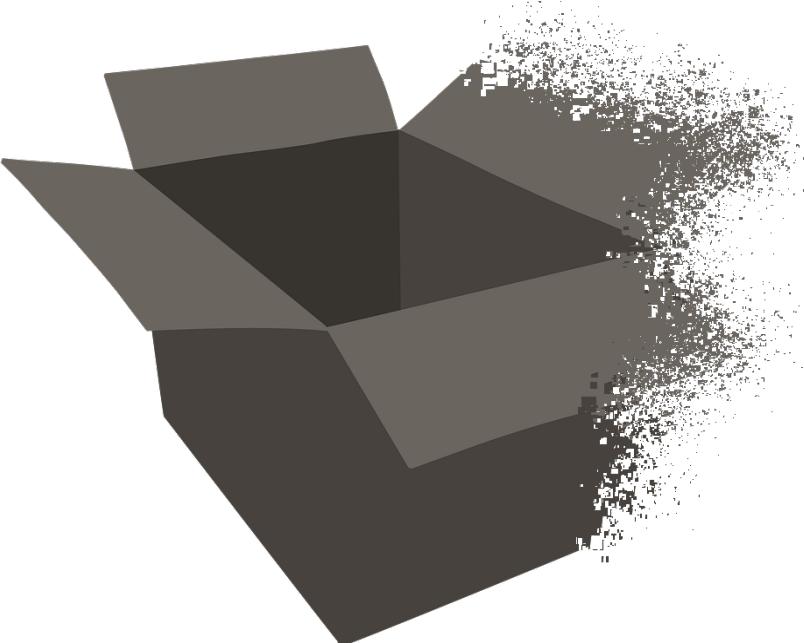


Faculty 1: Chair of IT-Security | Master Thesis Colloquium



# **Representation Learning for Content-Sensitive Anomaly Detection in Industrial Networks**

*B. Sc. Fabian Kopp*

***Supervisor:***

*M. Sc. Franka Schuster*

***Examination Board:***

*Prof. Dr.-Ing. Andriy Panchenko*

*Prof. Dr. rer. nat. Oliver Hohlfeld*

## 0 | Preamble

1 | Background

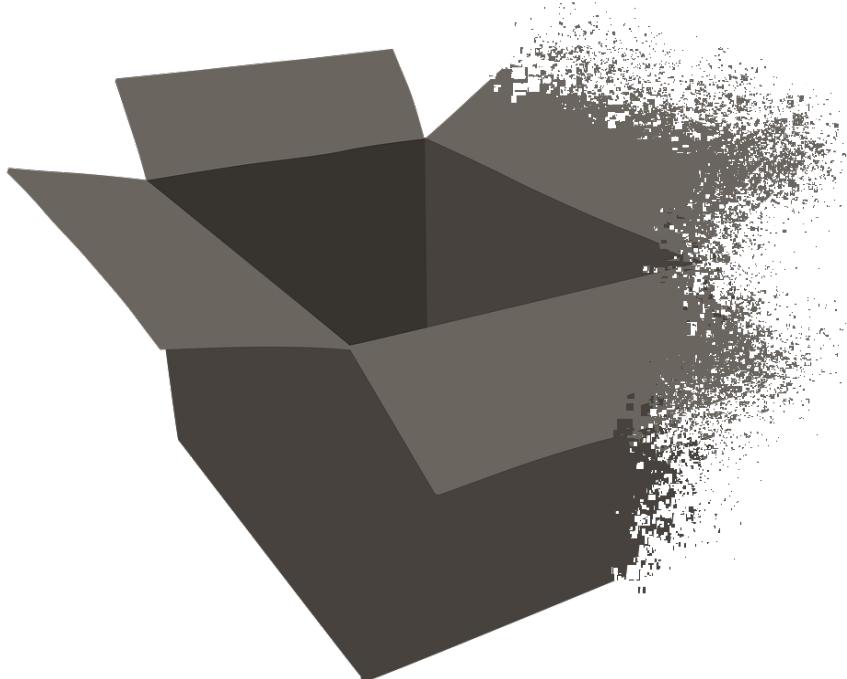
2 | Related Work

3 | Methodology

4 | Implementation

5 | Evaluation

6 | Summary



- Threat defense is an *open* problem
- Computer networks generate lots of data
- Artificial neural networks are *powerful* computational systems

- Threat defense is an *open* problem
- Computer networks generate lots of data
- Artificial neural networks are *powerful* computational systems

→ ***Can one use this data to detect attacks against a network in an automatic fashion ?***

- Threat defense is an *open* problem
- Computer networks generate lots of data
- Artificial neural networks are *powerful* computational systems

→ ***Can one use this data to detect attacks against a network in an automatic fashion ?***  
Yes, with *intrusion detection systems* that rely on *anomaly detection* !

- Threat defense is an *open* problem
- Computer networks generate lots of data
- Artificial neural networks are *powerful* computational systems

→ ***Can one use this data to detect attacks against a network in an automatic fashion ?***  
Yes, with *intrusion detection systems* that rely on *anomaly detection* !



Focus on **feature extraction**  
and model **interpretability**

- Threat defense is an *open* problem
- Computer networks generate lots of data
- Artificial neural networks are *powerful* computational systems

→ ***Can one use this data to detect attacks against a network in an automatic fashion ?***  
Yes, with *intrusion detection systems* that rely on *anomaly detection* !



Focus on **feature extraction**  
and model **interpretability**

- Threat defense is an *open* problem
- Computer networks generate lots of data
- Artificial neural networks are *powerful* computational systems

→ ***Can one use this data to detect attacks against a network in an automatic fashion ?***  
Yes, with *intrusion detection systems* that rely on *anomaly detection* !

Focus on **feature extraction**  
and model **interpretability**

- Automatic transformation of input data into a set of derived values (features)
- Features intended to be informative & non-redundant
- Extracted features are basis for any downstream task

**Representation Learning** for Content-Sensitive Anomaly Detection in Industrial Networks

- Threat defense is an *open* problem
- Computer networks generate lots of data
- Artificial neural networks are *powerful* computational systems

→ ***Can one use this data to detect attacks against a network in an automatic fashion ?***  
Yes, with *intrusion detection systems* that rely on *anomaly detection* !

Focus on **feature extraction**  
and model **interpretability**

– Automatic transformation of input data  
into a set of derived values (features)  
– Features intended to be informative & non-redundant  
– Extracted features are basis for any downstream task

– Crucial for the detection of  
unknown attacks  
– Detection based on raw data

- Threat defense is an *open* problem
- Computer networks generate lots of data
- Artificial neural networks are *powerful* computational systems

→ ***Can one use this data to detect attacks against a network in an automatic fashion ?***  
Yes, with *intrusion detection systems* that rely on *anomaly detection* !

Focus on **feature extraction**  
and model **interpretability**

– Automatic transformation of input data  
into a set of derived values (features)  
– Features intended to be informative & non-redundant  
– Extracted features are basis for any downstream task

– Crucial for the detection of  
unknown attacks  
– Detection based on raw data

– Valuable asset to protect  
– Fixed network topology

**Representation Learning** for **Content-Sensitive Anomaly Detection** in **Industrial Networks**

I. Are computing systems able to perform for a spatial-temporal extraction of feature from raw data ?

- I. Are computing systems able to perform for a spatial-temporal extraction of feature from raw data ?
- II. Can anomaly detection be improved when trained on extracted features ?

- I. Are computing systems able to perform for a spatial-temporal extraction of feature from raw data ?
- II. Can anomaly detection be improved when trained on extracted features ?
- III. What type of input data representations yields better results for different anomalies ?

- I. Are computing systems able to perform for a spatial-temporal extraction of feature from raw data ?
- II. Can anomaly detection be improved when trained on extracted features ?
- III. What type of input data representations yields better results for different anomalies ?
- IV. How can the semantic gap be closed with the help of explainable models ?

