

# ZeroWall: Detecting Zero-Day Web Attacks through Encoder-Decoder Recurrent Neural Networks

Ruming Tang\*, Zheng Yang\*, Zeyan Li\*, Weibin Meng\*, Haixin Wang<sup>+</sup>,  
Qi Li\*, Yongqian Sun<sup>#</sup>, Dan Pei\*, Tao Wei<sup>^</sup>, Yanfei Xu<sup>^</sup> and Yan Liu<sup>^</sup>



\*Tsinghua University, China



<sup>+</sup>University of Science and Technology Beijing, China



<sup>#</sup>Nankai University, China



<sup>^</sup>Baidu

---

IEEE International Conference on Computer Communications, 27-30 April 2020 // Beijing, China

# WAFs Do Not Capture Zero-Days

- **WAFs (Web Application Firewalls)** are **wildly deployed** in industry, however, such **signature-based** methods are not suitable to detect zero-day attacks.
- Zero-day attacks in general are hard to detect and zero-day Web attacks are particularly challenging because:
  1. have **not been previously seen**
    - most **supervised** approaches are inappropriate
  2. can be carried out by a **single** malicious HTTP request
    - **contextual** information is not helpful
  3. very **rare** within a large number of Web requests
    - **collective** and **statistical** information are not effective

## ZeroWall

An **unsupervised** approach, which can **work with an existing WAF in pipeline**, to effectively detecting a zero-day Web attack hidden in **an individual Web request**.

# What We Want

- WAF detects those **known** attacks effectively.
  - filter out **known** attacks
- **ZeroWall** detects **unknown** attacks **ignored by WAF rules**.
  - report **new attack patterns** to operators and security engineers to **update WAF rules**.

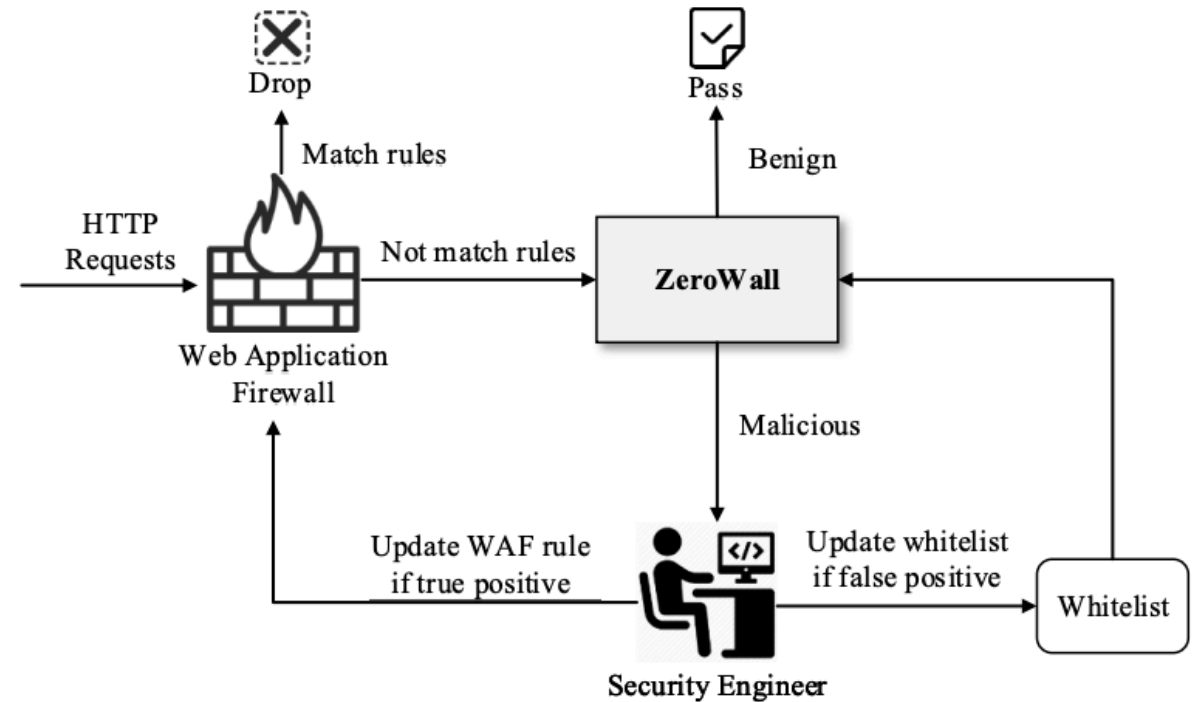
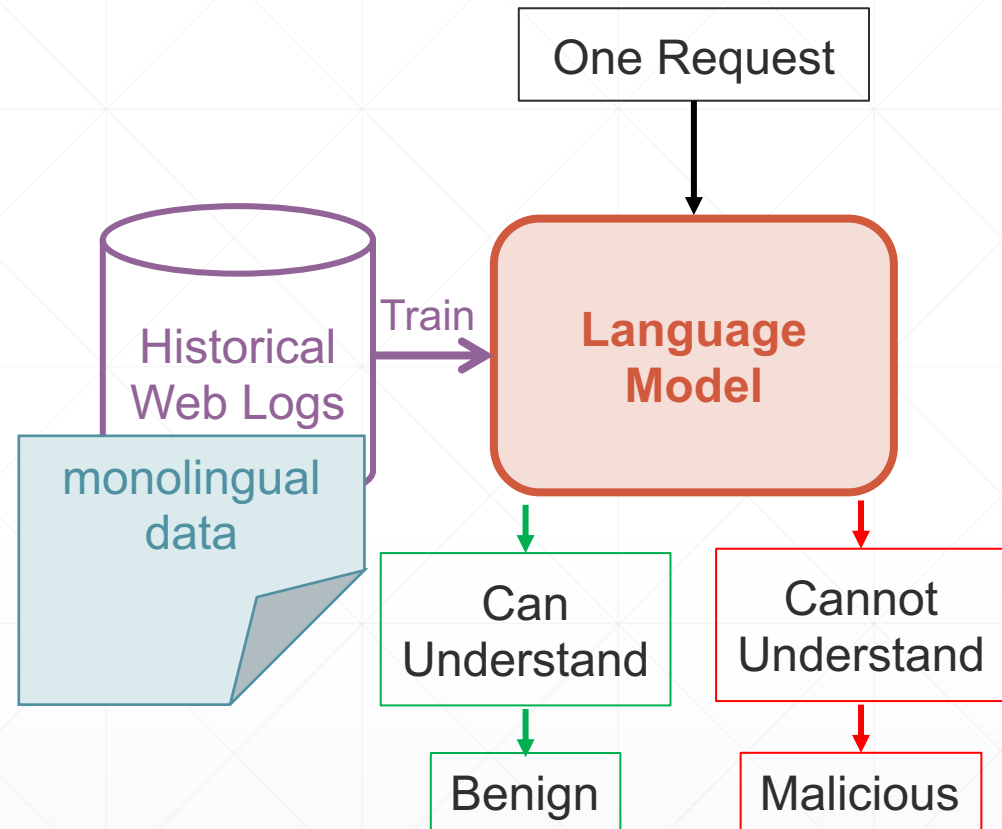


Figure 1: The workflow of *ZeroWall*.

