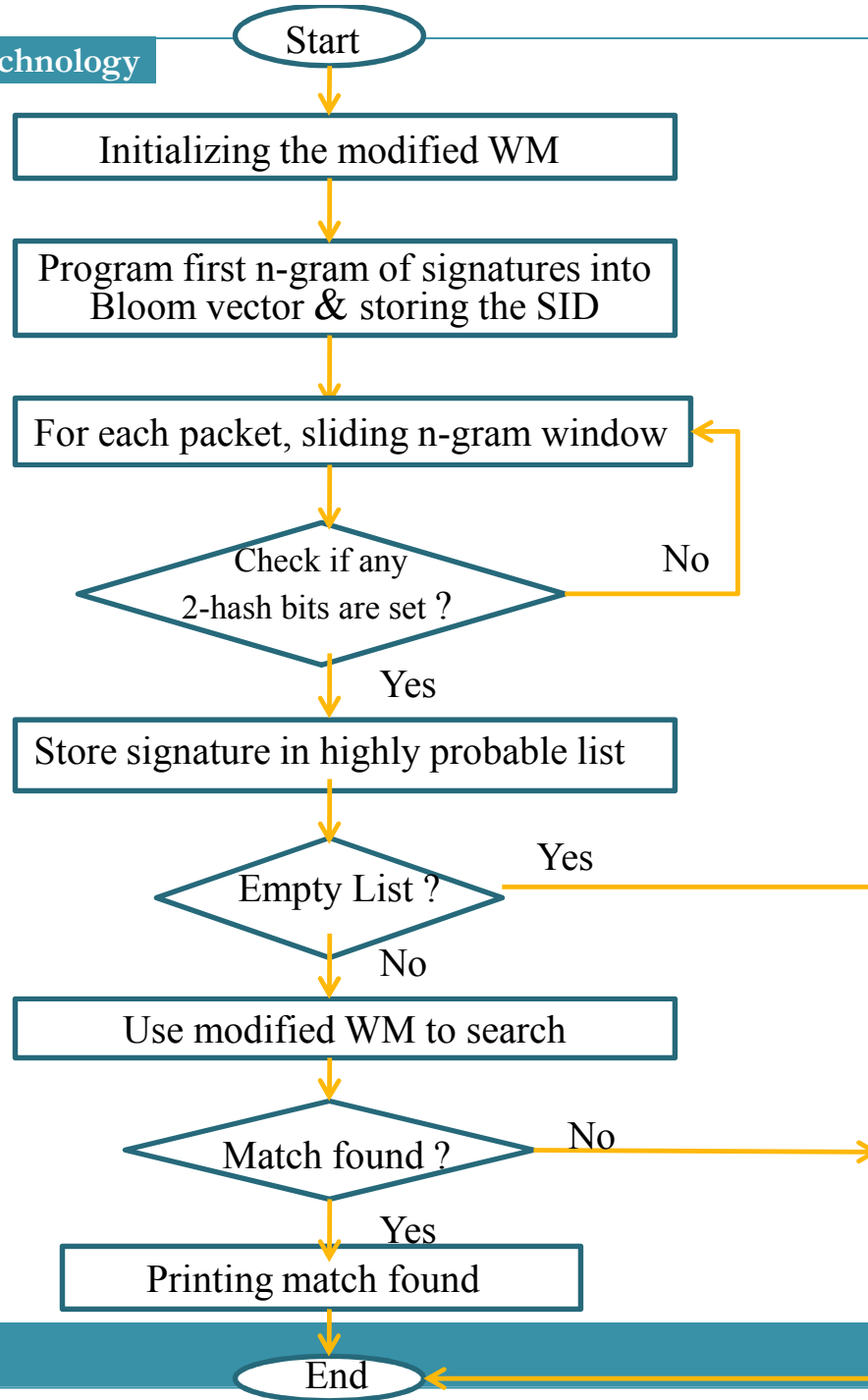
Exclusion-Inclusion Filter  
(Modified Bloom Filter)



Probable Matching  
First Match Index

Modified Wu-Manber (MWM)

1. Search for the probable matches
2. Search from the first probable match index



# Bloom Filters

12

- Exclusion filter through approximate membership queries on a set of strings
  - Bloom Vector: array vector of  $m$  bits
  - Programming: use  $k$  hash functions on signatures and set the corresponding bits
  - Querying: Compute the same  $k$  hash functions on the packet, then check corresponding hash bits?
    - ✦ If at least one bit is not set, 100% certainty the signature does not belong
    - ✦ If all  $k$  bits are set, the signature might belong to the set with a probability
  - Assume the signature “**cmd.exe**” hashes to **2** and **8**
  - Set bits **2** and **8** in a 10-bit Bloom filter



# Exscind Example

13

1. Initialize modified WM with all Snort signatures
2. Program the Bloom vector with the 4-gram prefix of all Snort signatures
3. Store the SID of each 4-gram corresponding to set bits

Serial	SID	Signature	Hash Bits
3135	1687	dba_tables	5, 16
3136	1688	user_tablespace	3, 7
3137	1689	sys.all_users	9, 14

			1688		1687		1688		1689					1689		1687				
0	0	0	1	0	1	0	1	0	1	0	0	0	0	1	0	1	0	0	0	0
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20

