

Generative Intrusion Detection and Prevention on Data Stream

HyungBin Seo and MyungKeun Yoon

Kookmin University

USENIX Security 2023

Motivation

- Different Types of Intrusion Detection

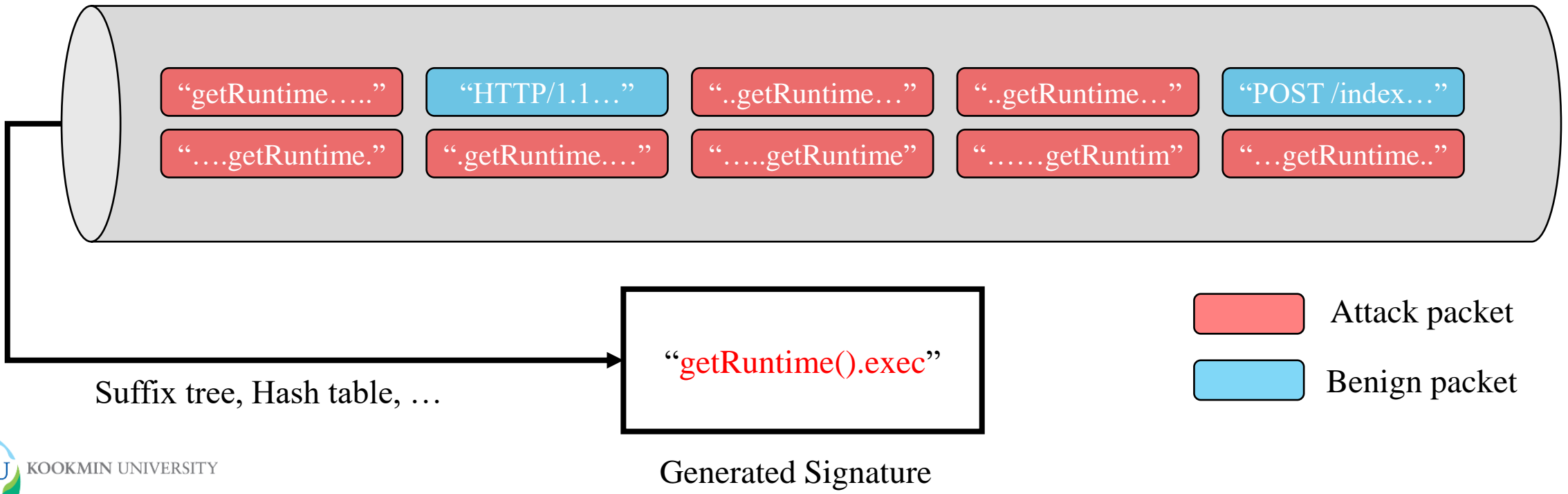
- Persistent zero-day attack, signature-group generation, evidence support

Type		Accuracy	Automation	Zero-day
Misuse-detection (signature)	Manual rule composition [49, 53]	High	Low	Low
	Signature generation [9, 48]	Medium	Medium	High
	Signature-group generation (GIPS)	High	High	High
Anomaly-detection (ML)	Supervised learning [19, 52]	Medium	Medium	Medium
	Unsupervised learning [54, 62]	Low	Medium	High

Motivation

- Previous signature-generation schemes [9], [24], [35], [39], [48]
 - high-volume attacks such as DDoS, worms

Attack data ratio = 8/10 → Big Group



Motivation

- Previous signature-generation schemes [9], [24], [35], [39], [48]
 - high-volume attacks such as DDoS, worms



Dataset with high attack ratio



signature	frequency
"getRunti..."	8
"exec()."	6
"POST"	1
"HTTP/1.1"	1
...	...

Hash table

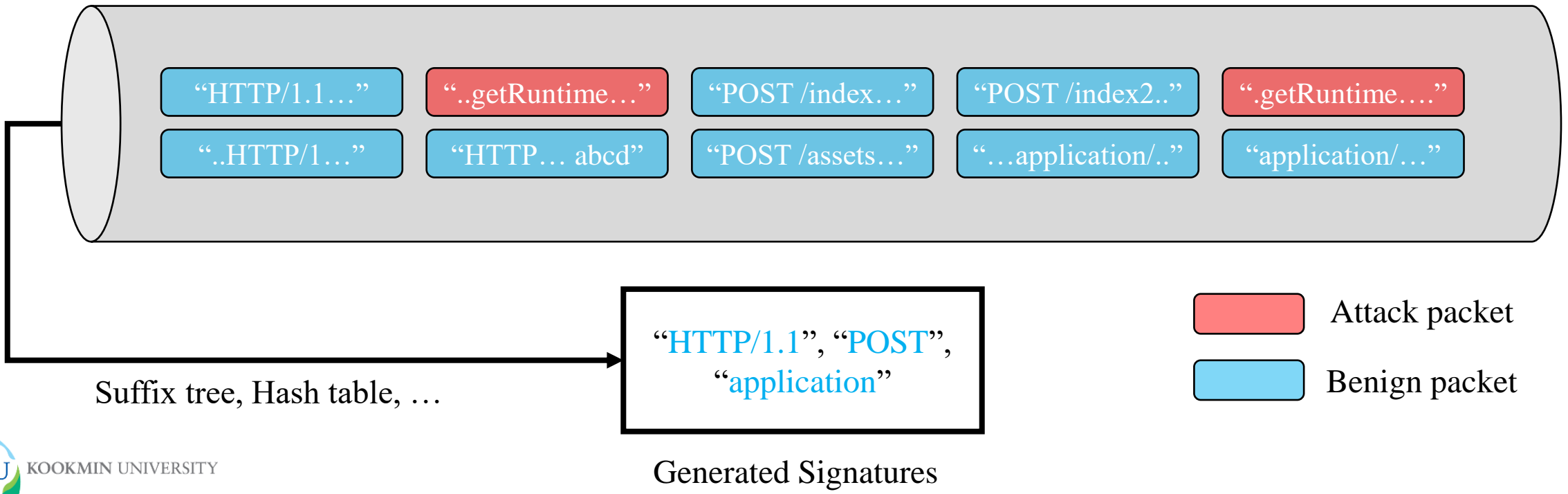


Generated Signatures

Motivation

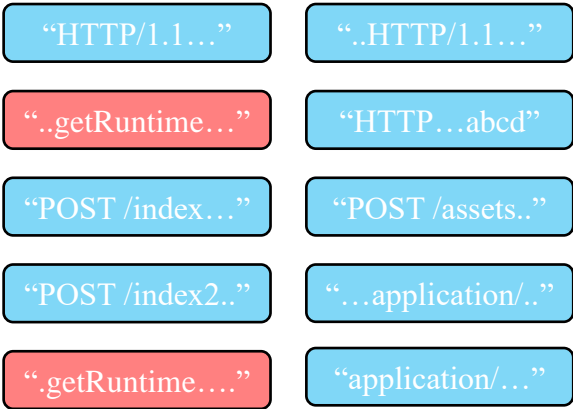
- Previous schemes would cause false-positives for **persistent zero-day attacks**.