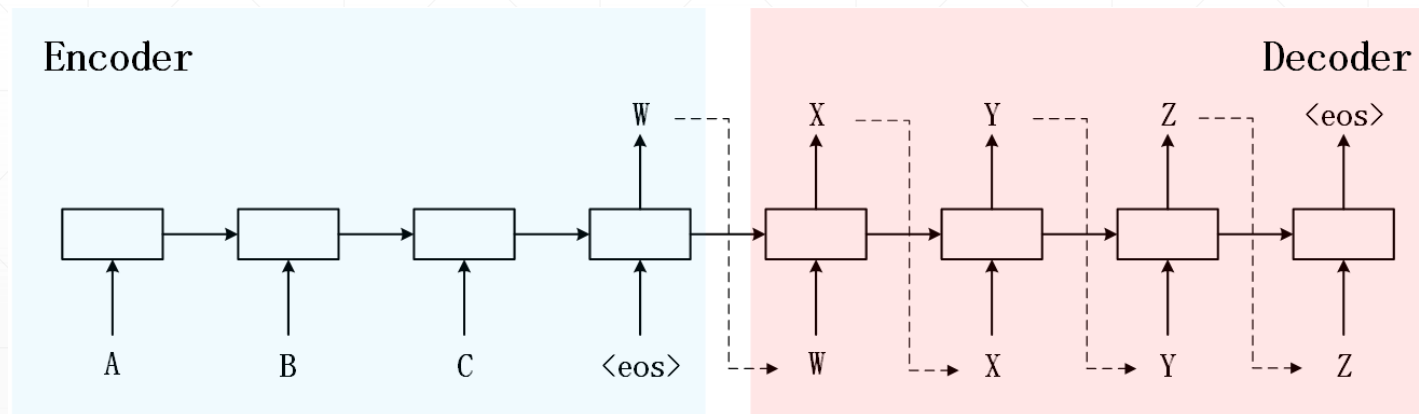
# Idea

- HTTP request is a **string following HTTP**, and we can consider an HTTP request as one **sentence** in the *HTTP request language*.
- **Most** requests are **benign**, and **malicious** requests are **rare**.
- Thus, we train a kind of **language model** based on historical logs, to **learn this language** from **benign requests**.

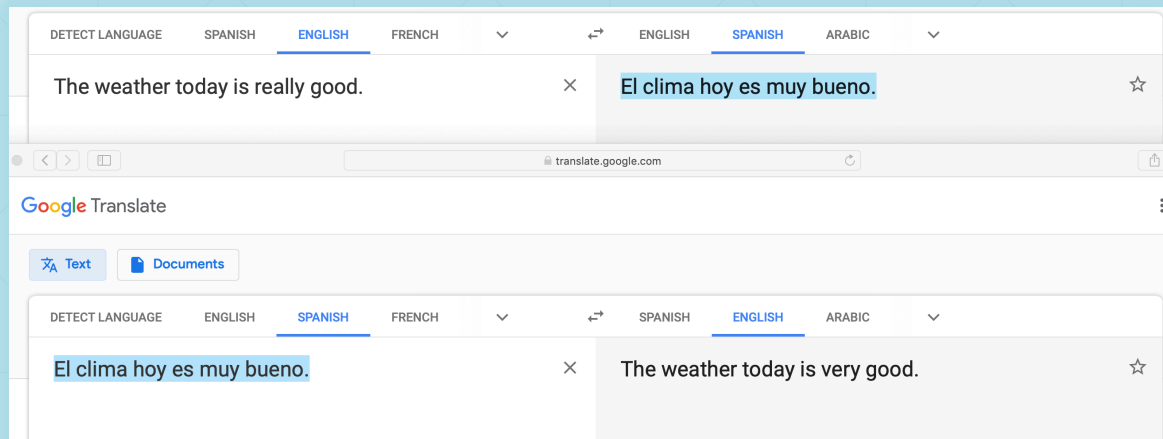


# Self-Translate Machine

- How to learn this “**Hyper-TEXT**” language?
- Use **Neural Machine Translation** model to train a **Self-Translate Machine**
  - **Encode** the original request into one *representation*
  - Then **Decode** it back