0 | Preamble

**1 | Background**

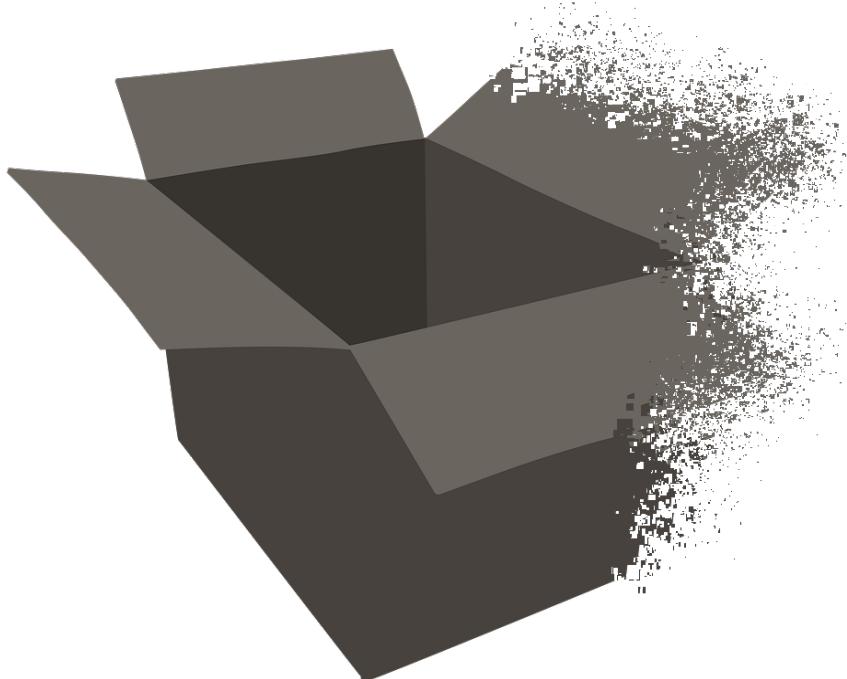
2 | Related Work

3 | Methodology

4 | Implementation

5 | Evaluation

6 | Summary

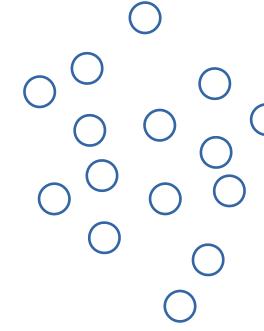


## Anomaly detection (AD)

- Unsupervised one class classification
- *Learn* normal behavior - alert deviations
- Used in: fraud detection, intrusion detection

## Anomaly detection (AD)

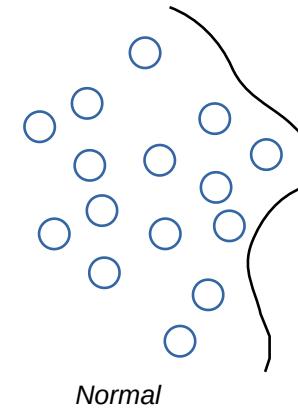
- Unsupervised one class classification
- *Learn* normal behavior - alert deviations
- Used in: fraud detection, intrusion detection



*Normal*

## Anomaly detection (AD)

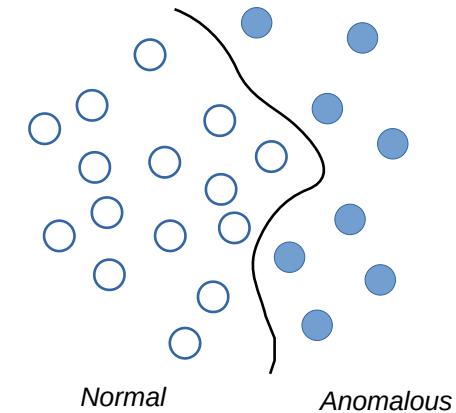
- Unsupervised one class classification
- *Learn* normal behavior - alert deviations
- Used in: fraud detection, intrusion detection



Normal

## Anomaly detection (AD)

- Unsupervised one class classification
- *Learn* normal behavior - alert deviations
- Used in: fraud detection, intrusion detection

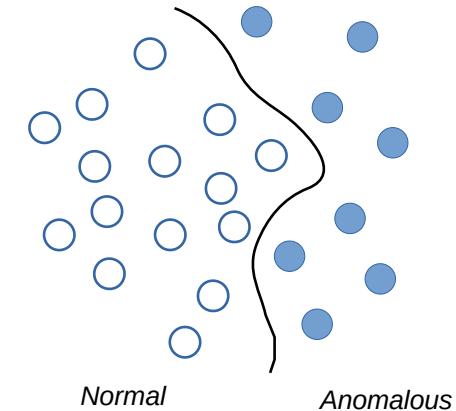


## Anomaly detection (AD)

- Unsupervised one class classification
- *Learn* normal behavior - alert deviations
- Used in: fraud detection, intrusion detection

## Types of anomalies

- (Un)known attacks
- Human error
- Insider abuse
- Faulty hardware

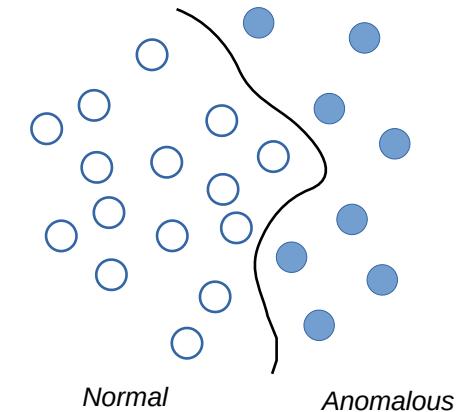


## Anomaly detection (AD)

- Unsupervised one class classification
- *Learn* normal behavior - alert deviations
- Used in: fraud detection, intrusion detection

## Types of anomalies

- (Un)known attacks
- Human error
- Insider abuse
- Faulty hardware



## Anomaly classes

- Intra-packet anomalies

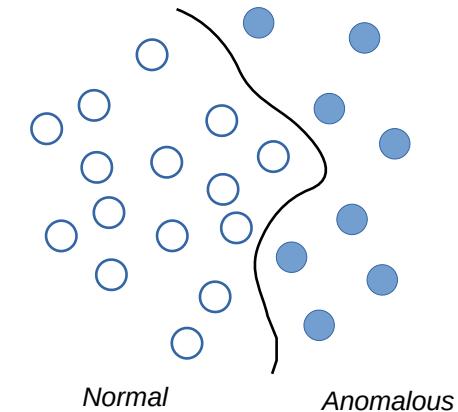


## Anomaly detection (AD)

- Unsupervised one class classification
- *Learn* normal behavior - alert deviations
- Used in: fraud detection, intrusion detection