Attack data ratio = 2/10 → **Still Big Group**



Motivation

- Previous schemes would cause false-positives for **persistent zero-day attacks**.



Dataset with low attack ratio



signature	frequency
"HTTP/1.1"	8
"POST"	5
"/index.."	3
"application"	3
...	...

Hash table

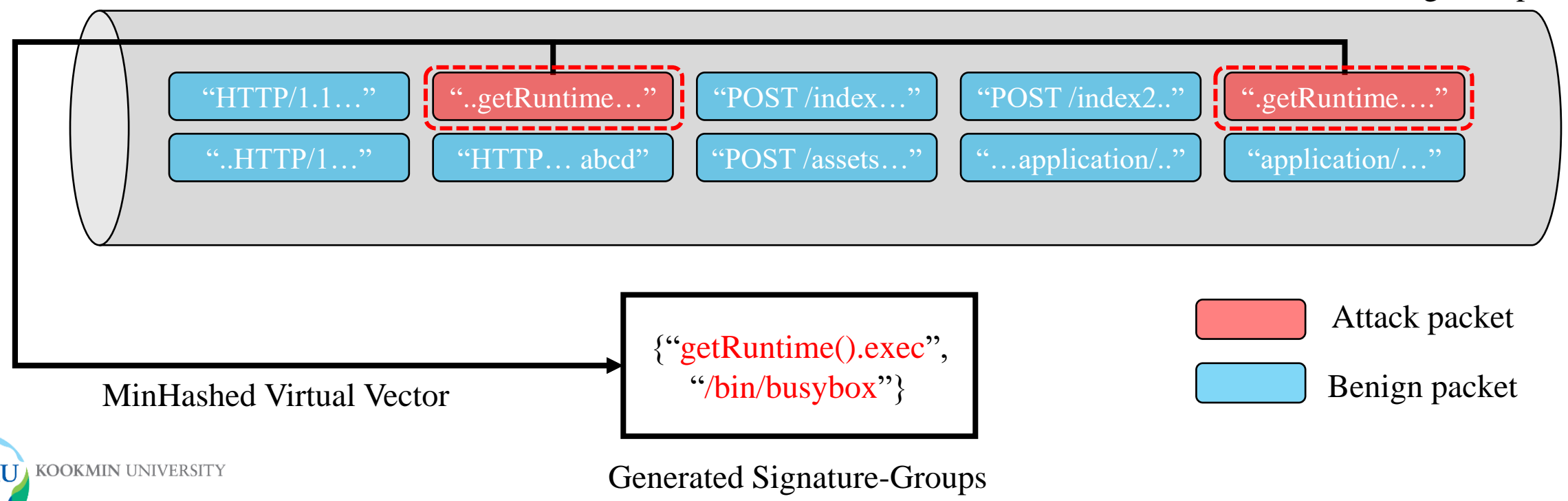


Generated Signatures

Motivation

- GIPS
 - Generative Intrusion detection and Prevention on data Stream
 - Data analysis tool