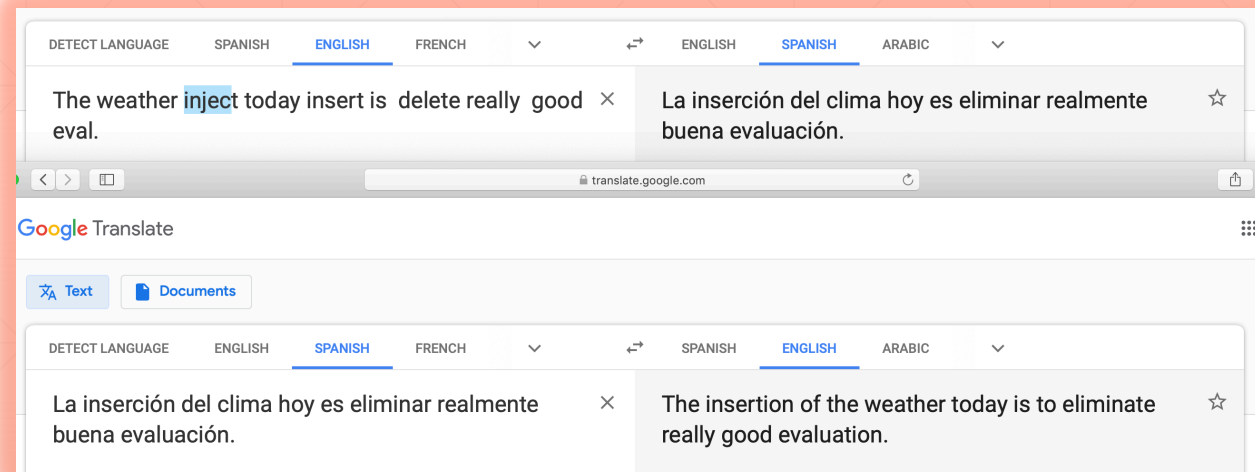


# Self-Translate Machine



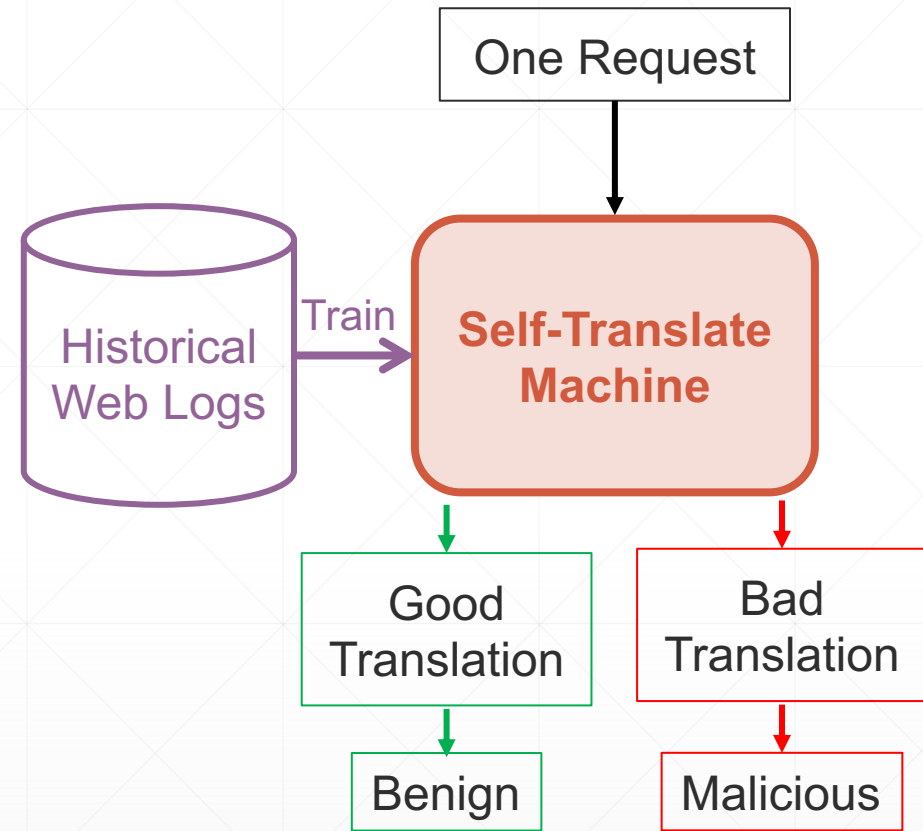
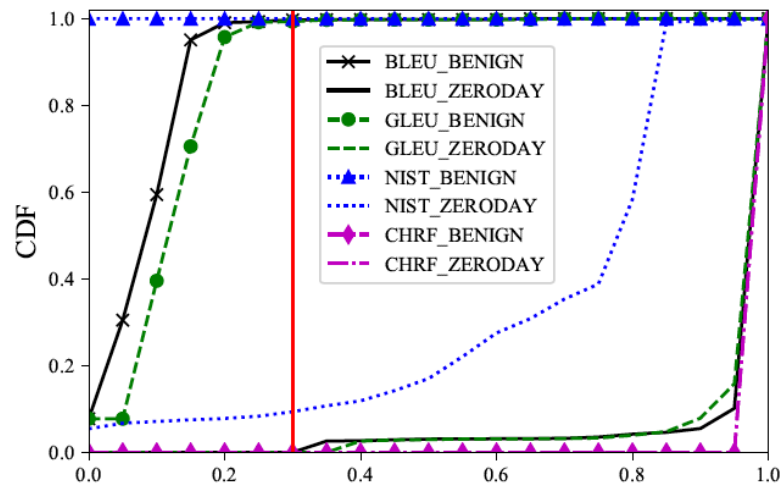
Self-translation works **well** for **normal** sentences

Output **deviates** significantly from the input, when the input is a sentence **not previously seen** in the training dataset of the self-translation models.



# Self-Translate Machine

- Translation Quality → Anomaly Score
- How to quantify the self-translation quality (anomaly score)?  
→ Use machine translation metrics



An attack detection problem → A machine translation quality assessment problem

# Self-Translated Sequence

- Translation Quality → Anomaly Score
  - Use BLEU as an example
  - Malicious Score =  $1 - BLEU\_Score$

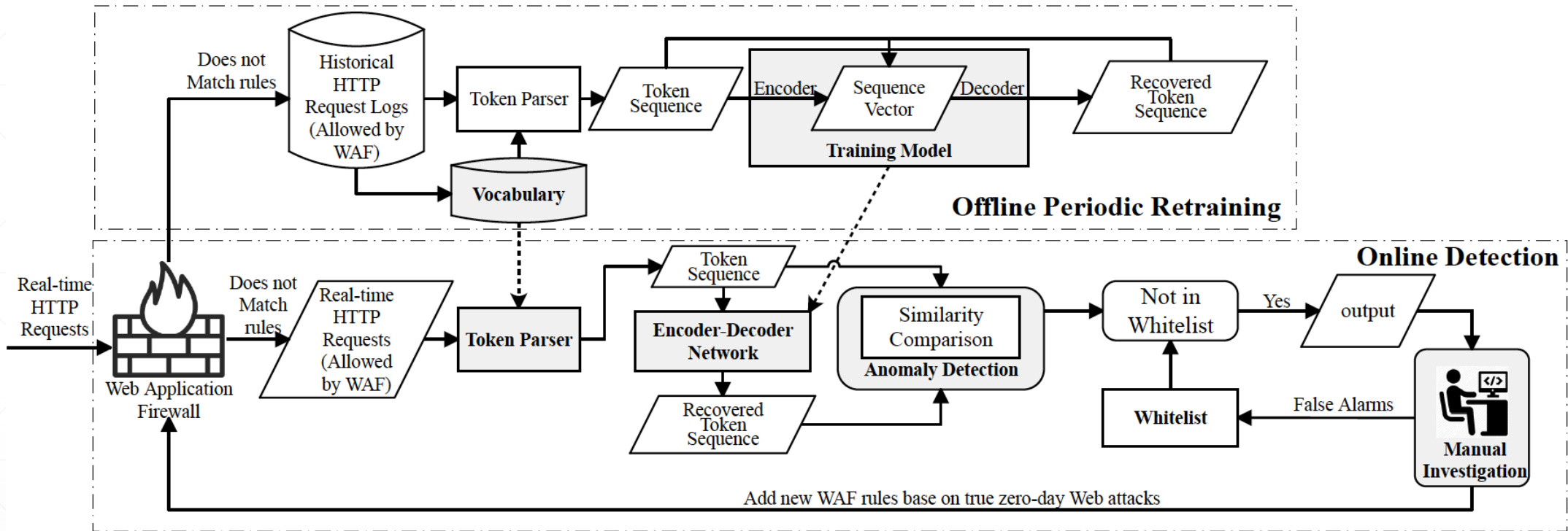
Original Request	POST http://localhost:8080/tienda1/publico/autenticar.jsp modo=entrar&login=caria&pwd=egipciaca&remember=off&B1=Entrar		
Tokenized	tienda1 publico autenticar jsp modo entrar login _OTHER_ pwd _OTHER_ remember off b1 entrar		
Translated	tienda1 publico autenticar jsp modo entrar login _OTHER_ pwd _OTHER_ remember on b1 entrar		
BLEU	0.8091	Malicious Score	0.1909