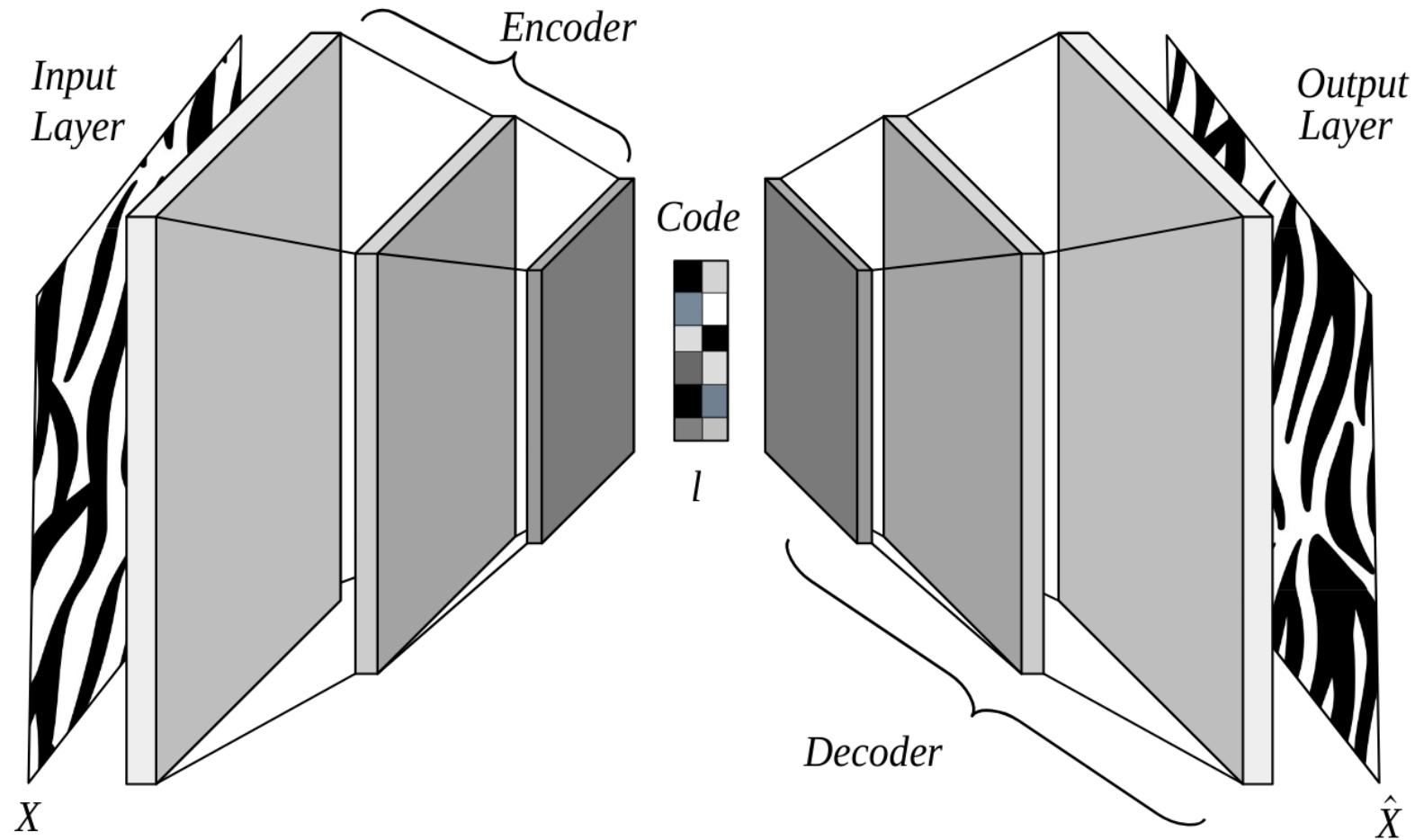
## Types of anomalies

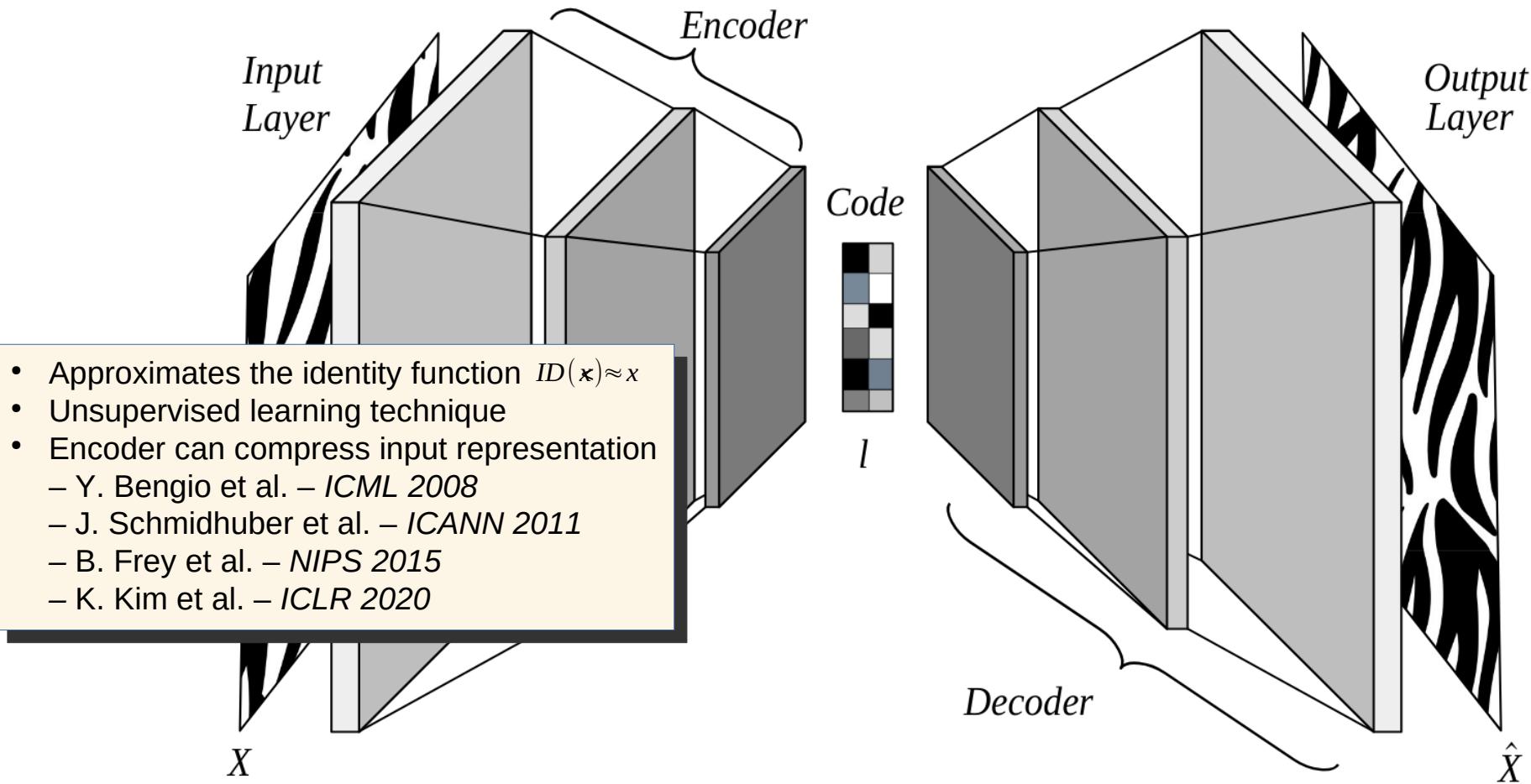
- (Un)known attacks
- Human error
- Insider abuse
- Faulty hardware



## Anomaly classes

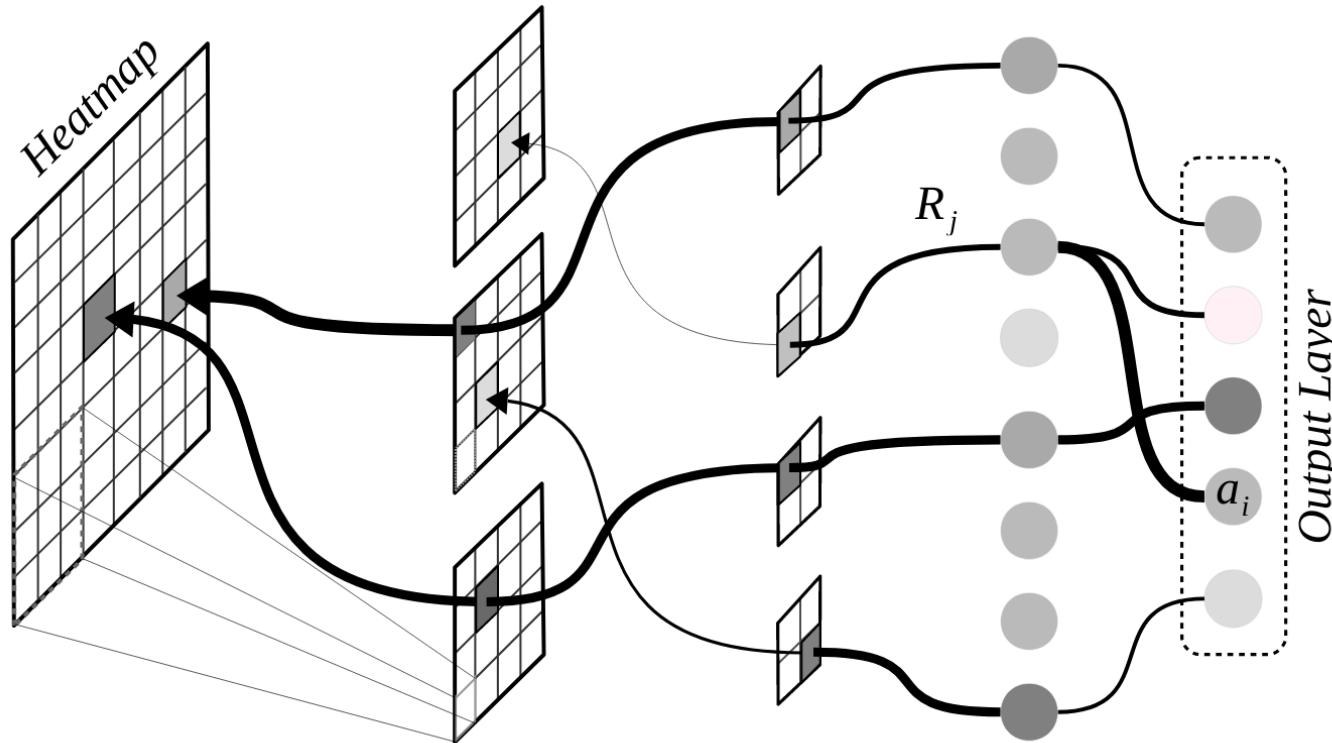
- Intra-packet anomalies
- Inter-packet anomalies