Signatures Bloom Vector

# Exscind Example

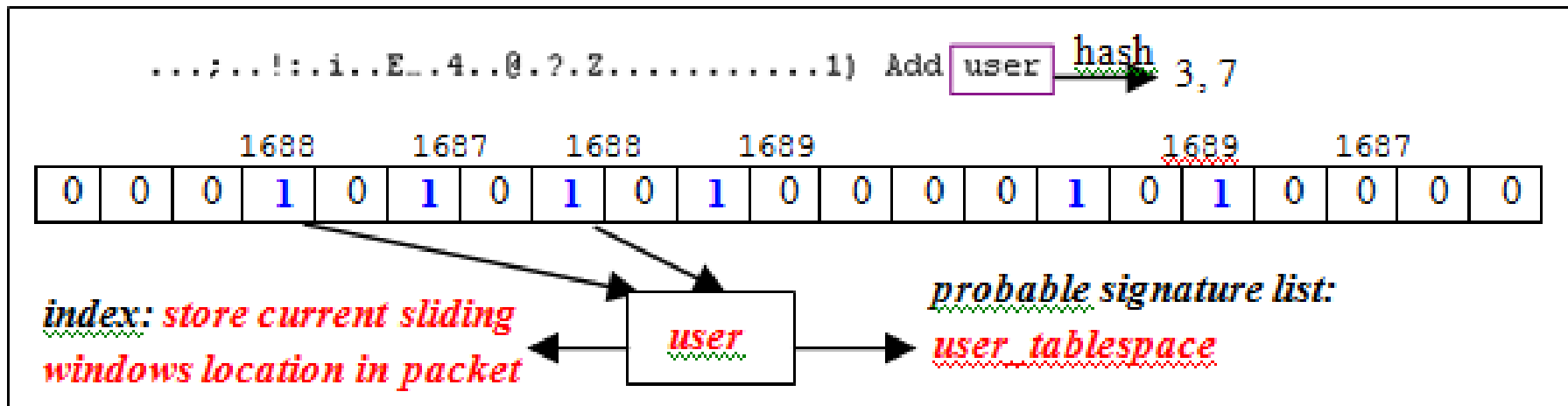
14

4. Slide a 4-gram window throughout the packet
5. Computes the same two hash functions

Sliding Window Moving Throughout the Packet	Hash Bits
<span style="border: 1px solid blue; padding: 2px;">...</span> ;...!:i..E...4..@.?..Z.....1) Add user	0, 8
... <span style="border: 1px solid red; padding: 2px;">;...!</span> ;...!:i..E...4..@.?..Z.....1) Add user	4, 12
... <span style="border: 1px solid green; padding: 2px;">;...!</span> ;...!:i..E...4..@.?..Z.....1) Add user	10, 8
... <span style="border: 1px solid purple; padding: 2px;">;...!</span> ;...!:i..E...4..@.?..Z.....1) Add user	14, 17
.	.
.	.
...;...!:i..E...4..@.?..Z.....1) Add <span style="border: 1px solid purple; padding: 2px;">user</span>	3, 7

# Exscind Example

4. Check the hash bits of each window against the signatures Bloom vector.
5. If they are set, this means that there is a probable match
6. The system retrieves the stored SIDs and adds the signatures to the probable signatures list.



# Exscind Example

16

8. Exscind checks the probable signatures list
  - a. If empty → packet is clean → 100% certainty it has no attack signatures → skip pattern matching
  - b. If not empty → packet is suspicious → modified WM → starting from first position of probable match → search only the probable signatures
  - c. If any match occurs within the packet, the system considers malicious packet and records the number of matches occurrences along with the index of each occurrence



# Exscind Savings and Overhead

17

- **Exscind reduces execution time**
  - Exclusion-Inclusion skips pattern matching for clean packets
  - MWM searches suspicious packets ONLY for small probable signatures matching list
  - MWM search starts from the index for the first probable match
- **Exscind over head on memory and time**
  - Time for BF programming added to preprocessing
  - Time for BF querying
  - Memory space for Bloom vector and the of signatures SIDs

# Outline

18

- Introduction to Intrusion Detection
- The Problem
- Related Work
- Exscind: The Big Picture
- Bloom Filters
- Exscind Example
- **Experimental Evaluation**
- Conclusions

# Experimental Setup and Metrics

19

- **Metrics**

1. Average execution time
2. Average memory usage
3. WMSC (Wu-Manber Skip Counter) means the number of clean packets skipped without pattern matching

- **Experimental Environment**

- Workstation: 3 GHz Intel Core™ 2, 4 GB of RAM
- execution time → ExecutionStopWatch class
- memory usage → Garbage Collection class
- 20 experiments average numbers reported

# Export Packet Contents

20

- Wireshark Network Protocol Analyzer v1.4.4

No.	Time	Source	Destination	Protocol	Info
1	0.000000	10.31.3.110	10.31.8.2	TCP	53131 > unicall [PSH, ACK] Seq=1 Ack=1 win=33
2	0.000113	10.31.3.110	10.31.8.2	TCP	53131 > unicall [FIN, ACK] Seq=72 Ack=1 win=3
3	0.000380	10.31.3.110	10.31.9.2	TCP	53132 > unicall [SYN] Seq=0 win=65535 Len=0 M
4	0.000584	10.31.6.100	10.31.8.2	TCP	57876 > 7331 [PSH, ACK] Seq=1 Ack=1 win=183 L
5	0.000793	10.31.6.100	10.31.8.2	TCP	57876 > 7331 [FIN, ACK] Seq=93 Ack=1 win=183
6	0.001649	10.31.9.2	10.31.3.110	TCP	unicall > 53132 [SYN, ACK] Seq=0 Ack=1 win=65
7	0.001720	10.31.8.2	10.31.3.110	TCP	unicall > 53131 [ACK] Seq=1 Ack=73 win=8317 L
8	0.002123	10.31.3.110	10.31.9.2	TCP	53132 > unicall [ACK] Seq=1 Ack=1 win=66608 L
9	0.002160	10.31.6.100	10.31.9.2	TCP	43918 > 7331 [SYN] Seq=0 win=5840 Len=0 MSS=1
10	0.002259	10.31.3.110	10.31.9.2	TCP	53132 > unicall [PSH, ACK] Seq=1 Ack=1 win=66
11	0.002325	10.31.3.110	10.31.9.2	TCP	53132 > unicall [FIN, ACK] Seq=72 Ack=1 win=6

Frame 1: 137 bytes on wire (1096 bits), 137 bytes captured (1096 bits)					
Ethernet II, Src: Apple_06:f4:fc (00:26:4a:06:f4:fc), Dst: IntelCor_3a:e2:69 (00:1b:21:3a:e2:69)					
Internet Protocol, Src: 10.31.3.110 (10.31.3.110), Dst: 10.31.8.2 (10.31.8.2)					
Transmission Control Protocol, Src Port: 53131 (53131), Dst Port: unicall (4343), Seq: 1, Ack: 1,					
Data (71 bytes)					