Original Request	<pre>POST http://m.thepaper.cn/admin_UploadDataHandler.ashx -----WebKitFormBoundaryRvkd1dbq3x1OJhUH\x0D\x0AContent-Disposition: form-data; name=\x22uploadify\x22; filename=\x2220170215180046.jpg\x22\x0D\x0AContent-Type: image/jpeg\x0D\x0A\x0D\x0A&lt;% eval request(\x22T\x22) %&gt;\x0D\x0A-----WebKitFormBoundaryRvkd1dbq3x1OJhUH\x0D\x0AContent-Disposition: form-data; name=\x22saveFile\x22\x0D\x0A\x0D\x0AAt.asp\x0D\x0A-----WebKitFormBoundaryRvkd1dbq3x1OJhUH\x0D\x0AContent-Disposition: form-data; name=\x22Upload\x22\x0D\x0A\x0D\x0ASubmit Query\x0D\x0A-----WebKitFormBoundaryRvkd1dbq3x1OJhUH--</pre>		
Tokenized	_OTHER_ ashx _OTHER_ content disposition form data name uploadify filename _pnum_0_ jpg content type image jpeg eval request onechr _OTHER_ content disposition form data name _OTHER_ onechr asp _OTHER_ content disposition form data name upload submit query _OTHER_		
Translated	_OTHER_ _OTHER_ do php _OTHER_ eval get_magic_quotes_gpc stripslashes _post chr _pnum_0_ chr _pnum_1_ _post chr _pnum_2_ chr _pnum_3+_ z0 _pnum_3+_ ini_set display_errors _pnum_3+_ set_time_limit _pnum_3+_ set_magic_quotes_runtime _pnum_3+_ echo onechr dirname _server script_filename if onechr onechr dirname _server path_translated		
BLEU	0	Malicious Score	1.0

An attack detection problem → A machine translation quality assessment problem

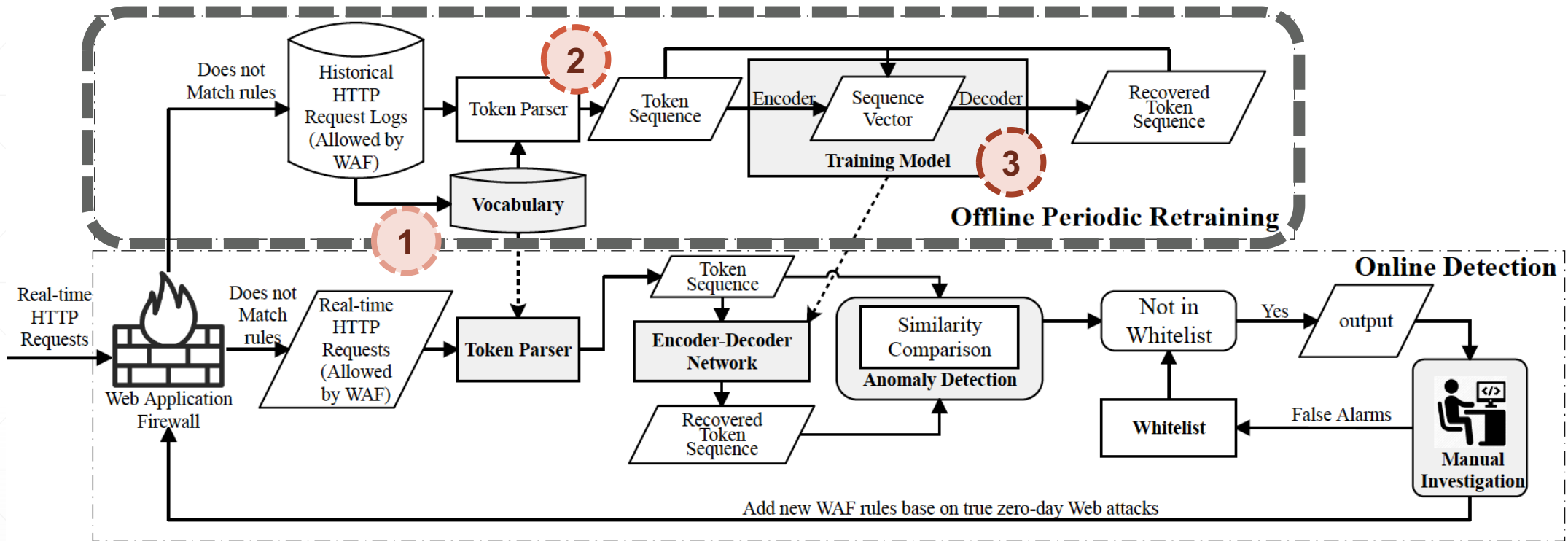


# ZeroWall Workflow



- Offline Periodic Retraining
  - Build and update **vocabulary** and re-train the **model**
- Online Detection
  - Detect **anomalies** in real-time requests for **manual investigation**