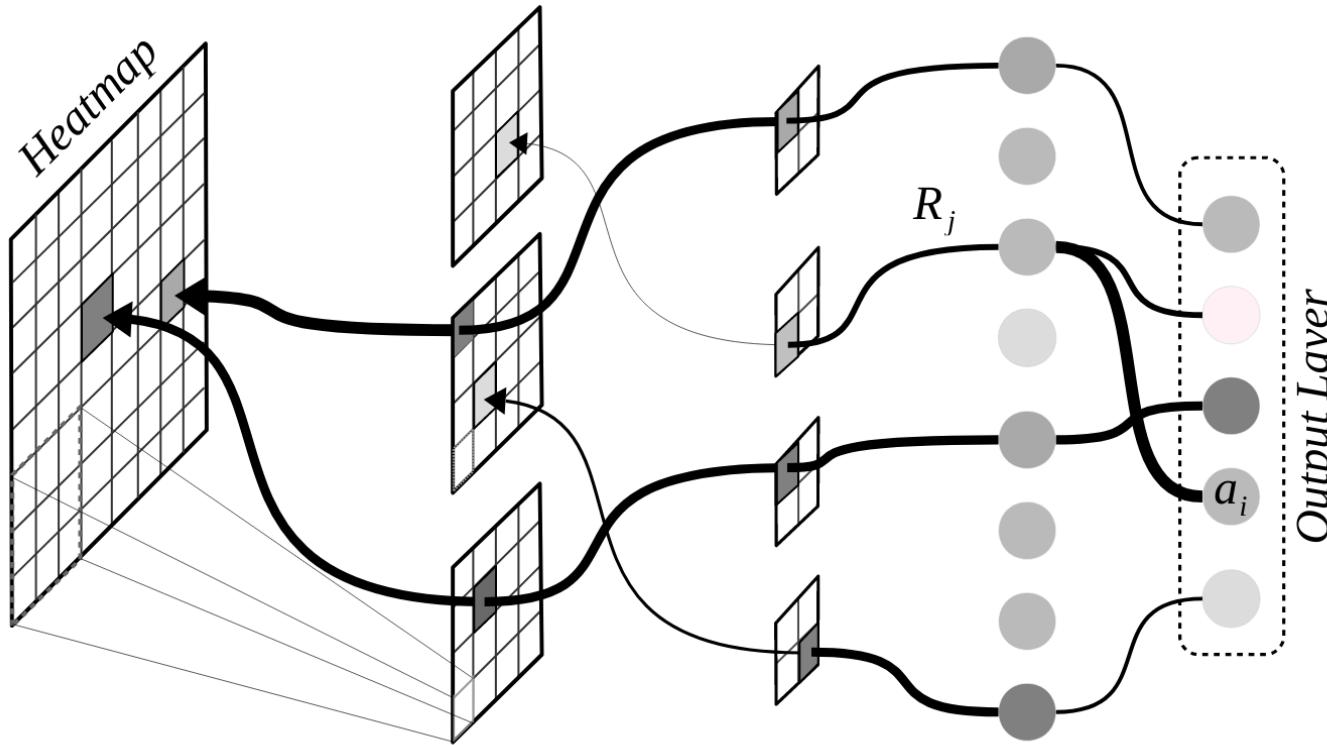
→ **Model Inference:** “*What class is present in a given image ?*”

← **Decision Analysis:** “*Why is a given image classified as a cat ?*”



→ **Model Inference:** “What class is present in a given image ?”

← **Decision Analysis:** “Why is a given image classified as a cat ?”



→ **Model Inference:** “What class is present in a given image ?”

← **Decision Analysis:** “Why is a given image classified as a cat ?”

$$a_k = \max(0, \sum_0^j a_j w_{jk}) \text{ with } a_0 = 1$$

$$R_j = \sum_k^j \frac{a_j w_{jk}}{\sum_0^j a_j w_{jk}} R_k \text{ where } R_0 = 1$$

Relevance Propagation Rule  
A. Binder et al. – ICANN 2016

0 | Preamble

1 | Background

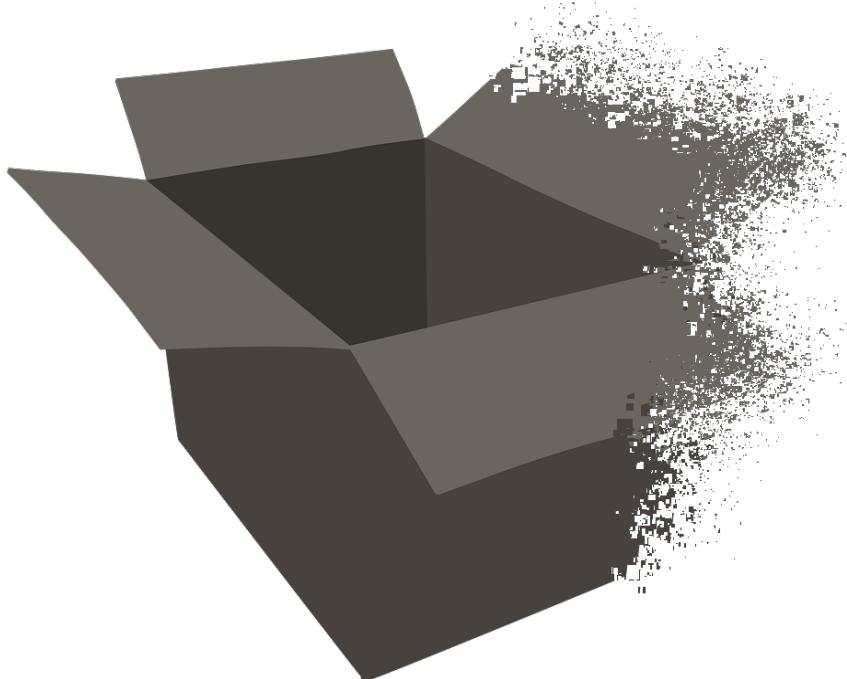
**2 | Related Work**

3 | Methodology

4 | Implementation

5 | Evaluation

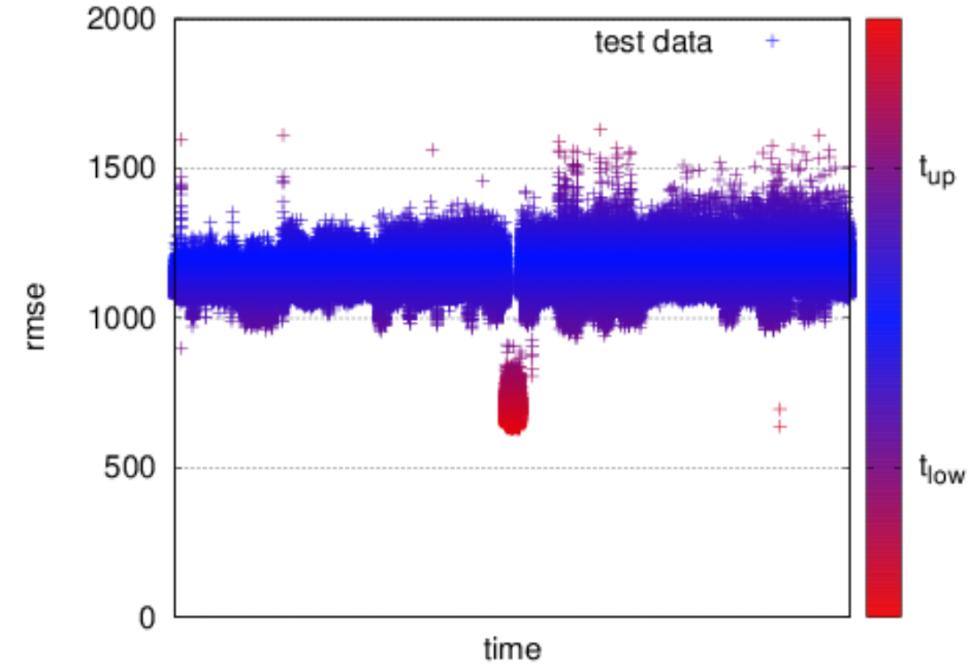
6 | Summary



## High-Performance Unsupervised Anomaly Detection for Cyber-Physical System Networks

[P. Schneider et al. – CPS-SPC 2018]

- 3 layer linear denosing autoencoder
- First 1000 bytes used for detection
- Loss-based anomaly detection
- Unsupervised end-to-end approach (+)
- No direct temporal learning (-)
- Drawbacks on detecting short anomalies (-)



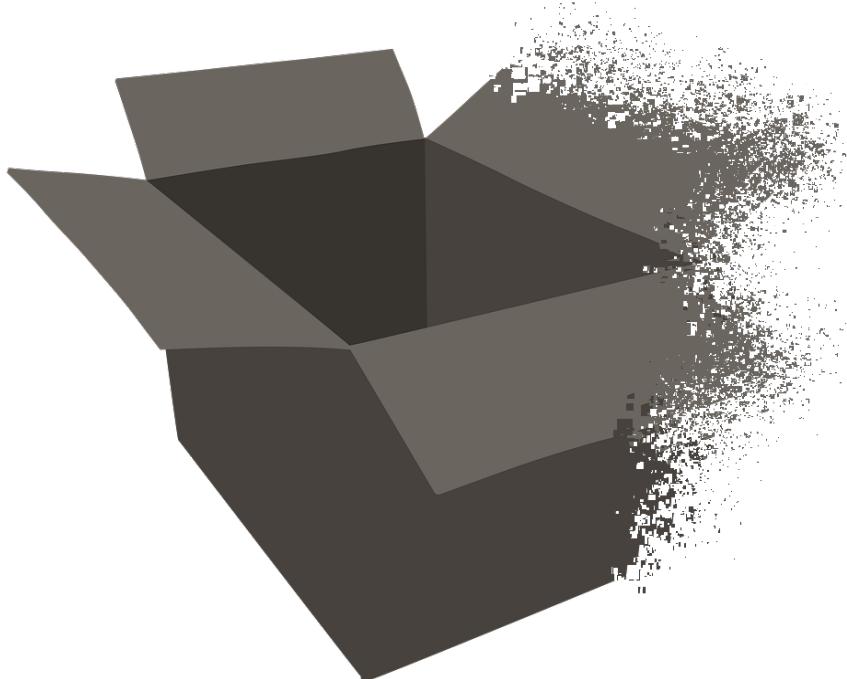
## Attentional Payload Anomaly Detector for Web Applications