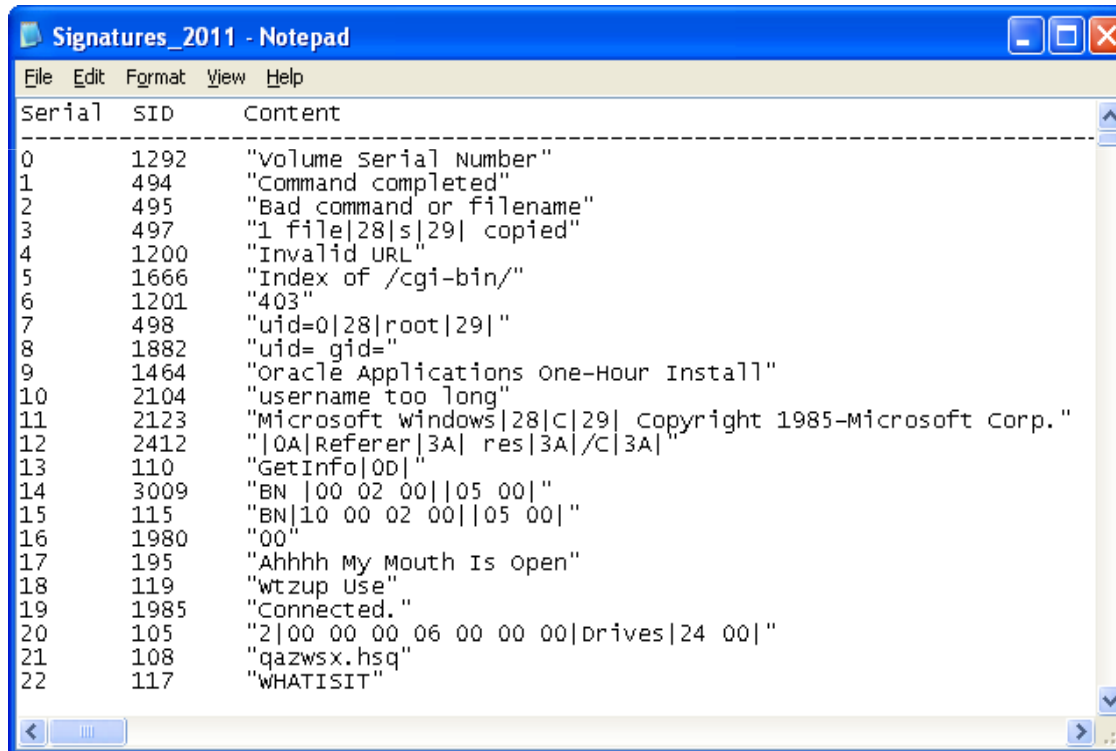
0000	00 1b 21 3a e2 69 00 26	4a 06 f4 fc 08 00 45 00	...!.i.& J.....E.
0010	00 7b 57 54 40 00 40 06	c3 7b 0a 1f 03 6e 0a 1f	..{wT@.@. .{...n..
0020	08 02 cf 8b 10 f7 bb 1d	25 ae 4e 7c 36 6b 80 18	..... %N 6k..
0030	82 18 9d 48 00 00 01 01	08 0a 0c 73 ba bb 44 cc	...H.... ...s..D.
0040	1e bf 68 4c 1c dc fc 45	14 08 76 30 f2 6e 86 af	..hL...E ..v0.n..
0050	75 b1 a4 c1 c1 bf 9e 26	08 d4 d5 00 a3 fc 78 ba	u.....& .....x.
0060	07 75 25 31 35 31 24 31	30 30 68 6e 25 34 38 24	.u%151\$1 00hn%48\$
0070	31 30 35 68 6e 25 31 34	37 24 31 30 39 68 6e 25	105hn%14 7\$109hn%
0080	31 33 32 24 31 31 38 68	6e	132\$118h n



# Signatures Extraction

22

- Snort 2.9.0.4 rules released in March 2011
- 56 rules files → C# to extract SIDs, content and uricontent  
→ concatenate same rule signatures with zero space

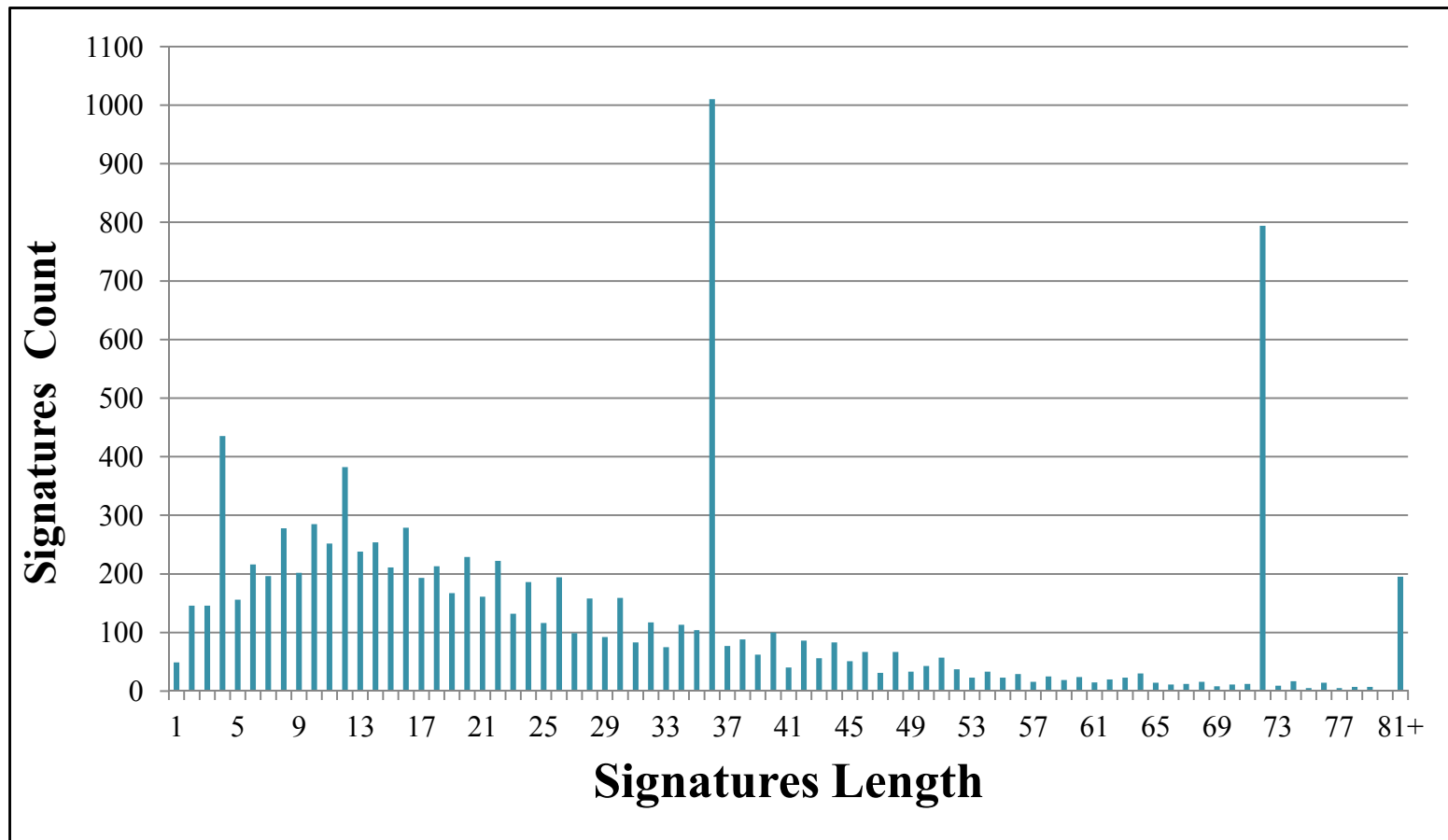


```
Signatures_2011 - Notepad
File Edit Format View Help
Serial SID Content
-----
0 1292 "Volume Serial Number"
1 494 "Command completed"
2 495 "Bad command or filename"
3 497 "1 file|28|s|29| copied"
4 1200 "Invalid URL"
5 1666 "Index of /cgi-bin/"
6 1201 "403"
7 498 "uid=0|28|root|29|"
8 1882 "uid= gid="
9 1464 "Oracle Applications One-Hour Install"
10 2104 "username too long"
11 2123 "Microsoft windows|28|c|29| copyright 1985-Microsoft Corp."
12 2412 "|0A|Referer|3A| res|3A|/C|3A|"
13 110 "GetInfo|0D|"
14 3009 "BN |00 02 00||05 00|"
15 115 "BN|10 00 02 00||05 00|"
16 1980 "00"
17 195 "Ahhhh My Mouth Is open"
18 119 "wtzup Use"
19 1985 "Connected."
20 105 "2|00 00 00 06 00 00 00|drives|24 00|"
21 108 "qazwsx.hsq"
22 117 "WHATISIT"
```