# Offline Training

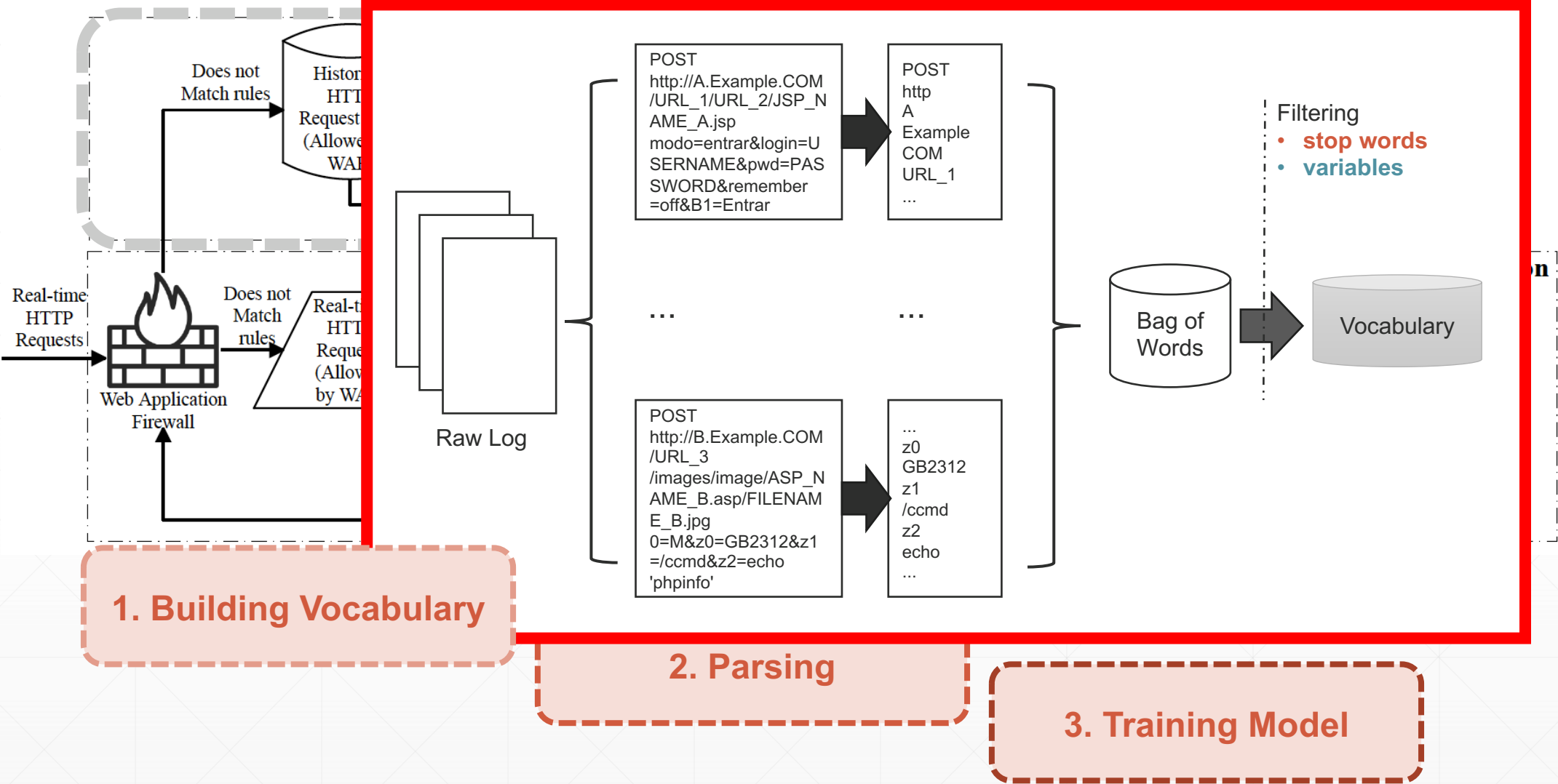


1. Building Vocabulary

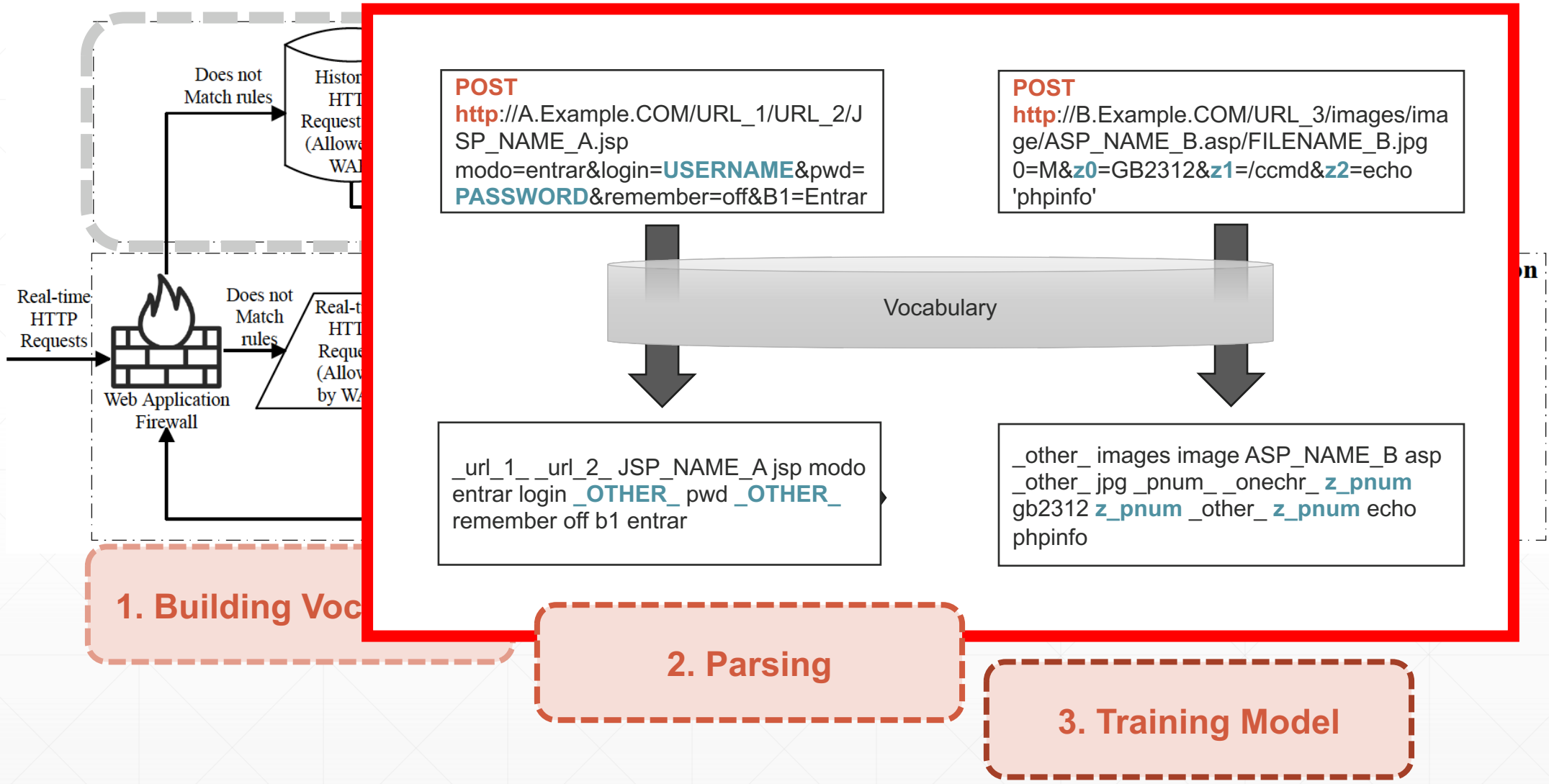
2. Parsing

3. Training Model

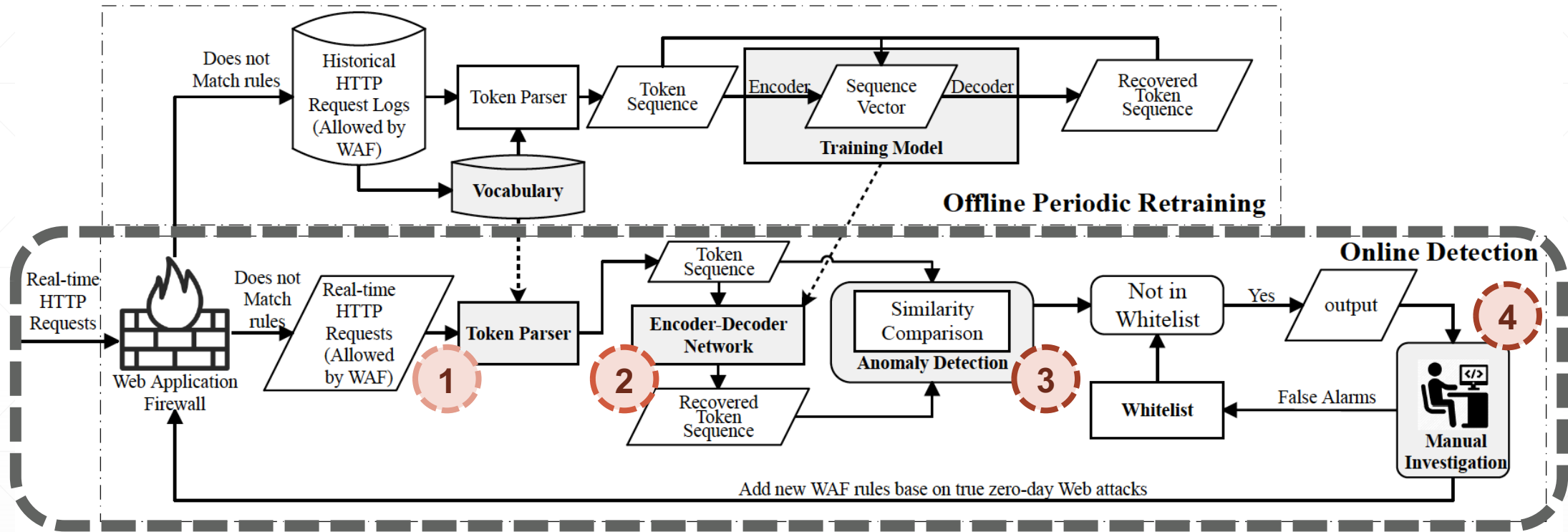
# Offline Training



# Offline Training



# Online Detection



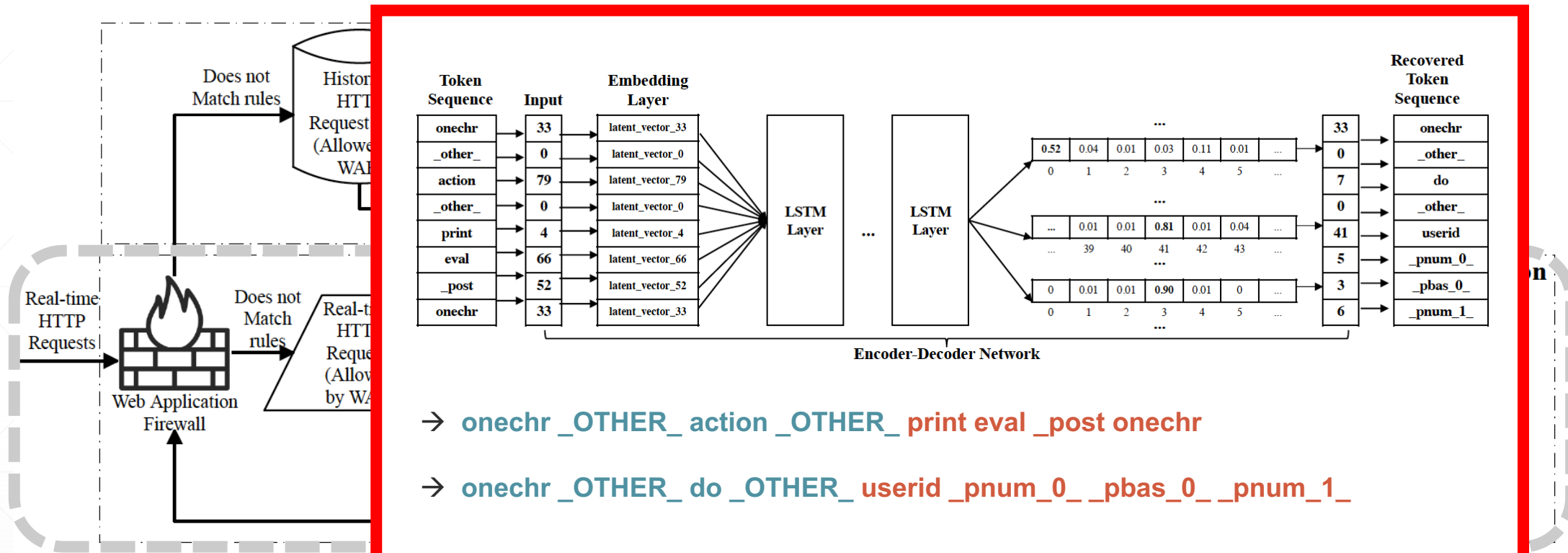
1. Parsing

2. Translation

3. Detection

4. Investigation