Attack data ratio = 2/10 → Still Big Group

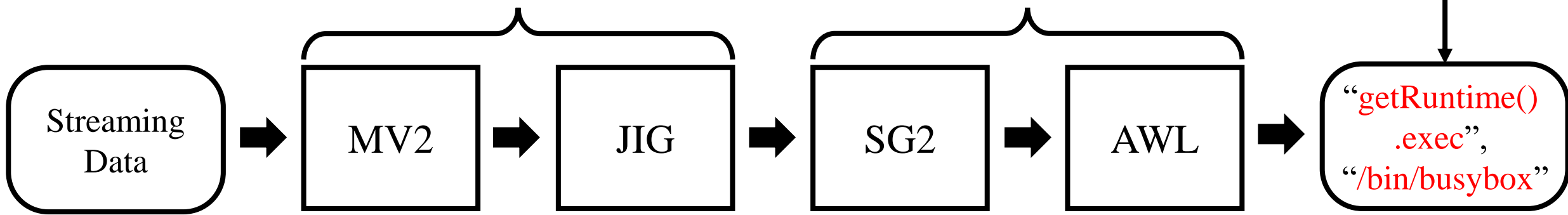


Overall Process of GIPS



Big-group identification

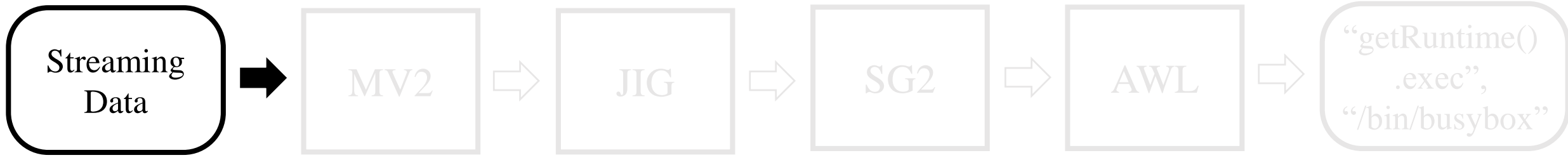
Signature-group generation



- MV2: Minhashed Virtual-Vector
- JIG: Jaccard-Index Grouping

- SG2: Signature-Group Generation
- AWL: Automatic WhiteListing

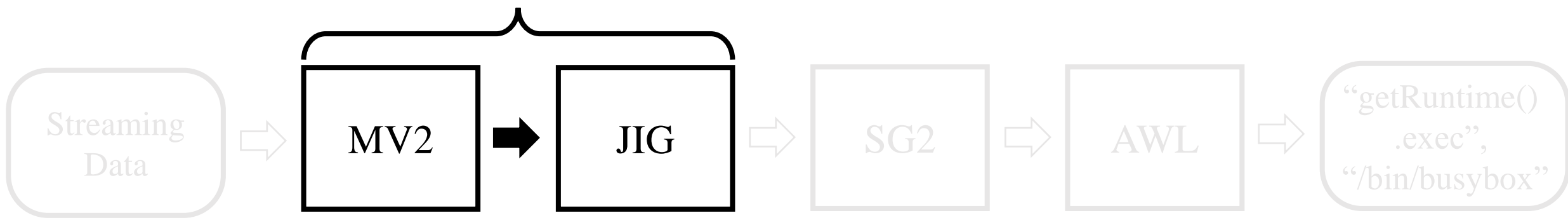
Overall Process of GIPS



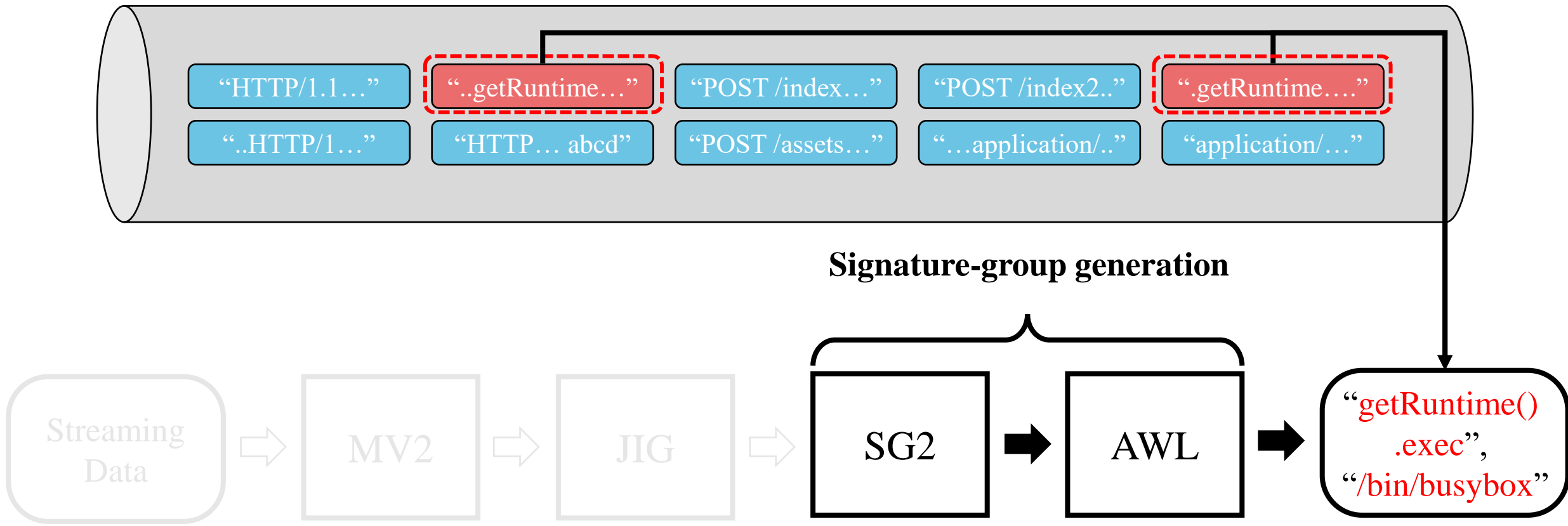
Overall Process of GIPS



Big-group identification



Overall Process of GIPS



MV2: Minhashed Virtual-Vector

“HTTProotadmin1234@”



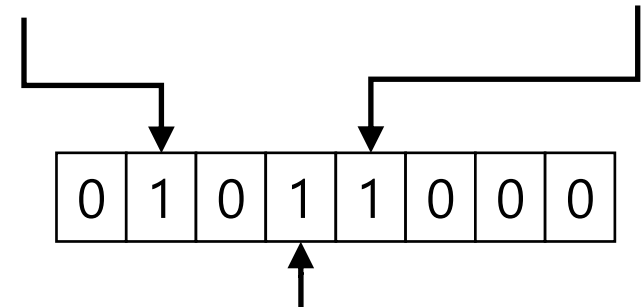
{“HTTP”, “root”, “admin”, “1234@”}

(a) Content-Defined Chunking*

$k=3$
 $h_0(\cdot), h_1(\cdot), h_2(\cdot),$



$$\left(\min_{x \in \text{chunks}} \{h_0(x)\} \right) \text{mod } m \quad \left(\min_{x \in \text{chunks}} \{h_2(x)\} \right) \text{mod } m$$



$$\left(\min_{x \in \text{chunks}} \{h_1(x)\} \right) \text{mod } m$$

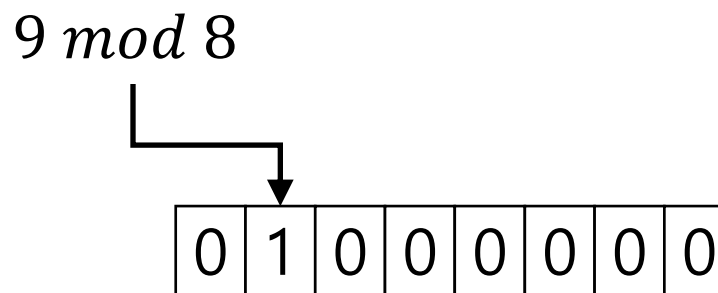
(b) Minhashed Virtual-Vector

MV2: Minhashed Virtual-Vector

{“HTTP”, “root”, “admin”, “1234@”}

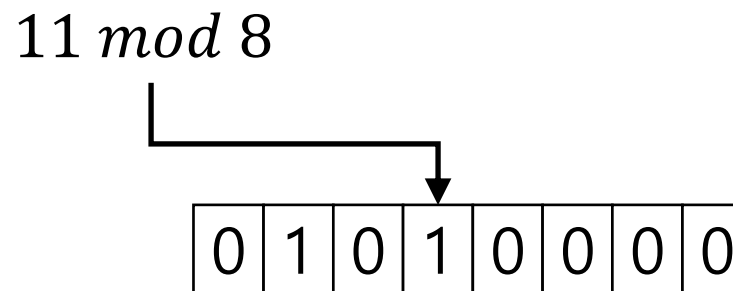
$h_0("HTTP") = 71$
 $h_0("root") = 48$
 $h_0("admin") = 9$
 $h_0("1234@") = 57$

→ $\min_{x \in \text{chunks}} h_0(x) = 9$



$h_1("HTTP") = 87$
 $h_1("root") = 11$
 $h_1("admin") = 95$
 $h_1("1234@") = 24$

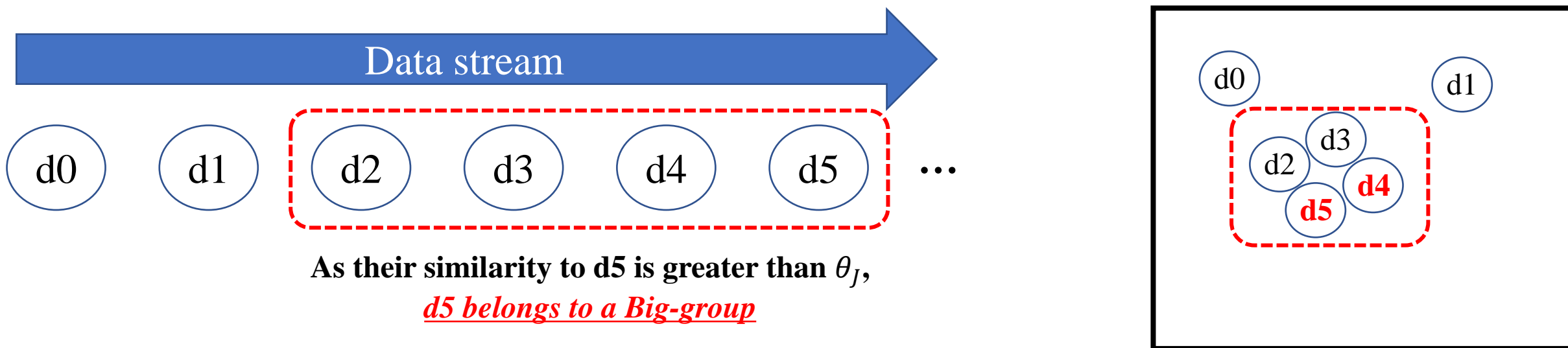
→ $\min_{x \in \text{chunks}} h_1(x) = 11$



JIG: Jaccard-Index Grouping

- Big-group of d_i from data stream

$$sg(d_i) = \{d_j | j(d_i, d_j) > \theta_J\} \cup \{d_i\} \quad b_r(d_i) = \frac{|sg(d_i)|}{i} > \theta_B$$

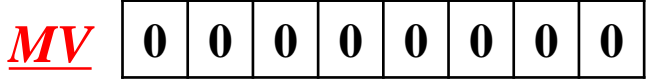


Naïve clustering on data stream requires a large time & space-complexity!