# Signatures Analysis

23

- 9,945 signatures



# Malicious Traffic Analysis

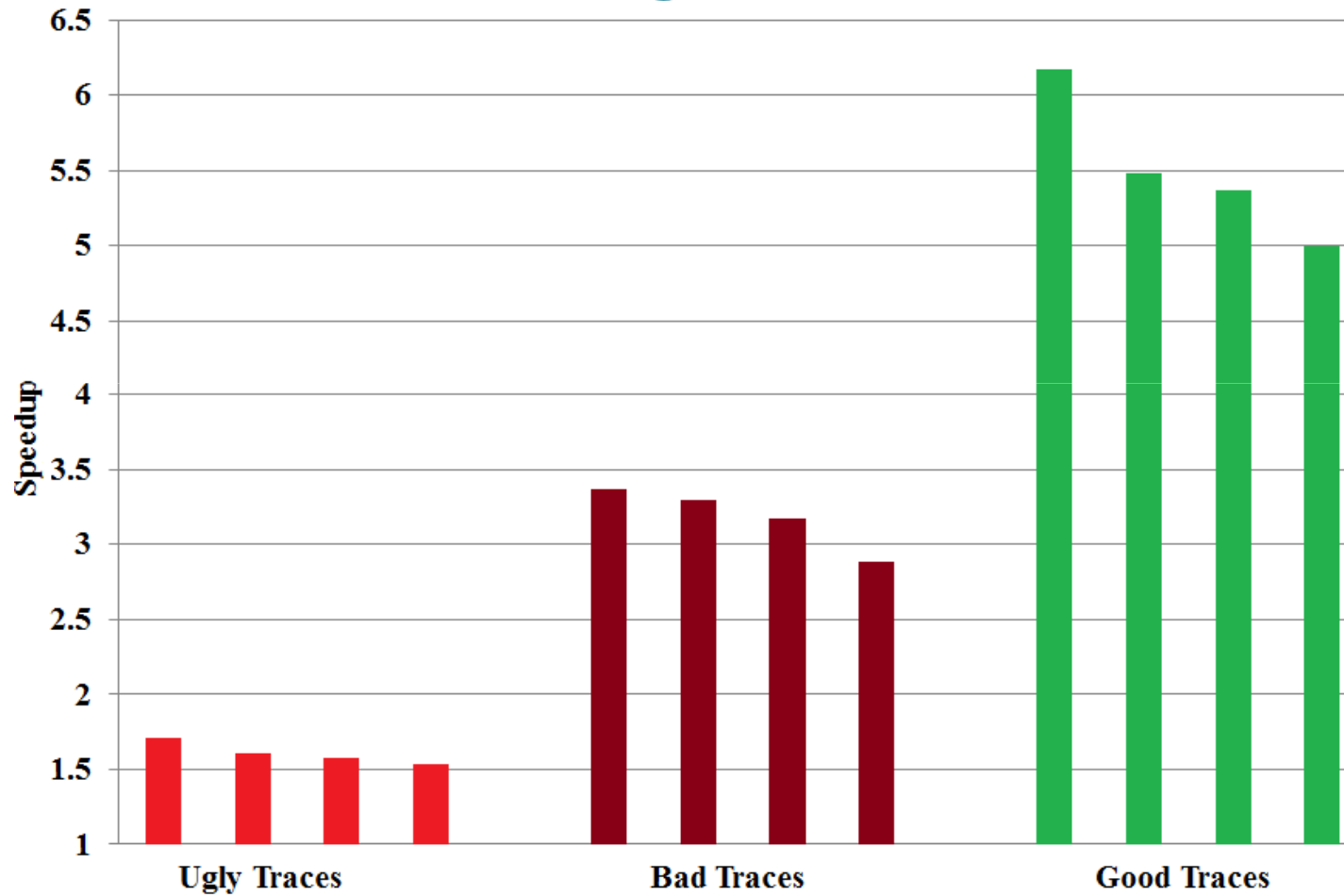
24

<b>Ugly Trace</b>	<b>Total Packets</b>	<b>Malicious Packets</b>	<b>Malicious Packets (%)</b>	<b>WM Time</b>	<b>Exscind Time</b>	<b>Speed up</b>
8	671116	258556	38.5	25.74	15.09	1.71
57	209188	86129	41.2	15.29	9.53	1.60
51	299713	129763	43.3	15.02	9.50	1.58
58	268000	119627	44.6	11.37	7.42	1.53
<b>Bad Trace</b>	<b>Total Packets</b>	<b>Malicious Packets</b>	<b>Malicious Packets (%)</b>	<b>WM Time</b>	<b>Exscind Time</b>	<b>Speed up</b>
1	688158	109311	15.9	75.44	23.00	3.28
22	659365	110363	16.7	76.22	22.90	3.33
0	771382	139866	18.1	66.18	20.87	3.17
2	642091	130901	20.4	55.04	19.11	2.88
<b>Good Trace</b>	<b>Total Packets</b>	<b>Malicious Packets</b>	<b>Malicious Packets (%)</b>	<b>WM Time</b>	<b>Exscind Time</b>	<b>Speed up</b>
VA	3606	54	1	105	17	6.18
GD	9354	175	2	345	63	5.48
LC	29474	1031	3	698	130	5.37
HM	2817	126	4	90	18	5.00