# Online Detection



→ onechr\_OTHER\_action\_OTHER\_print eval\_post onechr

→ onechr\_OTHER\_do\_OTHER\_userid\_pnum\_0\_\_pbas\_0\_\_pnum\_1\_

1. Token Parser

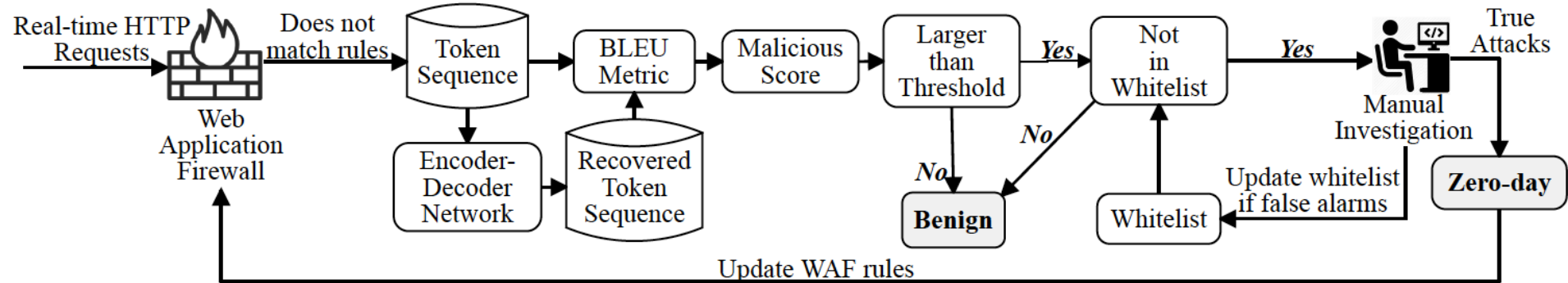
2. Translation

3. Detection

4. Investigation

# Online Detection

Compare the **original sequence** (token sequence) and the **translated sequence** (recovered token sequence).