JIG: Jaccard-Index Grouping

d0 = “GET /index.html HTTP/1.1”
 d1 = “HTTProotadmin1234@”
 d2 = “GET / HTTP/1.1 admin”
 d3 = “HTTP root directory”
 d4 = “GET root12, admin123@”
 d5 = “GET root12, admin123@”

(a) Data stream



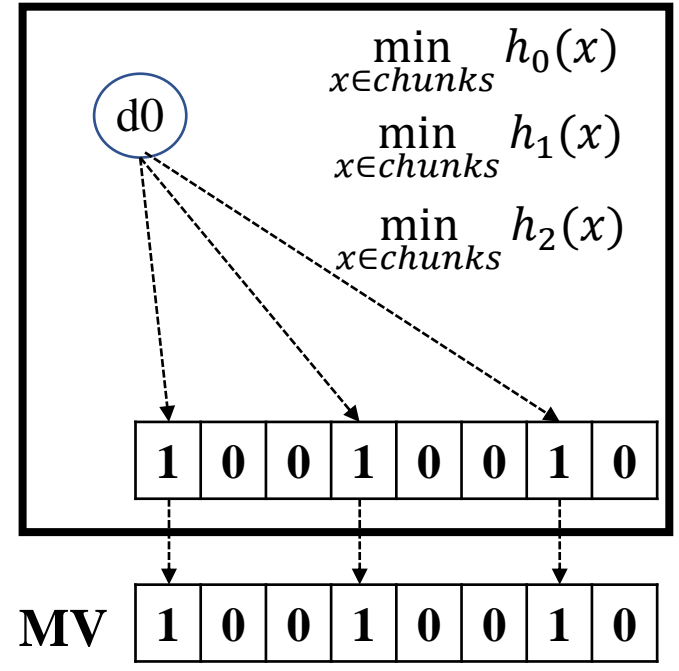
(b) Counter array

$\theta_J=0.6$

JIG: Jaccard-Index Grouping

- d0 = “GET /index.html HTTP/1.1”
- d1 = “HTTProotadmin1234@”
- d2 = “GET / HTTP/1.1 admin”
- d3 = “HTTP root directory”
- d4 = “GET root12, admin123@”
- d5 = “GET root12, admin123@”

(a) Data stream



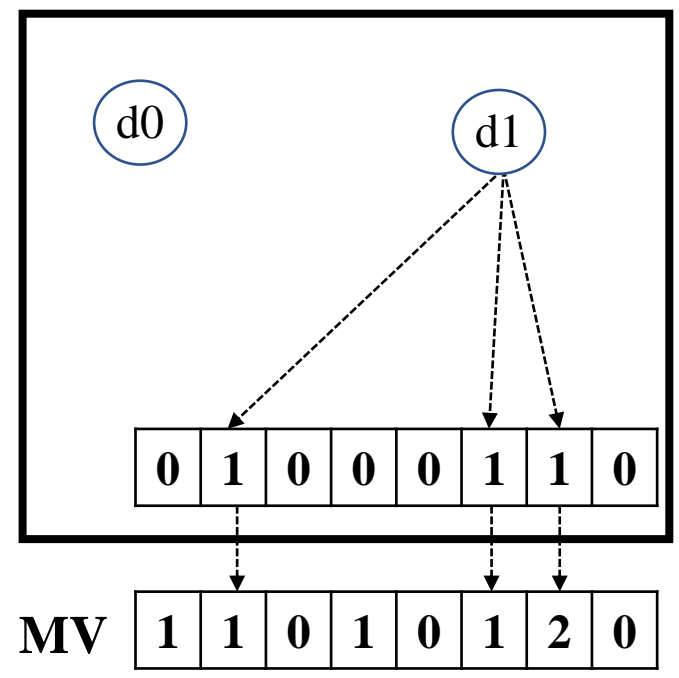
(b) Counter array

$\theta_J=0.6$

JIG: Jaccard-Index Grouping

d0 = “GET /index.html HTTP/1.1”
d1 = “HTTProotadmin1234@”
d2 = “GET / HTTP/1.1 admin”
d3 = “HTTP root directory”
d4 = “GET root12, admin123@”
d5 = “GET root12, admin123@”

(a) Data stream



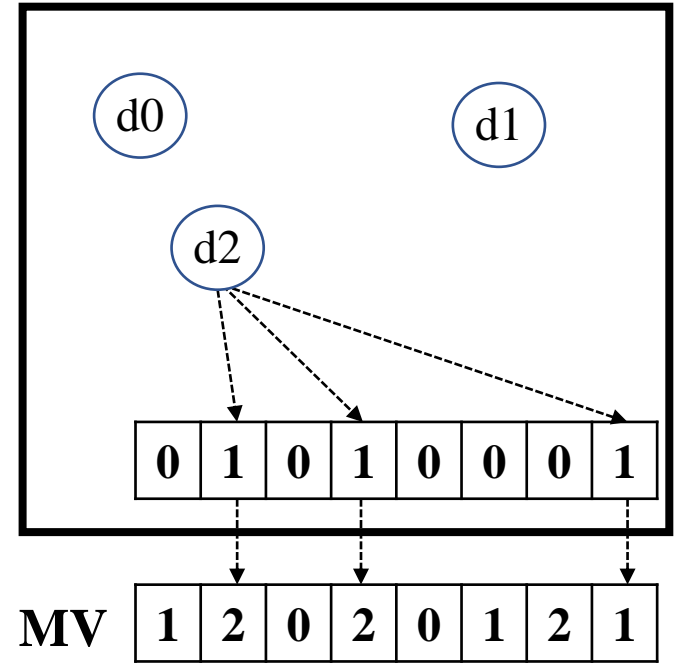
(b) Counter array

$\theta_J=0.6$

JIG: Jaccard-Index Grouping

d0 = “GET /index.html HTTP/1.1”
d1 = “HTTProotadmin1234@”
d2 = “GET / HTTP/1.1 admin”
d3 = “HTTP root directory”
d4 = “GET root12, admin123@”
d5 = “GET root12, admin123@”