# Speedup

25



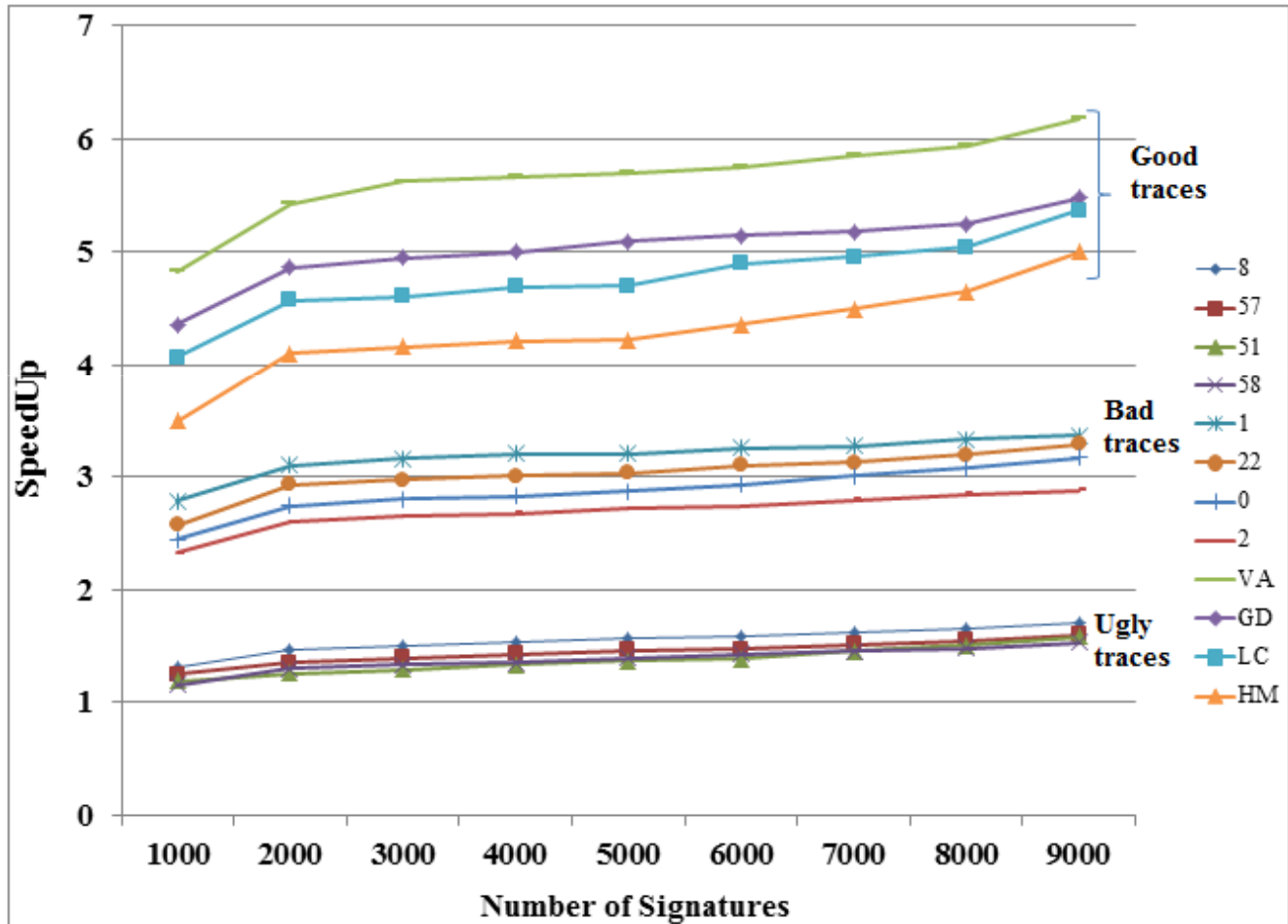
# Memory in KBs

26

<b>WM Memory</b>	<b>Our Memory</b>	<b>Added Overhead</b>	<b>Percentage overhead (%)</b>
537,178	537,766.1	588.1	1.095

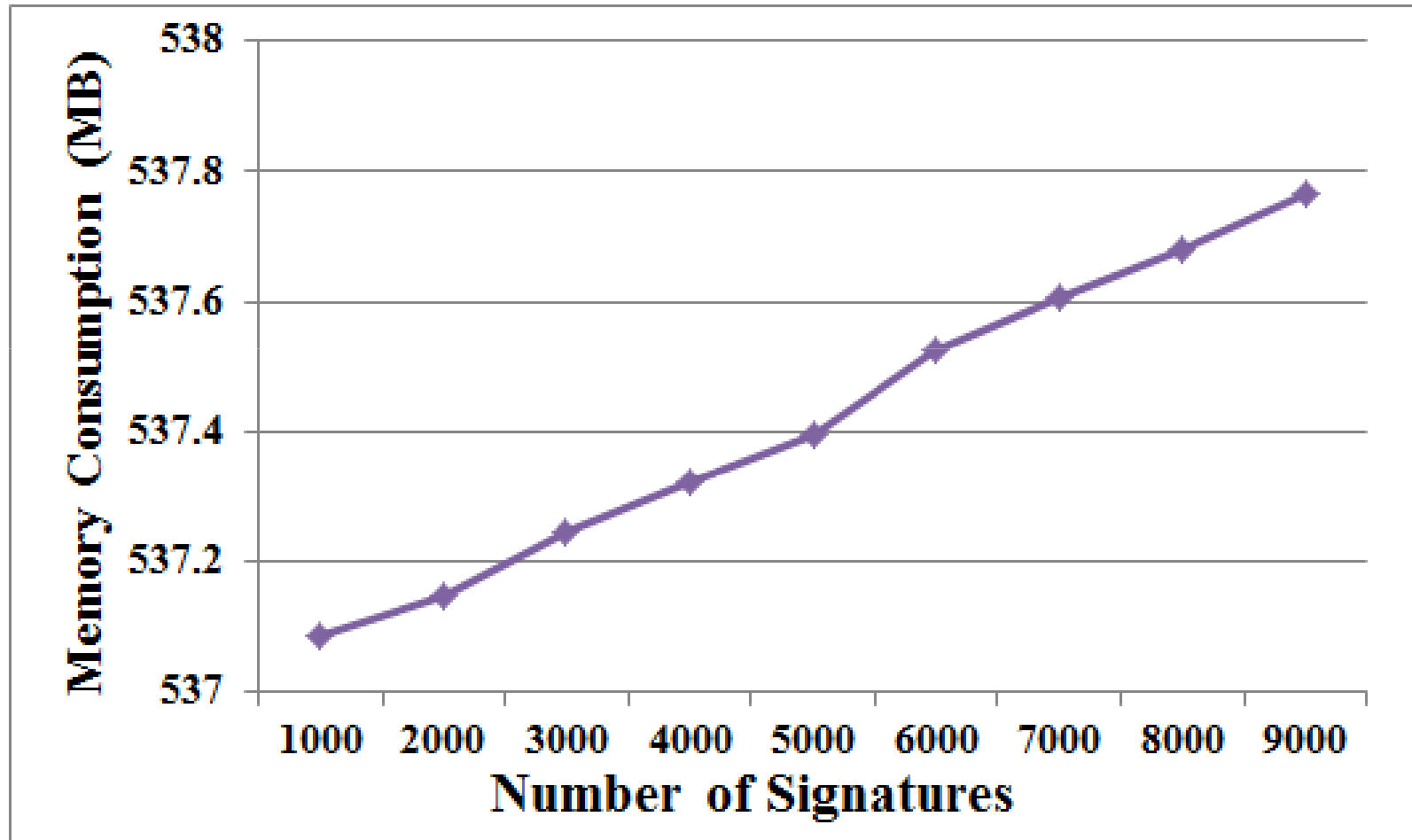
# Speedup Scaling

27



# Memory Scaling

28



# Conclusions

29

1. Average speedup of 3.4 times
2. Normal traffic speedup  $> 6$  times
3. Overhead: 0.11% increase in memory usage
4. Speedup is almost constant and undiminished with increasing number of signatures
5. Memory usage increases linearly