1. BLEU Metric
2. Threshold
3. Check whitelist
4. Investigation

[ Larger? **Yes** → Go to step 3; **No** → **Benign** ]  
 [ Not in whitelist? **Yes** → Go to step 4; **No** → **Benign** ]  
 [ **True Attacks** → Update **WAF/IDS** ; **False Alarms** → Update **whitelist** rules ]

1. Token Parser

2. Translation

3. Detection

4. Investigation

# Real-World Deployment

- Data Trace:
  - 8 real world trace from an Internet company.
  - Over 1.4 billion requests in a week.
- Overview
  - Captured 28 different types of zero-day attacks, which contribute to 10K of zero-day attack requests in total.
  - False positives: 0~6 per day

#	D-1	D-2	D-3	D-4	D-5	D-6	D-7	D-8	Total
<b>Malicious*</b>	51839	186066	19515	53394	33724	2136811	42088623	90982519	135552491
<b>Zero-Day</b>	25	1118	283	4209	1188	2003	49011	83746	141583
<b>Benign</b>	1576235	3142793	13572827	15618518	31718124	177993528	528158912	534048878	1305829815
<b>Total</b>	1628099	3329977	13592625	15676121	31753036	180132342	570296546	625115143	1441523889
<b>B2M<sup>(1)</sup></b>	30.4	16.9	695.5	292.5	940.5	83.3	12.5	5.9	9.6
<b>B2Z<sup>(2)</sup></b>	63049.4	2811.1	47960.5	3710.7	26698.8	88863.5	10776.3	6377.0	9223.1

\* Known malicious filtered by WAF. (1) Ratio of Benign to Malicious (in WAF); (2) Ratio of Benign to Zero-Day



# Baselines & Labels

- Unsupervised Approaches
  - SAE (stacked auto-encoder), HMM and DFA (Deterministic Finite Automata)
  - Use data filtered by WAF as training set.
- Supervised Approaches
  - CNN, RNN and DT (decision tree)
  - Use all data (allowed/dropped) as training set and WAF results as labels.

# Evaluation Results

Trace	Approach	Precision	Recall	F1-Score
<b>D-1</b> #WAF-Malicious: 51,839 #Zero-Day Attacks: 25 #Benign: 1,576,235 #Total: 1,628,099 B2M: 30.4 B2Z: 63049.4	<i>ZeroWall</i>	<b>0.9855</b>	<b>1.0000</b>	<b>0.9889</b>
	<i>DT-Token</i>	0.0010	1.0000	0.0019
	<i>CNN-Token</i>	0.0010	1.0000	0.0019
	<i>RNN-Token</i>	0.0000	1.0000	0.0000
	<i>SAE</i>	0.0001	1.0000	0.0002
	<i>Hmmpayl</i>	0.0000	0.0000	0.0000
<b>D-2</b> #WAF-Malicious:186,066 #Zero-Day: 1,118 #Benign: 3,142,793 #Total: 3,329,977 B2M: 16.9 B2Z: 2811.1	<i>ZeroWall</i>	<b>1.0000</b>	<b>1.0000</b>	<b>1.0000</b>
	<i>DT-Token</i>	0.0547	0.3712	0.0931
	<i>CNN-Token</i>	0.3300	0.7784	0.4593
	<i>RNN-Token</i>	0.0005	0.9760	0.0010
	<i>SAE</i>	0.0005	0.9820	0.0010
	<i>Hmmpayl</i>	0.0000	0.0000	0.0000
<b>D-3</b> #WAF-Malicious: 19,515 #Zero-Day: 283 #Benign: 13,572,827 #Total: 13,592,625 B2M: 695.5 B2Z: 47960.5	<i>ZeroWall</i>	<b>0.9925</b>	<b>0.9687</b>	<b>0.9805</b>
	<i>DT-Token</i>	0.7388	0.2463	0.3658
	<i>CNN-Token</i>	0.4230	0.6376	0.5039
	<i>RNN-Token</i>	0.0000	0.9999	0.0001
	<i>SAE</i>	0.0015	0.9130	0.0030
	<i>Hmmpayl</i>	0.0000	0.0000	0.0000
<b>D-4</b> #WAF-Malicious: 53,394 #Zero-Day: 4,209 #Benign: 15,618,518 #Total: 15,676,121 B2M: 292.5 B2Z: 3710.7	<i>ZeroWall</i>	<b>0.9879</b>	<b>1.0000</b>	<b>0.9939</b>
	<i>DT-Token</i>	0.0001	1.0000	0.0002
	<i>CNN-Token</i>	0.0001	1.0000	0.0002
	<i>RNN-Token</i>	0.0008	1.0000	0.0016
	<i>SAE</i>	1.0000	0.0000	0.0000
	<i>Hmmpayl</i>	0.0000	0.0000	0.0000
<b>D-5</b> #WAF-Malicious: 33,724 #Zero-Day: 1,188 #Benign: 31,718,124 #Total: 31,753,036 B2M: 940.5 B2Z: 26698.8	<i>ZeroWall</i>	<b>0.9928</b>	<b>1.0000</b>	<b>0.9964</b>
	<i>DT-Token</i>	0.2497	0.0082	0.0153
	<i>CNN-Token</i>	0.6567	0.5410	0.5883
	<i>RNN-Token</i>	0.9988	0.0328	0.0629
	<i>SAE</i>	0.0000	0.0492	0.0000
	<i>Hmmpayl</i>	-	-	-
<b>D-6</b> #WAF-Malicious:2,136K #Zero-Day: 2,003 #Benign: 177,993,528 #Total: 180,132,342 B2M: 83.3 B2Z: 88863.5	<i>ZeroWall</i>	<b>1.0000</b>	<b>0.9897</b>	<b>0.9948</b>
	<i>DT-Token</i>	0.1733	0.0365	0.0576
	<i>CNN-Token</i>	0.0204	0.0590	0.0269
	<i>RNN-Token</i>	0.0000	1.0000	0.0000
	<i>SAE</i>	0.0001	0.1461	0.0001
	<i>Hmmpayl</i>	-	-	-
<b>D-7</b> #WAF-Malicious:42,088K #Zero-Day: 49,011 #Benign: 528,158,912 #Total: 570,296,546 B2M: 12.5 B2Z: 10776.3	<i>ZeroWall</i>	<b>0.9943</b>	<b>1.0000</b>	<b>0.9971</b>
	<i>DT-Token</i>	0.0874	0.0267	0.0377
	<i>CNN-Token</i>	0.8094	0.3027	0.4366
	<i>RNN-Token</i>	0.6857	0.5608	0.6120
	<i>SAE</i>	0.0001	0.5691	0.0002
	<i>Hmmpayl</i>	-	-	-
<b>D-8</b> #WAF-Malicious:90,982K #Zero-Day: 83,746 #Benign: 534,048,878 #Total: 625,115,143 B2M: 5.9 B2Z: 6377.0	<i>ZeroWall</i>	<b>0.9966</b>	<b>0.9983</b>	<b>0.9974</b>
	<i>DT-Token</i>	0.2036	0.3054	0.2396
	<i>CNN-Token</i>	0.2525	0.0275	0.0479
	<i>RNN-Token</i>	0.5237	0.0718	0.1242
	<i>SAE</i>	0.0008	0.3476	0.0017
	<i>Hmmpayl</i>	-	-	-
<b>D-9</b> #WAF-Malicious: 1,000 #Zero-Day: 1,000 #Benign: 1,000 #Total: 2,000 B2M: 0.5 B2Z: 1.0	<i>ZeroWall</i>	<b>1.0000</b>	<b>1.0000</b>	<b>1.0000</b>
	<i>DT-Token</i>	0.0000	0.0000	0.0000
	<i>CNN-Token</i>	0.0000	0.0000	0.0000
	<i>RNN-Token</i>	0.0000	0.0000	0.0000
	<i>SAE</i>	0.0000	0.0000	0.0000
	<i>Hmmpayl</i>	0.0000	0.0000	0.0000