(a) Data stream



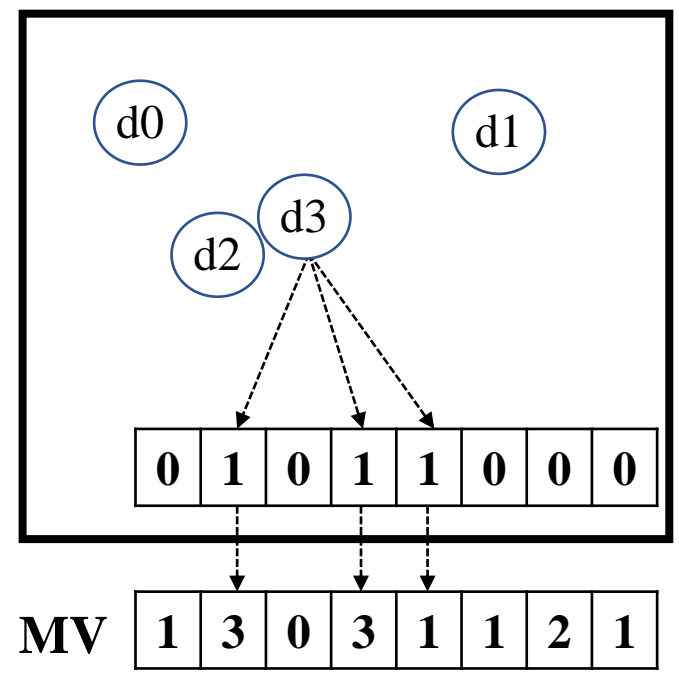
(b) Counter array

$\theta_J=0.6$

JIG: Jaccard-Index Grouping

d0 = “GET /index.html HTTP/1.1”
 d1 = “HTTProotadmin1234@”
 d2 = “GET / HTTP/1.1 admin”
d3 = “HTTP root directory”
 d4 = “GET root12, admin123@”
 d5 = “GET root12, admin123@”

(a) Data stream



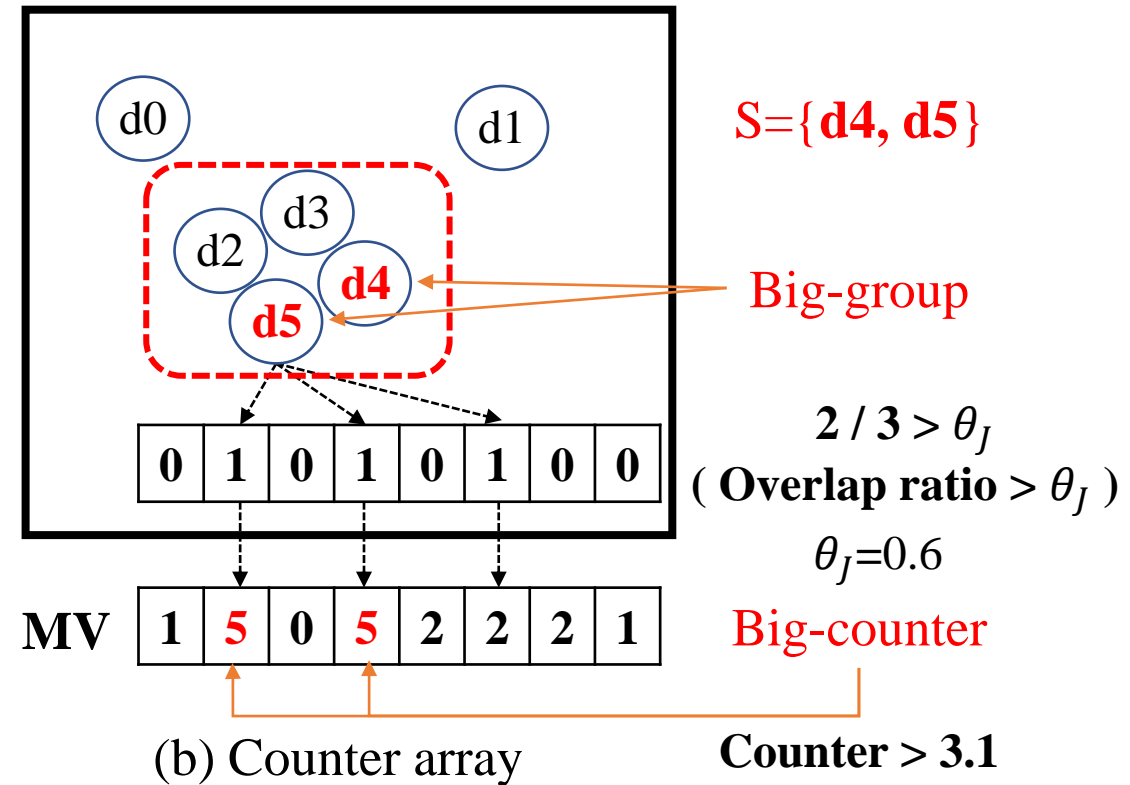
(b) Counter array

$\theta_J=0.6$

JIG: Jaccard-Index Grouping

d0 = "GET /index.html HTTP/1.1"
 d1 = "HTTProotadmin1234@"
 d2 = "GET / HTTP/1.1 admin"
 d3 = "HTTP root directory"
 d4 = "GET root12, admin123@"
d5 = "GET root12, admin123@"

(a) Data stream



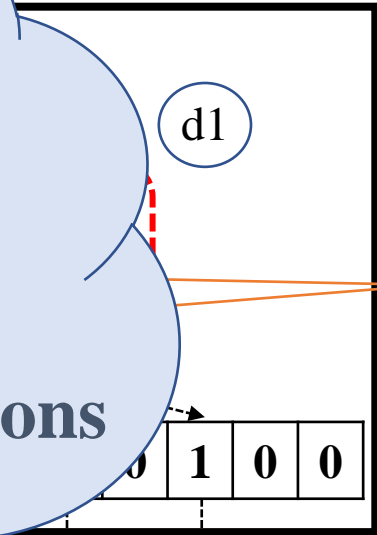
Mean + 6 * Sigma = 3.1
 (Iterative Outlier Removal Algorithm)

JIG: Jaccard-Index Grouping

d0 = "GET
 d1 = "HE
 d2 = "
 d3 = "
 d4 = "C
d5 = "C

Big-group identification

- small fixed amount of memory
- constant number of hash operations



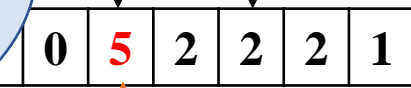
$S = \{d4, d5\}$

Big-group

$2 / 3 > \theta_J$
 (Overlap ratio $> \theta_J$)