[ZQ. Qin et al. – ICONIP 2018]

- Anomaly detection based on raw bytes sequences
- TCP payloads are basis for detection
- Attention-based RNN model
- Supervised training (-)
- Heuristic to visualize the model's interpretation (+)

```
GET /dv/vulnerabilities/sqli/?id=1%27+and+1%3D1%23&Submit=Submit HTTP/1.1\r\nHost: 205.174.165.68\r\nUser-Agent: Mozilla/5.0 (X11; Linux x86_64; rv:45.0)\r\nGecko/20100101 Firefox/45.0\r\nAccept: text/html,application/xhtml+xml,application/xml;q=0.9,*/*;q=0.8\r\nAccept-Language: en-US,en;q=0.5\r\nAccept-Encoding: gzip, deflate\r\nReferer: http://205.174.165.68/dv/vulnerabilities/sqli/\r\nCookie: security=low; PHPSESSID=5dfcuh85kg0vvvidf8nrsjtbob5\r\nConnection: keep-alive\r\n\r\n
```

- 0 | Preamble
- 1 | Background
- 2 | Related Work
- 3 | Methodology**
- 4 | Implementation
- 5 | Evaluation
- 6 | Summary



- Extract features from raw frames for a content sensitive AD
  - Enables a protocol-agnostic processing
  - Comparison between *byte*- and *flow*-oriented input data

54	50	21	31	20	31	0d	0a	48	6f	73	74	3a	20	78	6e
2d	2d	6d	62	69	75	73	2d	6a	75	01	20	62	61	6e	64
0d	0a	55	73	65	72	2d	41	67	05	6e	74	3a	20	4d	6f
7a	69	6c	6c	61	2f	35	2c	30	20	20	58	31	31	3b	20
55	62	75	6e	74	75	3b	20	4c	69	6e	75	78	20	78	38
36	51	36	34	3b	20	72	76	3a	38	31	2e	30	20	20	47
65	63	6b	6f	2f	32	30	31	30	30	31	30	31	20	4d	69
72	65	66	6f	78	2f	38	31	2c	30	0d	0a	41	63	63	65
70	74	3a	20	74	65	78	74	21	68	74	6d	6c	2c	6	70
70	6c	69	63	61	74	69	6f	6e	21	78	68	74	6d	6c	2b
78	6d	6c	2d	61	70	70	6c	69	63	61	74	69	6f	6e	21
78	6d	6c	3b	71	2d	30	2c	39	20	69	6d	61	67	65	21
77	65	62	70	2c	2a	2f	2a	3b	71	3c	30	2c	38	0d	0a
41	63	63	65	70	74	2d	4c	64	6e	67	75	61	67	65	3a
20	65	6e	2d	47	2c	65	6e	3b	71	3d	30	2c	35	0d	0a
0a	41	63	63	65	70	74	2d	45	6e	63	6f	64	69	6e	67
3a	20	67	7a	69	70	2c	20	64	65	66	6c	61	74	65	0d
0a	44	4e	54	3a	20	31	0d	0a	43	6f	6e	6e	65	63	74
69	6f	6e	3a	20	6b	65	65	70	2d	61	6c	69	76	65	72
0c	55	70	67	72	61	64	65	2d	49	6e	73	65	63	75	72
65	2d	52	65	71	75	65	73	74	73	3a	20	31	0d	0a	49
66	2d	4d	6f	64	69	66	69	65	64	2d	53	69	6e	63	65
3a	20	46	72	69	2c	20	31	31	20	53	65	70	20	32	30
32	30	20	30	38	3a	34	32	3a	32	37	20	47	4d	54	0d
0a	49	66	2d	4e	6f	6e	65	2d	4d	61	74	63	68	3a	20
57	27	22	35	63	61	2d	35	61	66	30	35	61	61	32	37
36	34	39	39	22	0d	0a	43	62	63	68	65	2d	43	6f	6e
74	72	6f	6c	3a	20	6d	6	78	2d	61	67	65	3c	30	0d
0a	0d	0a	cfc	7f	90	5t	10	1b	08	00	0d	00	00	00	00

TP/1.1 · Host: xn--mbius-jua.band · User-Agent: Mozilla/5.0 (X11; Ubuntu; Linux x86\_64; rv:81.0) Gecko/20100101 Firefox/81.0 · Accept: text/html,application/xhtml+xml,application/xml,application/xml;q=0.9,image/webp,\*/\*;q=0.8 · Accept-Language: en-GB,en;q=0.5 · Accept-Encoding: gzip, deflate · DNT: 1 · Connection: keep-alive · Upgrade-Insecure-Requests: 1 · If-Modified-Since: Fri, 11 Sep 2020 08:42:27 GMT · If-None-Match: W/"5ca-5af05aa276499" · Cache-Control: max-age=0 ·

- Extract features from raw frames for a content sensitive AD
  - Enables a protocol-agnostic processing
  - Comparison between *byte*- and *flow*-oriented input data
- Mixture of different types of neurons for versatile features
  - *Convolutional* neurons for spatial extraction  
→ intra-packet anomalies
  - *Recurrent* neurons for temporal extraction  
→ inter-packet anomalies

54	50	21	31	20	31	0d	0a	48	6f	73	74	3a	20	78	6e
2d	2d	6d	62	69	75	73	2d	6a	75	01	26	62	61	6e	64
0d	0a	55	73	65	72	2d	4	67	05	6e	74	3a	20	4d	6f
7a	69	6c	6d	61	35	2c	30	20	26	58	31	31	3b	20	47
55	62	75	6e	74	75	3b	20	4c	69	6e	75	78	20	78	38
36	51	36	34	3b	20	72	76	3a	38	31	2e	30	20	20	47
6d	63	6b	6f	2f	32	30	31	30	30	31	30	31	20	4d	69
72	65	66	6f	78	2f	38	31	2c	30	0d	0a	4	63	63	65
70	74	3a	20	74	65	78	74	21	68	74	6d	6c	2c	6	70
70	6c	69	63	61	74	69	6f	6e	21	78	68	74	6d	6c	2b
78	6d	6c	2d	61	70	70	6c	69	63	61	74	69	6f	6e	21
78	6d	6c	3b	71	2d	30	2	39	20	69	6d	61	67	65	21
77	65	62	70	2c	2a	2f	2a	3b	71	3c	30	2c	38	0d	0a
4	63	63	65	70	74	2d	4c	6	6e	67	75	61	67	65	3a
20	65	6e	2d	47	2c	65	6e	3b	71	3d	30	2c	35	0d	0a
0a	71	63	65	65	70	74	2d	45	6e	63	6f	64	69	6e	67
3a	20	67	7a	69	70	2c	20	64	65	66	6c	61	74	65	0d
0a	44	4e	54	3a	20	31	0d	0a	43	6f	6e	6e	65	63	74
69	6f	6e	3a	20	6b	65	65	70	2d	61	6c	69	76	65	72
0	55	70	67	72	61	64	65	2d	49	6e	73	65	63	75	72
65	2d	52	65	71	75	65	73	74	73	3a	20	31	0d	0a	49
66	2d	4d	6f	64	69	66	69	65	64	2d	53	69	6e	63	65
3a	20	46	72	69	2c	20	31	31	20	53	65	70	20	32	30
32	30	20	30	38	3a	34	32	3a	32	37	20	47	4d	54	0d
0a	49	66	2d	4e	6f	6e	65	2d	4d	61	74	63	68	3a	20
57	27	22	35	63	61	2d	35	61	66	30	35	61	61	32	37
36	34	39	39	22	0d	0a	43	6	63	68	65	2d	43	6f	6e
74	72	6f	6c	3a	20	6d	6	78	2d	61	67	65	3	30	0d
0a	0d	0a	cfc	7f	90	5t	10	1b	08	00	00	00	00	00	00

TP/1.1 · Host: xn--mbius-jua.band · User-Agent: Mozilla/5.0 (X11; Ubuntu; Linux x86\_64; rv:81.0) Gecko/20100101 Firefox/81.0 · Accept: text/html,application/xhtml+xml,application/xml;q=0.9,image/webp,\*/\*;q=0.8 · Accept-Language: en-GB,en;q=0.5 · Accept-Encoding: gzip, deflate · DNT: 1 · Connection: keep-alive · Upgrade-Insecure-Requests: 1 · If-Modified-Since: Fri, 11 Sep 2020 08:42:27 GMT · If-None-Match: W/"5ca-5af05aa276499" · Cache-Control: max-age=0 ·

- Extract features from raw frames for a content sensitive AD
  - Enables a protocol-agnostic processing
  - Comparison between *byte*- and *flow*-oriented input data
- Mixture of different types of neurons for versatile features
  - *Convolutional* neurons for spatial extraction  
→ intra-packet anomalies
  - *Recurrent* neurons for temporal extraction  
→ inter-packet anomalies
- Unsupervised feature learning based on autoencoders
  - Leverage the encoder to extract features