Qs

# Signatures Collisions

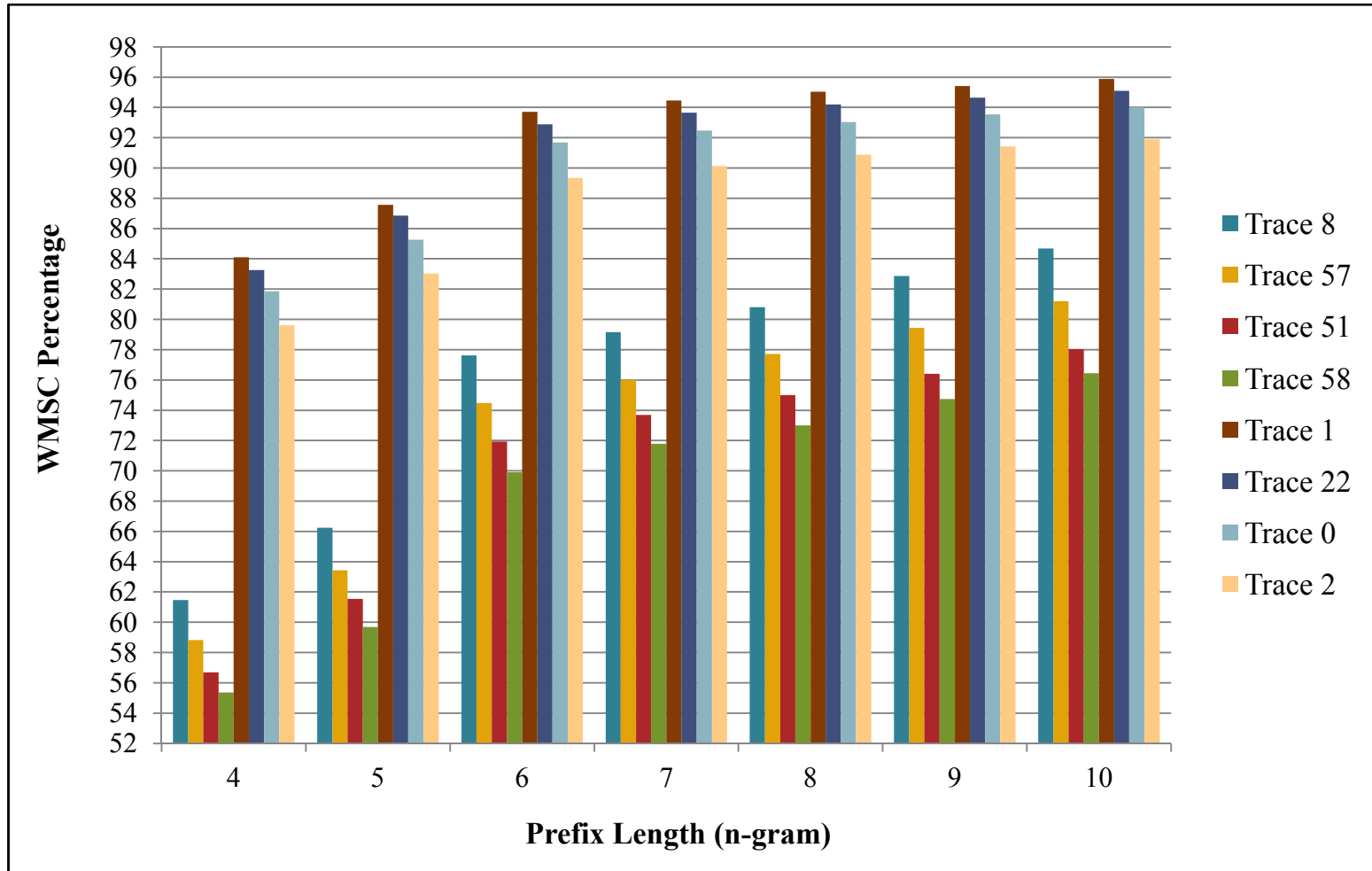
31

- Collision are the number of signatures with same hash bits!
- Bloom vector length 700,000 bit

Min Signature Length (n-gram)	Num of Signatures	n-grams of Signatures	Bloom Collisions	Bloom Collisions Percentage
4	9605	5017	1	0.000199322
5	9170	5662	1	0.000176616
6	9014	5934	1	0.00016852
7	8798	6080	0	0
8	8602	6193	0	0
9	8324	6208	0	0
10	8122	6277	0	0