$\theta_J = 0.6$

Big-counter



Counter > 3.1

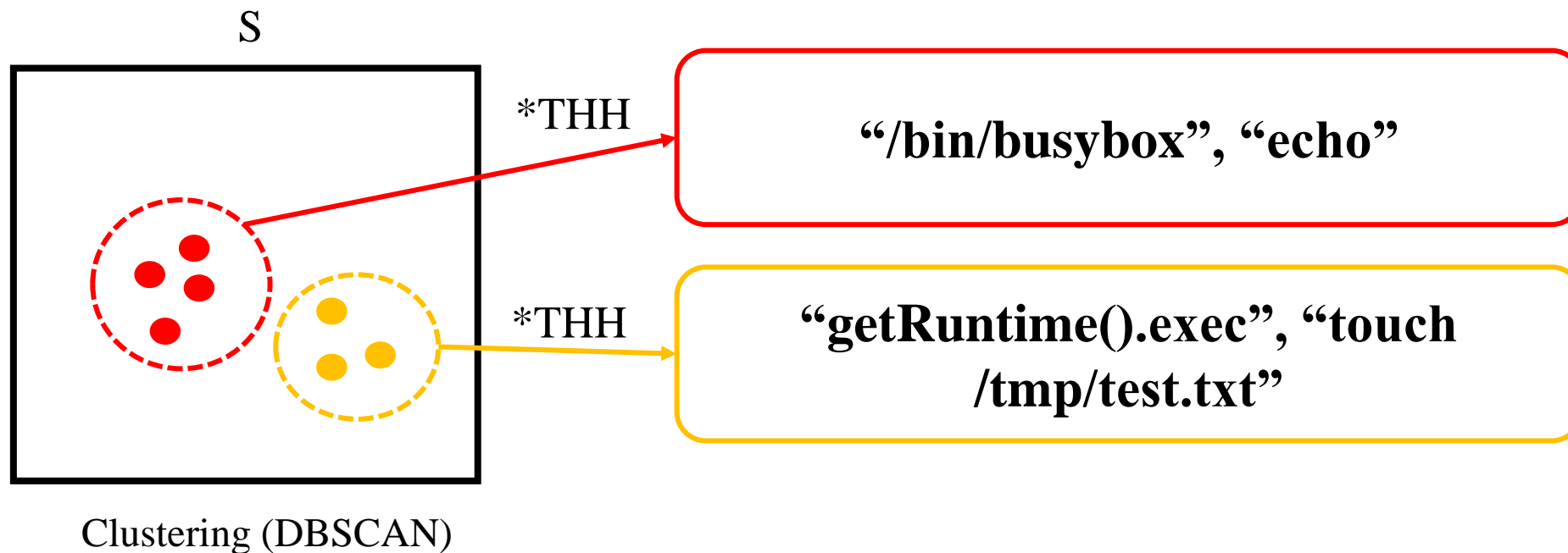
(a) Data stream

(b) Counter array

Mean + 6 * Sigma = 3.1
 (Iterative Outlier Removal Algorithm)

SG2: Signature-Group Generation

- Signature-group generation from S
 - Clustering & Signature extraction (Triple-Heavy-Hitter*)
 - Other schemes can be used instead



AWL: Automatic WhiteListing

- AWL
 - (SG2 THH) – (Global THH)

{“HTTP/1.1”, “application”,
“/bin/busybox”,
“getRuntime().exec”}

(a) SG2



{“HTTP/1.1”, “POST”,
“application”}

(b) AWL



{“/bin/busybox”,
“getRuntime().exec”}

(c) GIPS

Experiments

- Dataset
 - SIM1-1 ~ SIM3-3
 - Simulated dataset for big-group identification
 - <https://tinyurl.com/2gfz4f7v>
 - IoT1~8
 - A.K.A. IoT23, <https://www.stratosphereips.org/datasets-iot23>
 - IDS1~8
 - A.K.A. CICIDS2017, <https://www.unb.ca/cic/datasets/ids-2017.html>
 - ISP1~3
 - Private dataset
 - Suspicious packets from IPS

Experiments

- SIM1-1 ~ SIM3-3
 - <https://tinyurl.com/2gfz4f7v>

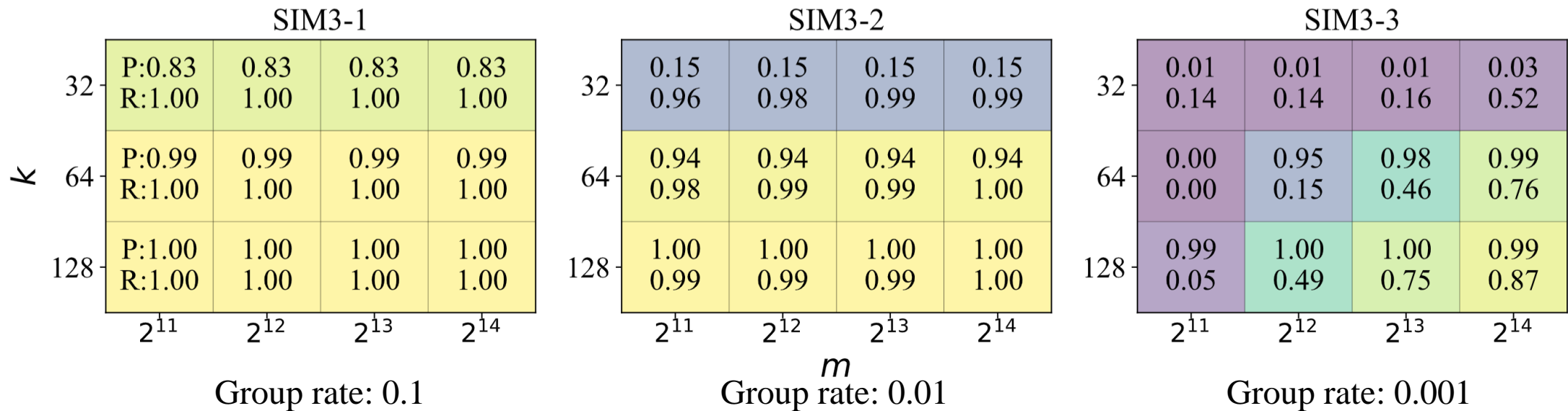


Figure 6: Precision (P) and Recall (R) of JIG for SIM datasets with different big-group ratios, k , and m .

Experiments

- IoT1~8
 - IoT23, <https://www.stratosphereips.org/datasets-iot23>
 - Table4 for Precision & Recall