54	50	21	31	20	31	0d	0a	48	6f	73	74	3a	20	78	6e
2d	2d	6d	62	69	75	73	2d	6a	75	01	20	62	61	6e	64
0d	0a	55	73	65	72	2d	41	67	05	6e	74	3a	20	4d	6f
7a	69	6c	6d	61	35	2e	30	20	20	58	31	31	3b	20	47
55	62	75	6e	74	75	3b	20	4e	69	6e	75	78	20	78	38
36	51	36	34	3b	20	72	76	3a	38	31	2e	30	20	20	47
65	63	6b	6f	2f	32	30	31	30	30	31	30	31	20	4d	69
72	65	66	6f	78	2f	38	31	2e	30	0d	0a	41	63	63	65
70	74	3a	20	74	65	78	74	21	68	74	6d	6c	2c	6	70
70	6c	69	63	61	74	69	6f	6e	21	78	68	74	6d	6c	2b
78	6d	6c	2d	61	70	70	6c	69	63	61	74	69	6f	6e	21
78	6d	6c	3b	71	3d	30	2e	39	20	69	6d	61	67	65	21
77	65	62	70	2c	2a	2f	2a	3b	71	3d	30	2e	38	0d	0a
41	63	63	65	70	74	2d	4c	61	6e	67	75	61	67	65	3a
20	65	6e	2d	47	2c	65	6e	3b	71	3d	30	2e	35	0d	0a
0a	41	63	63	65	70	74	2d	45	6e	63	6f	64	69	6e	67
3a	20	67	7a	69	70	2c	20	64	65	66	6c	61	74	65	0d
0a	44	4e	54	3a	20	31	0d	0a	43	6f	6e	65	63	74	0d
69	6f	6e	3a	20	6b	65	65	70	2d	61	6c	69	76	65	0d
0a	55	70	67	72	61	64	65	2d	49	6e	73	65	63	75	72
65	2d	52	65	71	75	65	73	74	73	3a	20	31	0d	0a	49
66	2d	4d	6f	64	69	66	69	65	64	2d	53	69	6e	63	65
3a	20	46	72	69	2c	20	31	31	20	53	65	70	20	32	30
32	30	20	30	38	3a	34	32	3a	32	37	20	47	4d	54	0d
0a	49	66	2d	4e	6f	6e	65	2d	4d	61	74	63	68	3a	20
57	27	22	35	63	61	2d	35	61	66	30	35	61	61	32	37
36	34	39	39	22	0d	0a	43	62	63	68	65	2d	43	6f	6e
74	72	6f	6c	3a	20	6d	61	78	2d	61	67	65	3e	30	0d
0a	0d	0a	cfc	7f	90	5t	10	1b	08	00	00	00	00	00	00

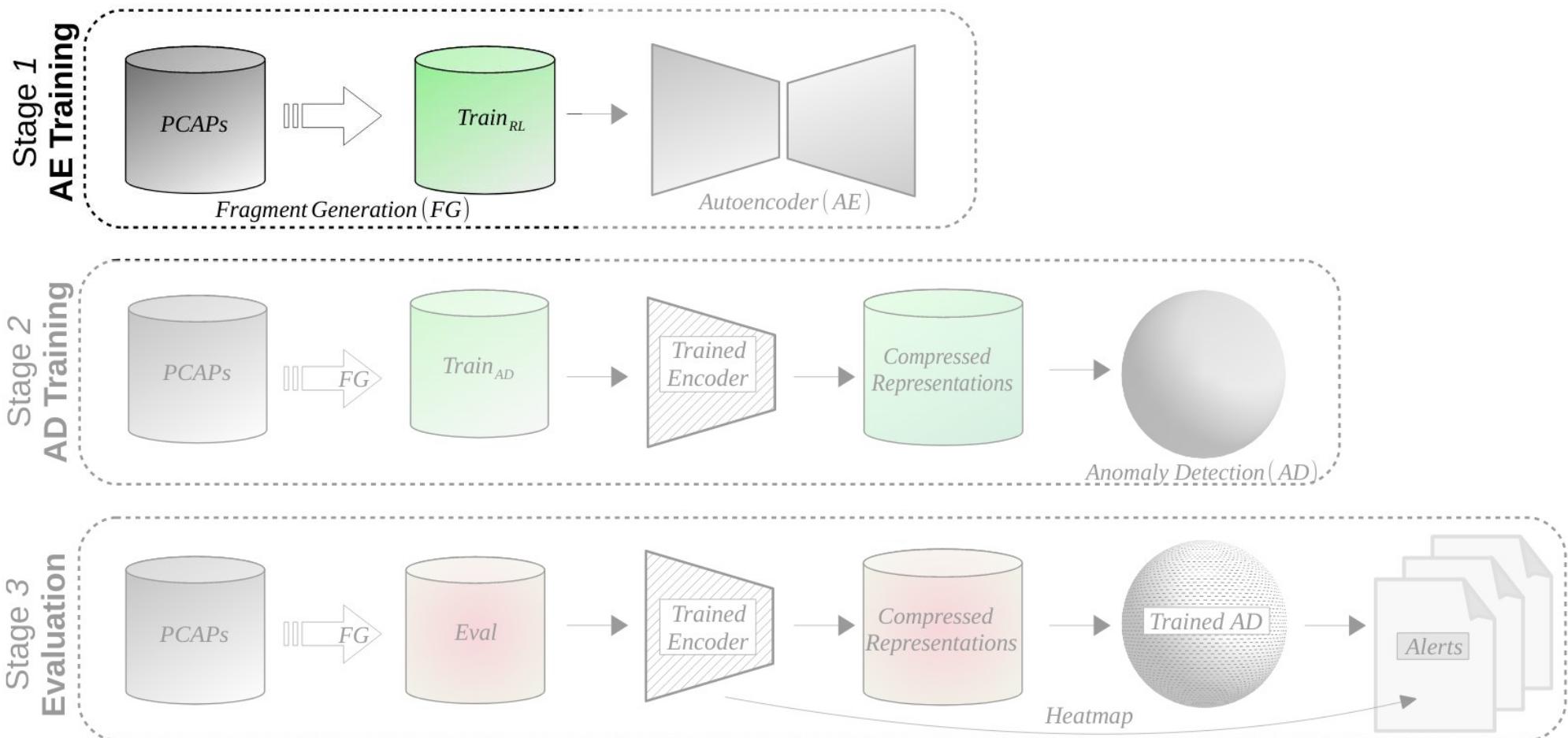
TP/1.1 · Host: xn--mbius-jua.band · User-Agent: Mozilla/5.0 (X11; Ubuntu; Linux x86\_64; rv:81.0) Gecko/20100101 Firefox/81.0 · Accept: text/html,application/xhtml+xml,application/xml;q=0.9,image/webp,\*/\*;q=0.8 · Accept-Language: en-GB,en;q=0.5 · Accept-Encoding: gzip, deflate · DNT: 1 · Connection: keep-alive · Upgrade-Insecure-Requests: 1 · If-Modified-Since: Fri, 11 Sep 2020 08:42:27 GMT · If-None-Match: W/"5ca-5af05aa276499" · Cache-Control: max-age=0 ·

- Extract features from raw frames for a content sensitive AD
  - Enables a protocol-agnostic processing
  - Comparison between *byte*- and *flow*-oriented input data
- Mixture of different types of neurons for versatile features
  - *Convolutional* neurons for spatial extraction  
→ intra-packet anomalies
  - *Recurrent* neurons for temporal extraction  
→ inter-packet anomalies
- Unsupervised feature learning based on autoencoders
  - Leverage the encoder to extract features
- Use layer-wise relevance propagation for raised alerts
  - Helps to interpret the model's decisions

54	50	21	31	20	31	0d	0a	48	6f	73	74	3a	20	78	6e
2d	2d	6d	62	69	75	73	2d	6a	75	01	20	62	61	6e	64
0d	0a	55	73	65	72	2d	4	67	05	6e	74	3a	20	4d	6f
7a	69	6c	6d	61	3f	35	2e	30	20	20	58	31	31	3b	20
55	62	75	6e	74	75	3b	20	4e	69	6e	75	78	20	78	38
36	51	36	34	3b	20	72	76	3a	38	31	2e	30	20	20	47
65	63	6b	6f	2f	32	30	31	30	30	31	30	31	20	4d	69
72	65	66	6f	78	2f	38	31	3e	30	0d	0a	4	63	63	65
70	74	3a	20	74	65	78	74	21	68	74	6d	6c	2c	6	70
70	6c	69	63	61	74	69	6f	6e	21	78	68	74	6d	6c	2b
78	6d	6c	2d	61	70	70	6c	69	63	61	74	69	6f	6e	21
78	6d	6c	3b	71	3d	30	2	39	20	69	6d	61	67	65	21
77	65	62	70	2c	2a	2f	2a	3b	71	3d	30	2	38	0d	0a
4	63	63	65	70	74	2d	4c	6	6e	67	75	61	67	65	3a
20	65	6e	2d	47	2c	65	6e	3b	71	3d	30	2	35	0d	0a
0a	71	63	63	65	70	74	2d	45	6e	63	6f	64	69	6e	67
3a	20	67	7a	69	70	2c	20	64	65	66	6c	61	74	65	0d
0a	44	4e	54	3a	20	31	0d	0a	43	6f	6e	65	63	74	0d
69	6f	6e	3a	20	6b	65	65	70	2d	61	6c	69	76	65	0d
0	55	70	67	72	61	64	65	2d	49	6e	73	65	63	75	72
65	2d	52	65	71	75	65	73	74	73	3a	20	31	0d	0a	49
66	2d	4d	6f	64	69	66	69	65	64	2d	53	69	6e	63	65
3a	20	46	72	69	2c	20	31	31	20	53	65	70	20	32	30
32	30	20	30	38	3a	34	32	3a	32	37	20	47	4d	54	0d
0a	49	66	2d	4e	6f	6e	65	2d	4d	61	74	63	68	3a	20
57	27	22	35	63	61	2d	35	61	66	30	35	61	61	32	37
36	34	39	39	22	0d	0a	43	6	63	68	65	2d	43	6f	6e
74	72	6f	6c	3a	20	6d	6	78	2d	61	67	65	3	30	0d
0a	0d	0a	cfc	71	90	5t	10	1b	08	00	00	00	00	00	0d