Trace	Approach	Precision	Recall	F1-Score
<b>D-9</b> #WAF-Malicious: 1,000 #Zero-Day: 1,000 #Benign: 1,000 #Total: 2,000 B2M: 0.5 B2Z: 1.0	<i>ZeroWall</i>	<b>1.0000</b>	<b>1.0000</b>	<b>1.0000</b>
	<i>DT-Token</i>	0.0000	0.0000	0.0000
	<i>CNN-Token</i>	0.0000	0.0000	0.0000
	<i>RNN-Token</i>	0.0000	0.0000	0.0000
	<i>SAE</i>	0.0000	0.0000	0.0000
	<i>Hmmpayl</i>	0.0000	0.0000	0.0000

# A Zero-Day Case

Captured 28 different types of zero-day attacks, including webshell, SQL injection, probing, trojan and other exploiting against particular applications. For each category, the security engineers have already composed a new WAF rule to detect these attacks in the future.

- These attack is detected by **ZeroWall, CNN and RNN**.
- WAF** are usually based on **keywords**, e.g., **eval**, **request**, **select** and **execute**.
- ZeroWall** is based on the “**understanding**” of benign requests. The structure of this zero-day attack request is more like a programming language.

```

...
searchword=d&order=} {end if} {if:1)print_r(
$_POST[func]($_POST[cmd]));//}
{end if}&func=assert&cmd=phpinfo();

```

---

Token Sequence: search php searchtype \_pnum\_0\_ \_OTHER\_ onechr order end if if \_pnum\_1\_ \_OTHER\_ \_post \_OTHER\_ \_post cmd end if \_OTHER\_ assert cmd phpinfo

contains **none** of **WAF keywords**

overlap with tokens in training set for **CNN and RNN**

1	plus ad_js php aid _pnum_0_ onechr assert _pnum_1_ execute execute function bd byval onechr for onechr _pnum_2_ to len onechr step _pnum_3+_ onechr mid onechr _pnum_3+_ if isnumeric mid onechr _pnum_3+_ then execute bd bd chr onechr else execute bd bd chr onechr mid onechr _pnum_3+_ onechr _pnum_3+_ end if chr _pnum_3+_ next end function response write execute on error resume next bd _phex_0_ response write response end
2	preview php _OTHER_ php assert _OTHER_ onechr
3	lib _OTHER_ module inc php _OTHER_ eval _OTHER_ onechr class _OTHER_ onechr <b>phpinfo</b>
4	cms _OTHER_ uploads _OTHER_ php id assert _OTHER_ eval base64_decode _post z0 z0 _pbas_0_
5	myship php cmd eval base64_decode _post z0 z0 _pbas_0_

# Whitelist

- To mitigate **False Alarms**, we add **whitelist** to our approach.
- The **numbers of whitelist rules** refer to how many whitelist rules are added each day, based on the FPs labeled on that day. (No rules applied on 0602 since it is the first day of testing set.)
- The results shows that the whitelist **reduces the number of FPs with low overhead** (numbers of rules are very small).
- Based on these results, we believe ZeroWall is practical in real-world deployment.

Date	Precision		# of FP		# of white-list rules
	No WL	WL	No WL	WL	
0602	0.9972	-	16	-	5
0603	0.9643	0.9753	222	152	3
0604	0.9580	0.9999	310	1	1
0605	0.9731	0.9944	320	65	6
0606	0.9845	0.9993	121	5	1
0607	0.9672	1.0000	194	0	0

# Overhead

- Training and testing speed with and without hash table (requests/s)

Trace	Incoming Requests	Training		Testing	
		No Hash	Hash	No Hash	Hash
<i>D-1</i>	2.60	1.09	256.89	229.24	39634.40
<i>D-2</i>	5.19	3.72	202.13	785.75	65556.80
<i>D-3</i>	22.44	7.09	835.43	514.33	50420.17
<i>D-4</i>	25.83	5.42	1014.67	305.42	50913.24
<i>D-5</i>	52.45	12.48	1046.55	414.86	38132.88
<i>D-6</i>	294.30	1.47	4001.95	70.04	176255.90
<i>D-7</i>	873.36	3.23	4262.48	53.77	88989.06
<i>D-8</i>	883.16	6.67	6389.23	142.29	106692.90

\*The incoming requests refer to the average number of requests received by the customer per second.

Intel(R) Xeon(R) Gold 6148 CPU 2.40GHz \* 2  
512GB RAM

# Summary

- Present a zero-day web attack detection system **ZeroWall**
  - **Augmenting** existing **signature-based WAFs**
  - Use **Encoder-Decoder Network** to learn patterns from normal requests
  - Use **Self-Translate Machine & BLEU Metric**
- **Deployed** in the wild
  - Over **1.4** billion requests
  - Captured **28** different types of zero-day attacks (**10K** of zero-day attack requests)
  - Low overhead

**An attack detection problem →  
A machine translation quality  
assessment problem**

**Thanks!  
And Questions**

Ruming Tang: trm14@mails.tsinghua.edu.cn