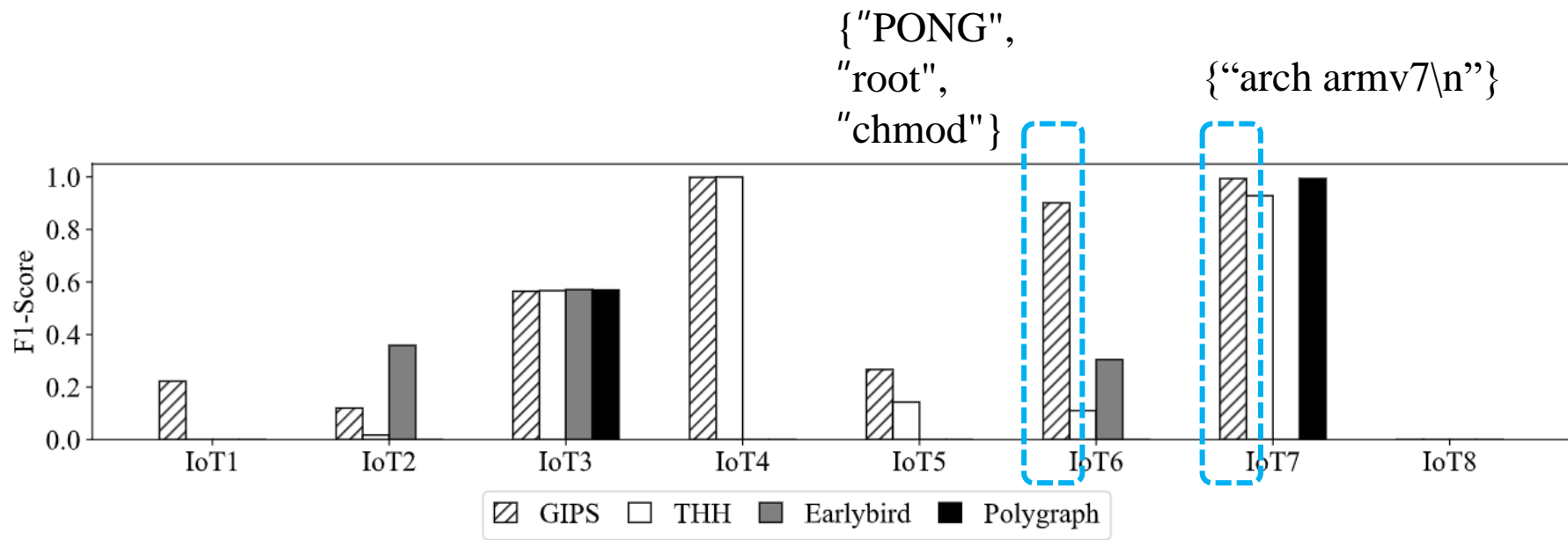


Figure 7: F1-Score of GIPS, THH [9], Earlybird [48], and Polygraph [39] for IoT datasets.

Experiments

- IoT1~8
 - IoT23, <https://www.stratosphereips.org/datasets-iot23>
 - Table4 for Precision & Recall

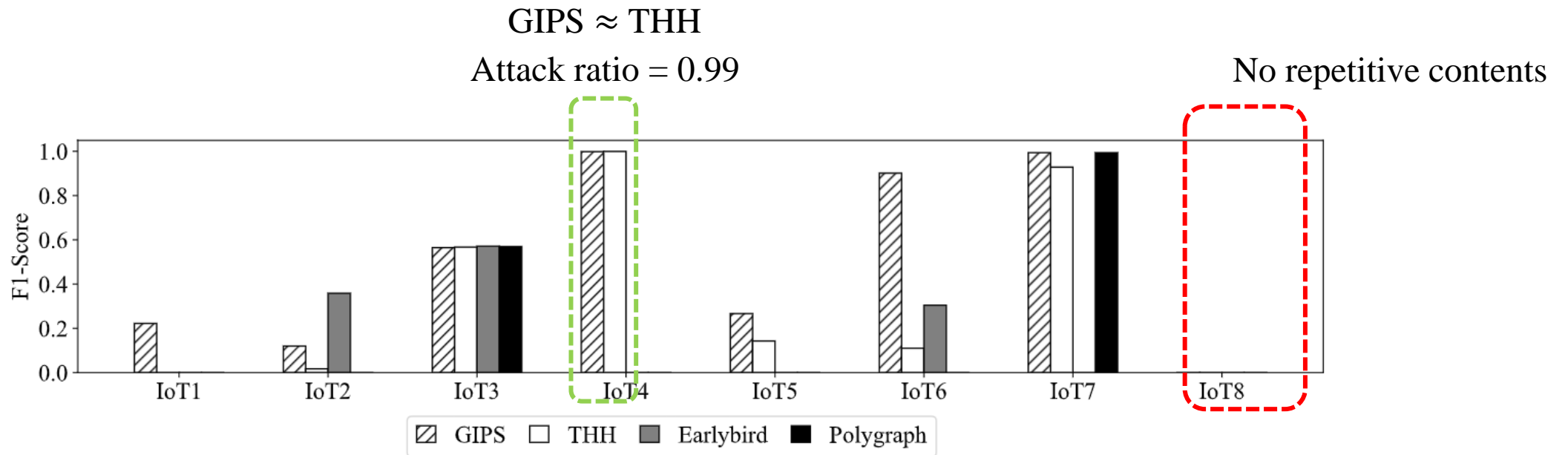


Figure 7: F1-Score of GIPS, THH [9], Earlybird [48], and Polygraph [39] for IoT datasets.

Experiments

- IDS1~8
 - CICIDS2017, <https://www.unb.ca/cic/datasets/ids-2017.html>
 - Table5 for Precision & Recall

{"/dv/vulnerabilities", "Cookie: security=low"}

{"GET / HTTP/1.0\r\n\r\n\r\n\r\n"}

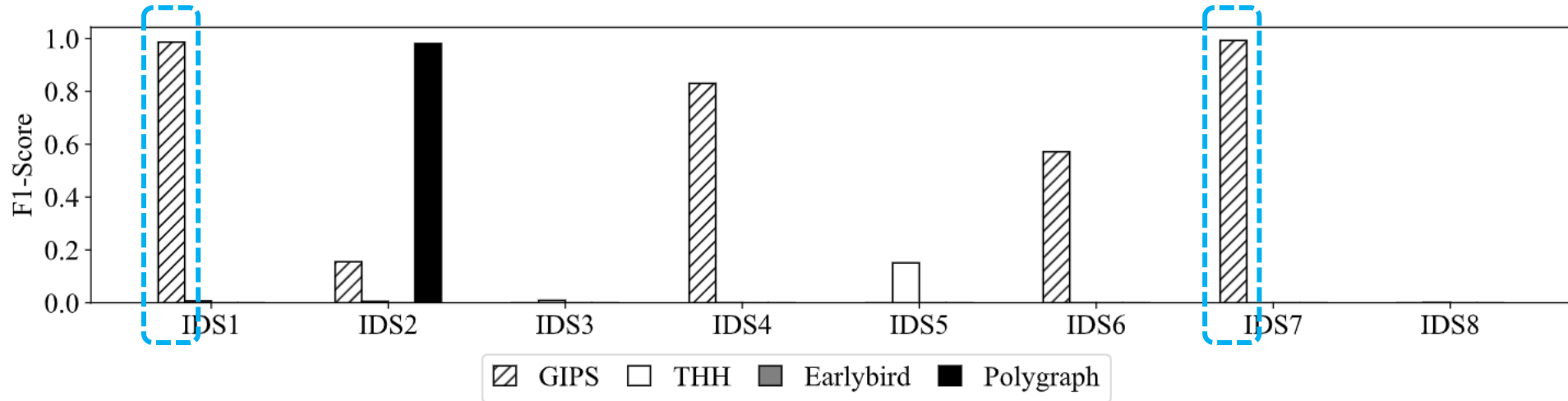


Figure 10: F1-Score of GIPS, THH [9], Earlybird [48], and Polygraph [39] for IDS datasets.