# WMSC Worst Case Performance

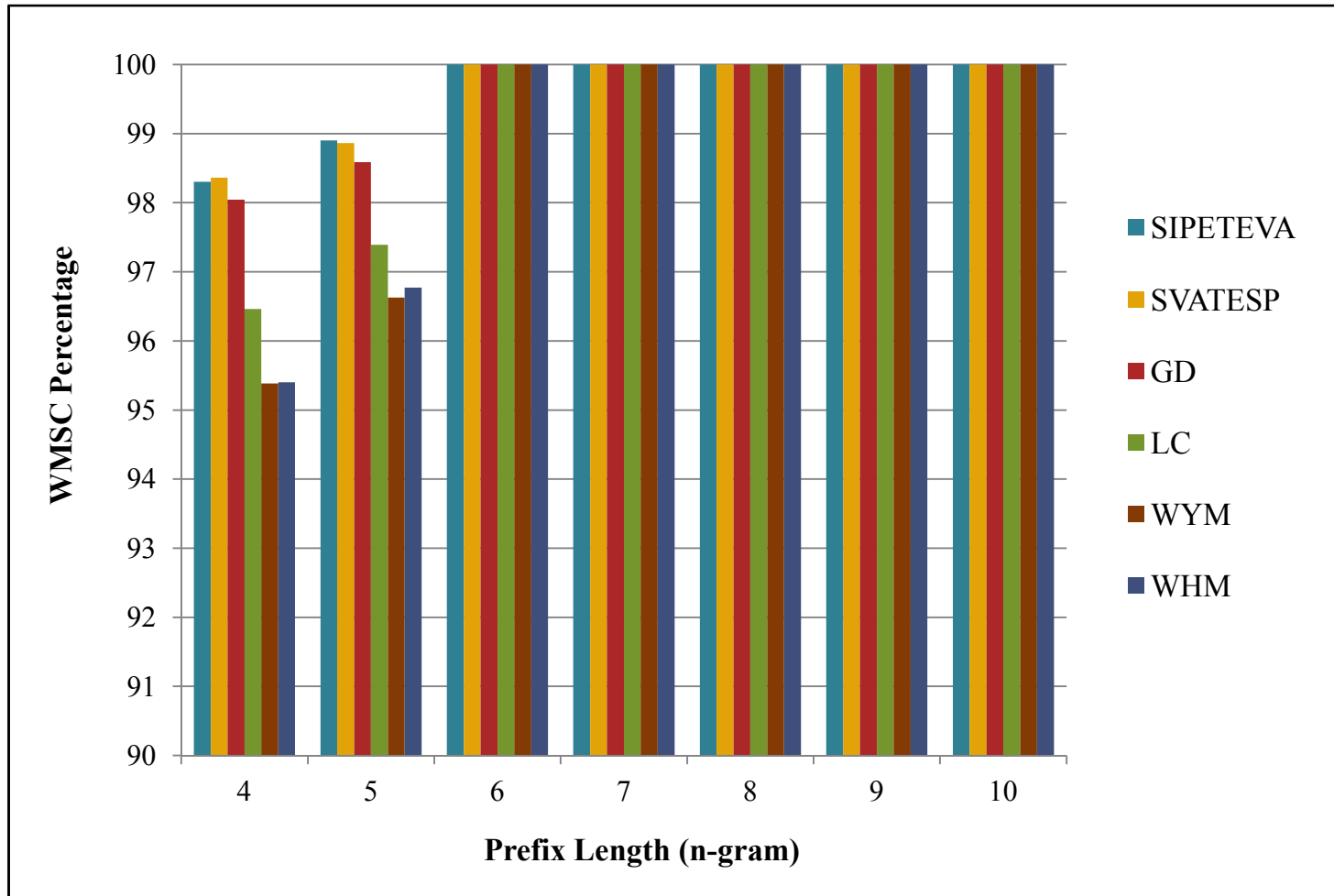
32





# WMSC Best Case Performance

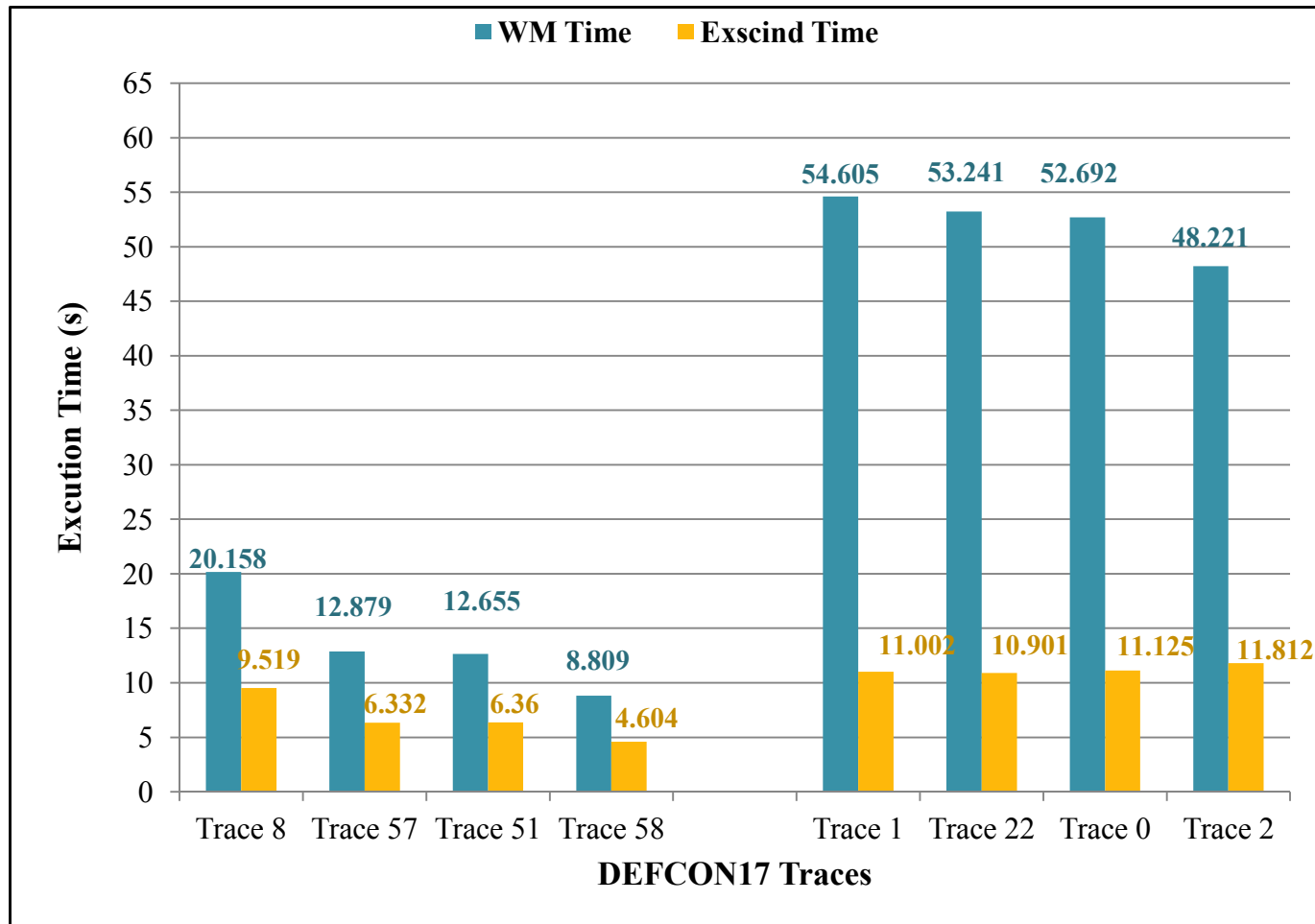
33



# Worst Case Execution Time Comparison

34

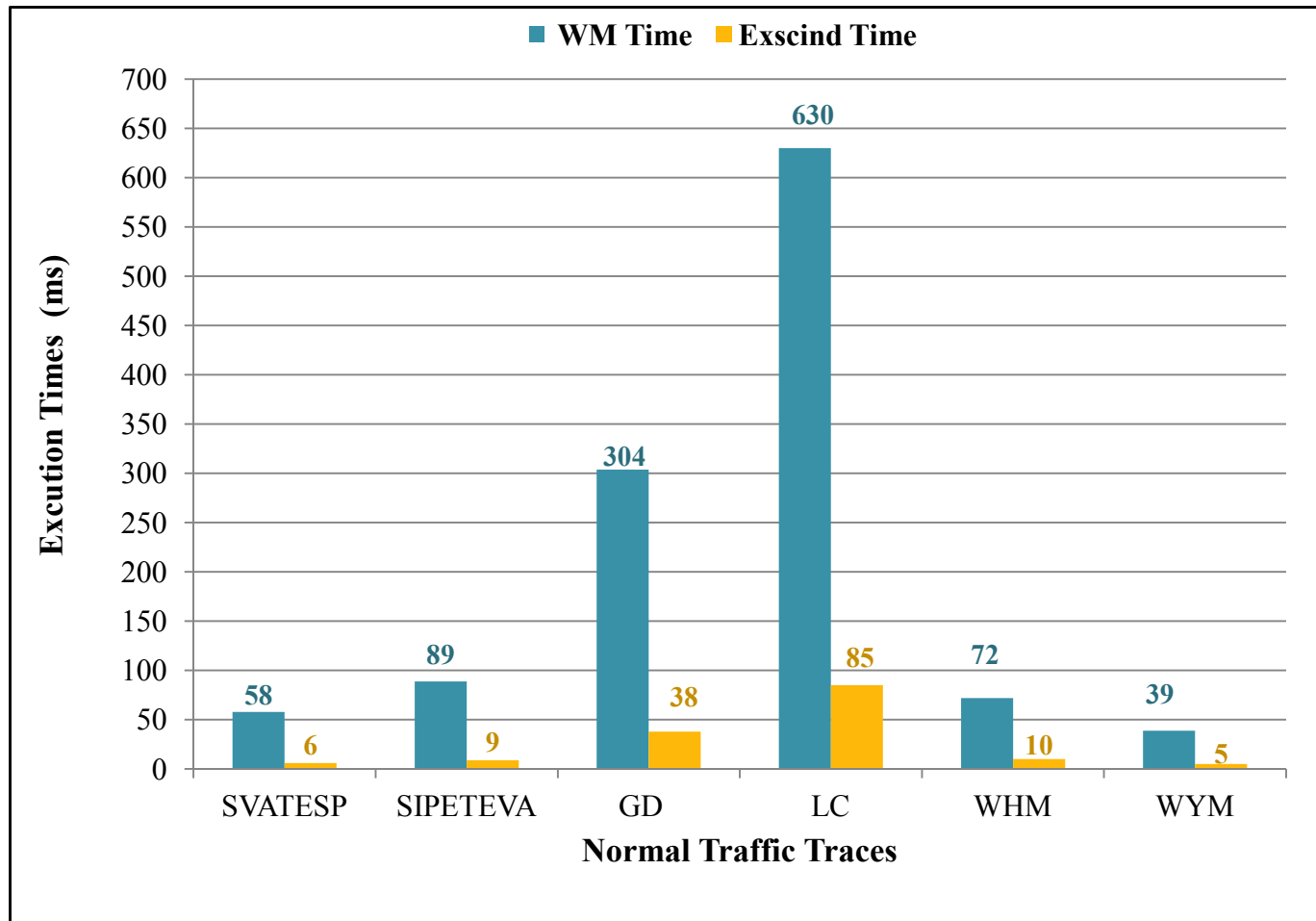
- For  $n=6$



# Best Case Execution Time Comparison

35

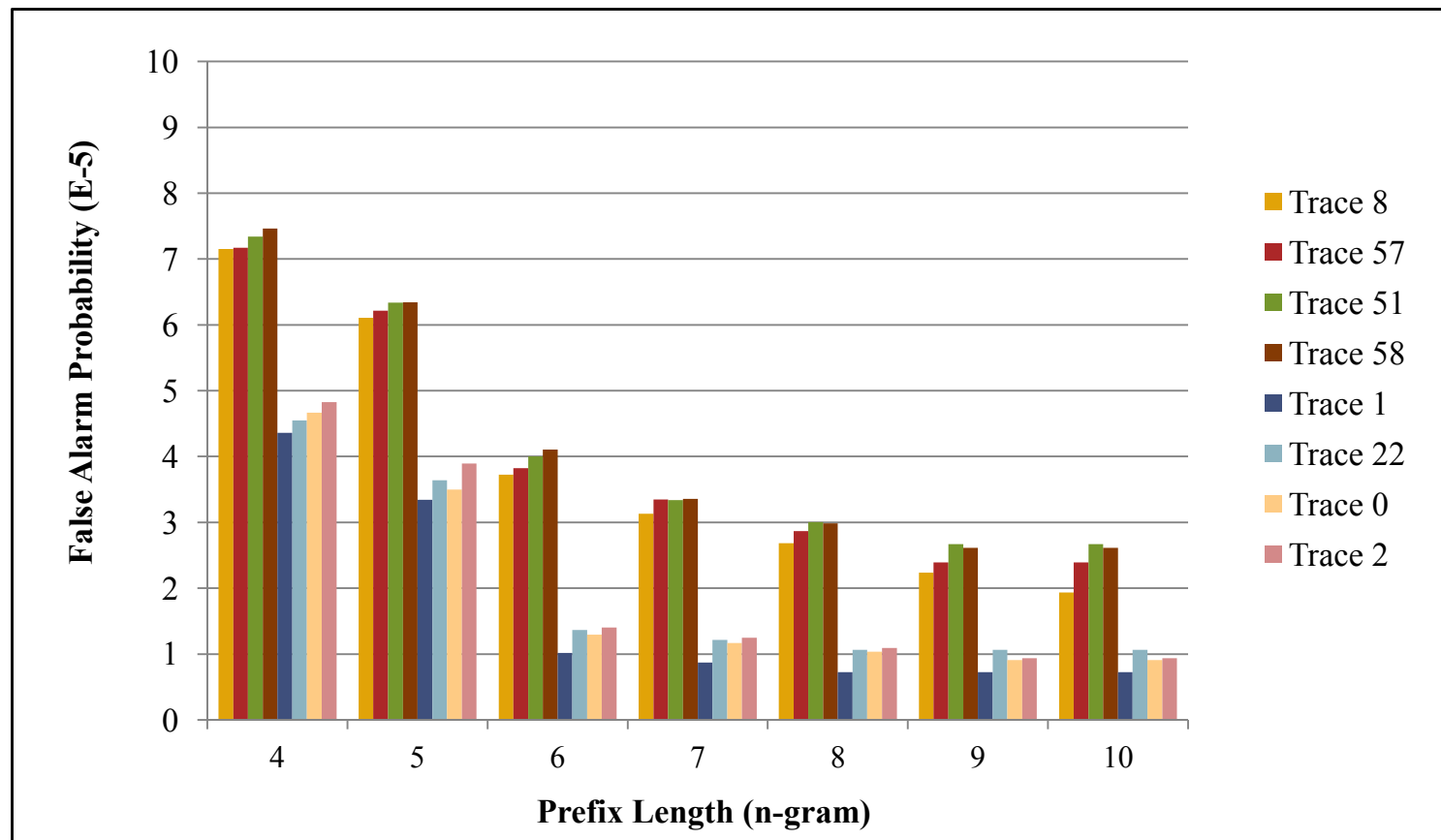
- For  $n=6$



# Worst Case False Alarms Probability

36

Additional pattern matching performed due to the fact the BF only inspects n-grams



# Hash Functions

37

1. **Rolling hash:** it is used for n-gram sliding windows move throughout a text. There are two popular rolling functions: Cyclic Polynomial and Rabin-Karp. We use Cyclic Polynomial because it is faster than Rabin-Karp
2. **SAX hash:** it is also called Shift-add-XOR, because it combines two functions; addition and XOR with a shift operation. It works as follows for each character in the string:

$$\text{SAX hash} = (\text{hash} \ll 5) + (\text{hash} \gg 2) + \text{character}$$