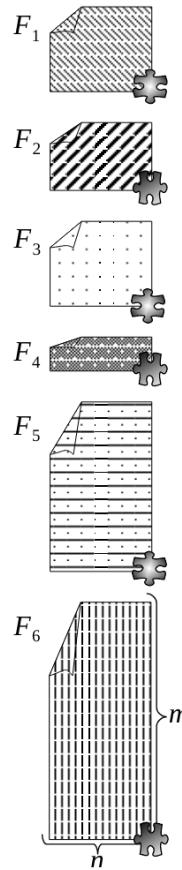
TP/1.1 · Host: xn--mbius-jua.band · User-Agent: Mozilla/5.0 (X11; Ubuntu; Linux x86\_64; rv:81.0) Gecko/20100101 Firefox/81.0 · Accept: text/html,application/xhtml+xml,application/xml;q=0.9,image/webp,\*/\*;q=0.8 · Accept-Language: en-GB,en;q=0.5 · Accept-Encoding: gzip, deflate · DNT: 1 · Connection: keep-alive · Upgrade-Insecure-Requests: 1 · If-Modified-Since: Fri, 11 Sep 2020 08:42:27 GMT · If-None-Match: W/"5ca-5af05aa276499" · Cache-Control: max-age=0 ·

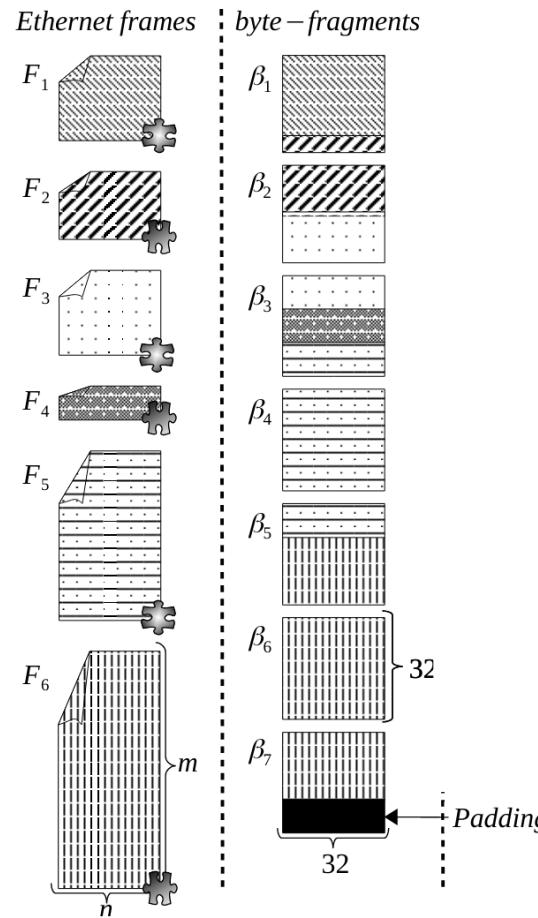
### 3 | Methodology: Overview



### 3 | Methodology: Fragment Generation

*Ethernet frames*

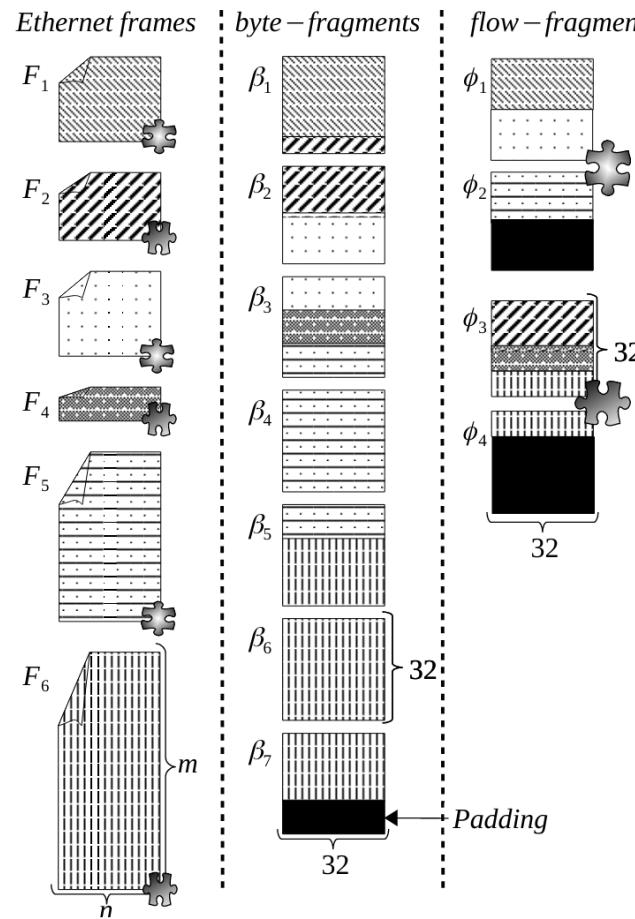




- Byte-Fragments

- Jumping sliding-window extracted fixed fragments
- Processes the whole input stream of frames

### 3 | Methodology: Fragment Generation

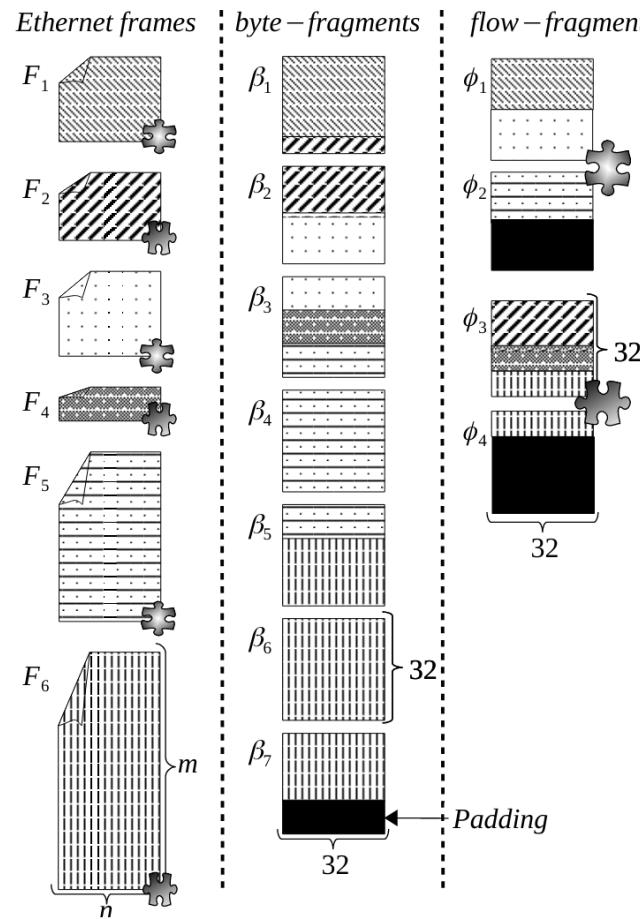


- Byte-Fragments

- Jumping sliding-window extracted fixed fragments
- Processes the whole input stream of frames

- Flow-Fragments

- Sorts input stream by *network flow*
  - (Transport Protocol, src:IP, src:Port, dst:IP, dst:Port)

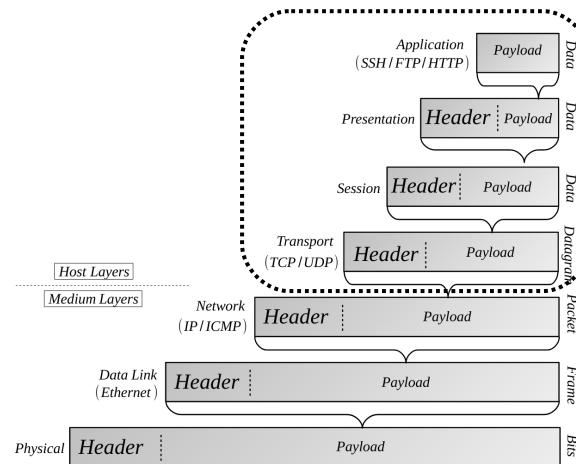


- Byte-Fragments

- Jumping sliding-window extracted fixed fragments
- Processes the whole input stream of frames