Experiments

- IDS1~8
 - CICIDS2017, <https://www.unb.ca/cic/datasets/ids-2017.html>
 - Table5 for Precision & Recall

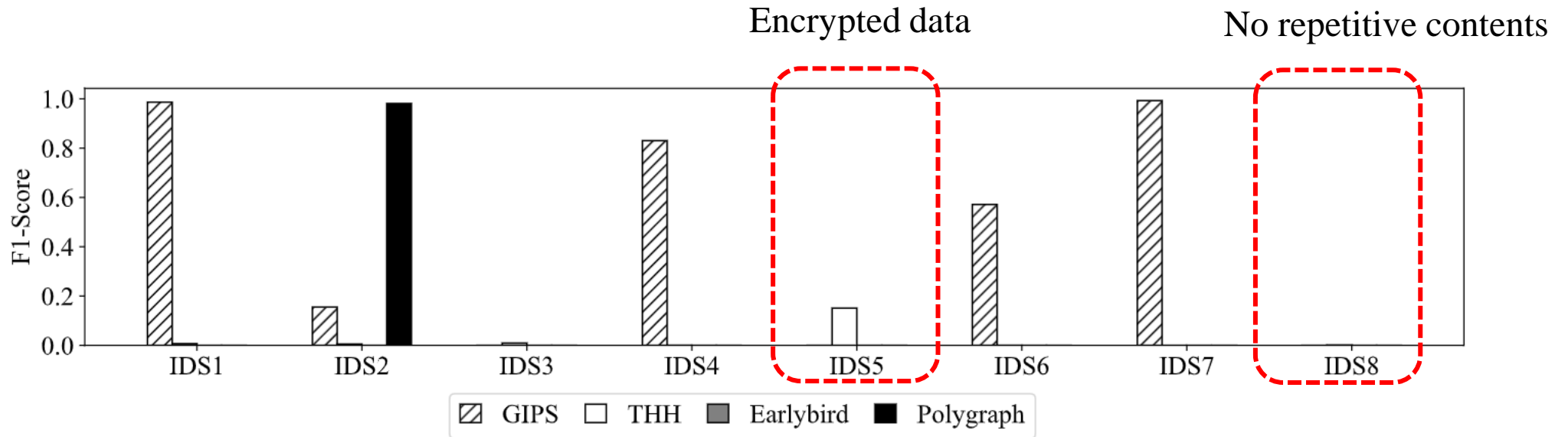


Figure 10: F1-Score of GIPS, THH [9], Earlybird [48], and Polygraph [39] for IDS datasets.

Experiments

- ISP1~3
 - Private dataset
 - Table6 for Precision & Recall

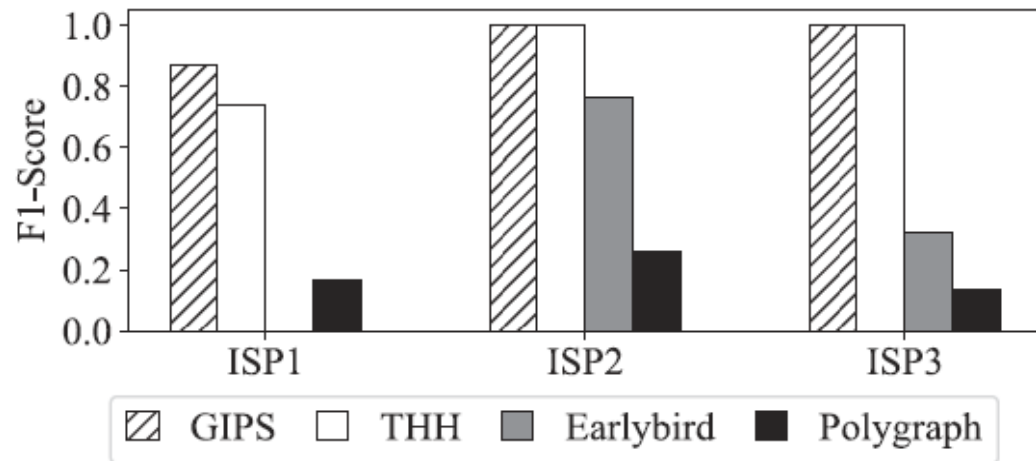
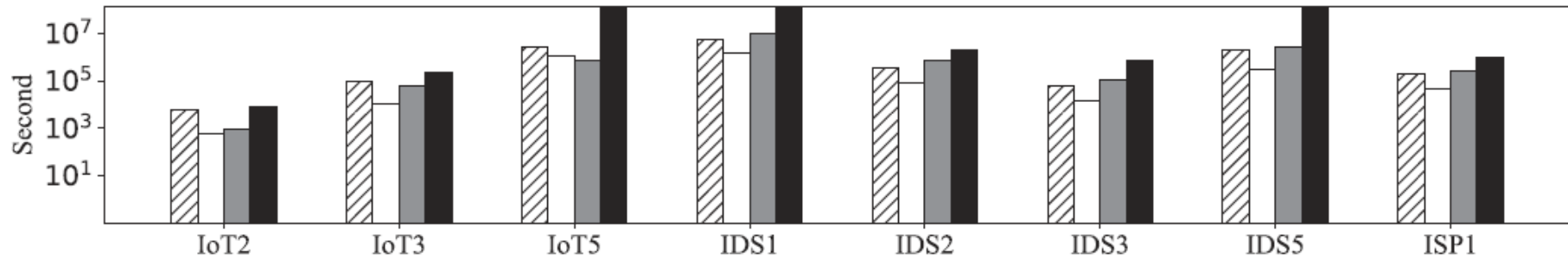


Table 7: GIPS Signatures and Attacks for ISP Datasets

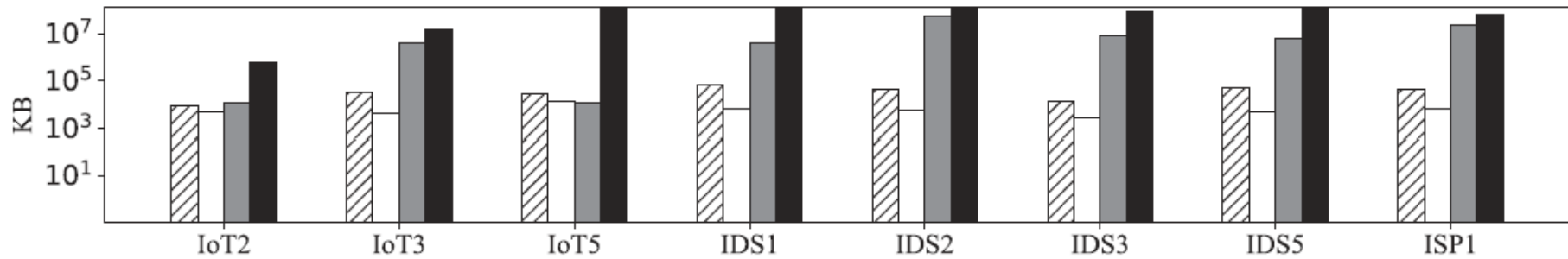
Dataset	Signature	Related Attack
ISP1	<code>getRuntime().exec("touch /tmp/test.txt"), /bin/busybox</code>	CVE-2022-22963 [40] CVE-2017-17215 [37]
ISP2	<code>0002736c0000ff</code>	CVE-2010-1574 [36]
ISP3	<code>\x00\x00\x00DHIP\x00... "Port" : 37777,"RemoteVideo...</code>	Reflection attack [41]

Experiments

- Processing time and memory usage



(a) Comparison of Processing Time



(b) Comparison of Memory Usage



Conclusion

- A new generative IDS/IPS to mitigate persistent zero-day attacks
- Fast and compact identification of big-groups from data streams
- Automatic generation of robust group-signatures
- Work from a high to a low attack ratio, for example from 0.99 to 0.001
- Limitations
 - Encrypted data, repetitive contents
 - Data analysis tool rather than real-time system

Thank You