- Flow-Fragments

- Sorts input stream by *network flow*  
→ (Transport Protocol, src:IP, src:Port, dst:IP, dst:Port)
- Use *first 64 bytes of every network layer payload*



### 3 | Methodology: Fragment Generation

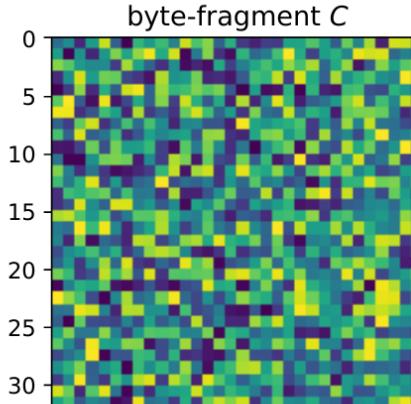
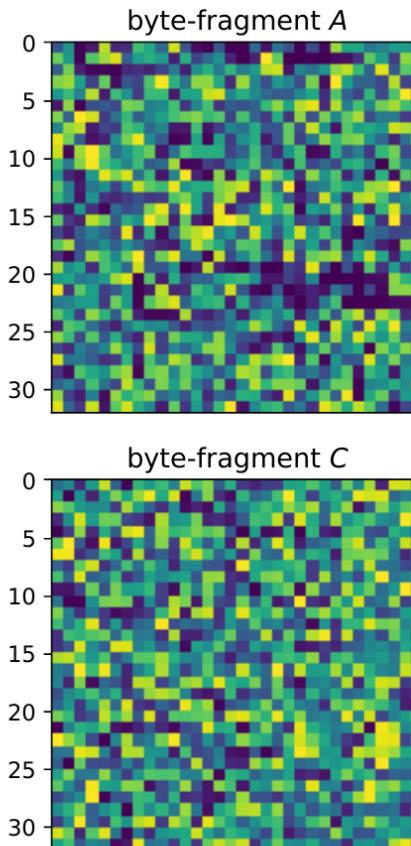


Figure 5.2: Four consecutive examples of byte-fragments shaded using a *viridis* color map

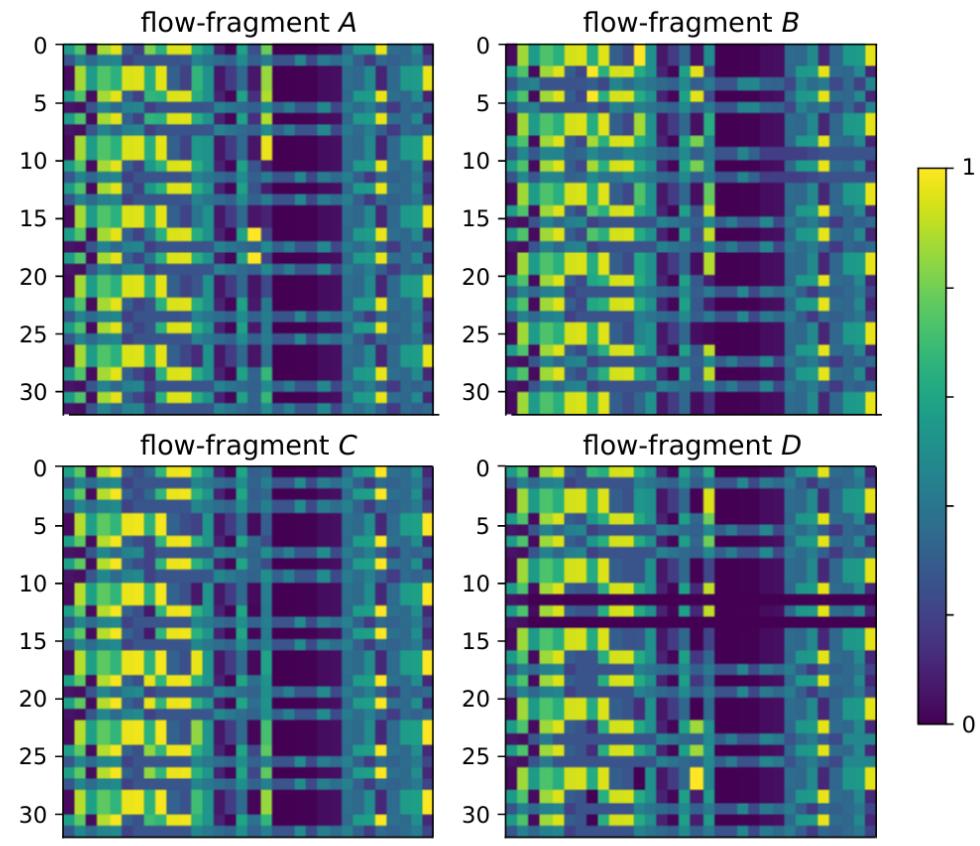
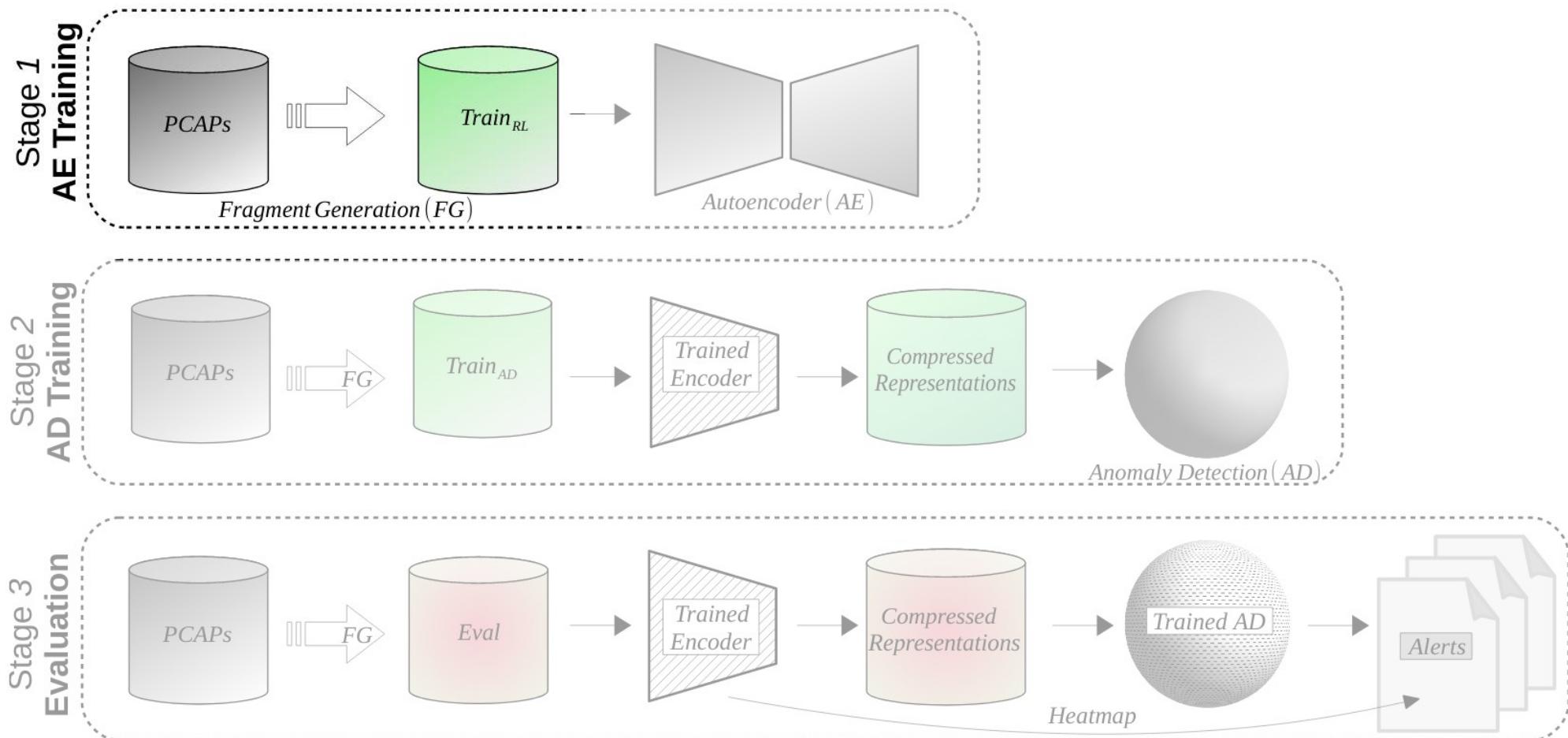
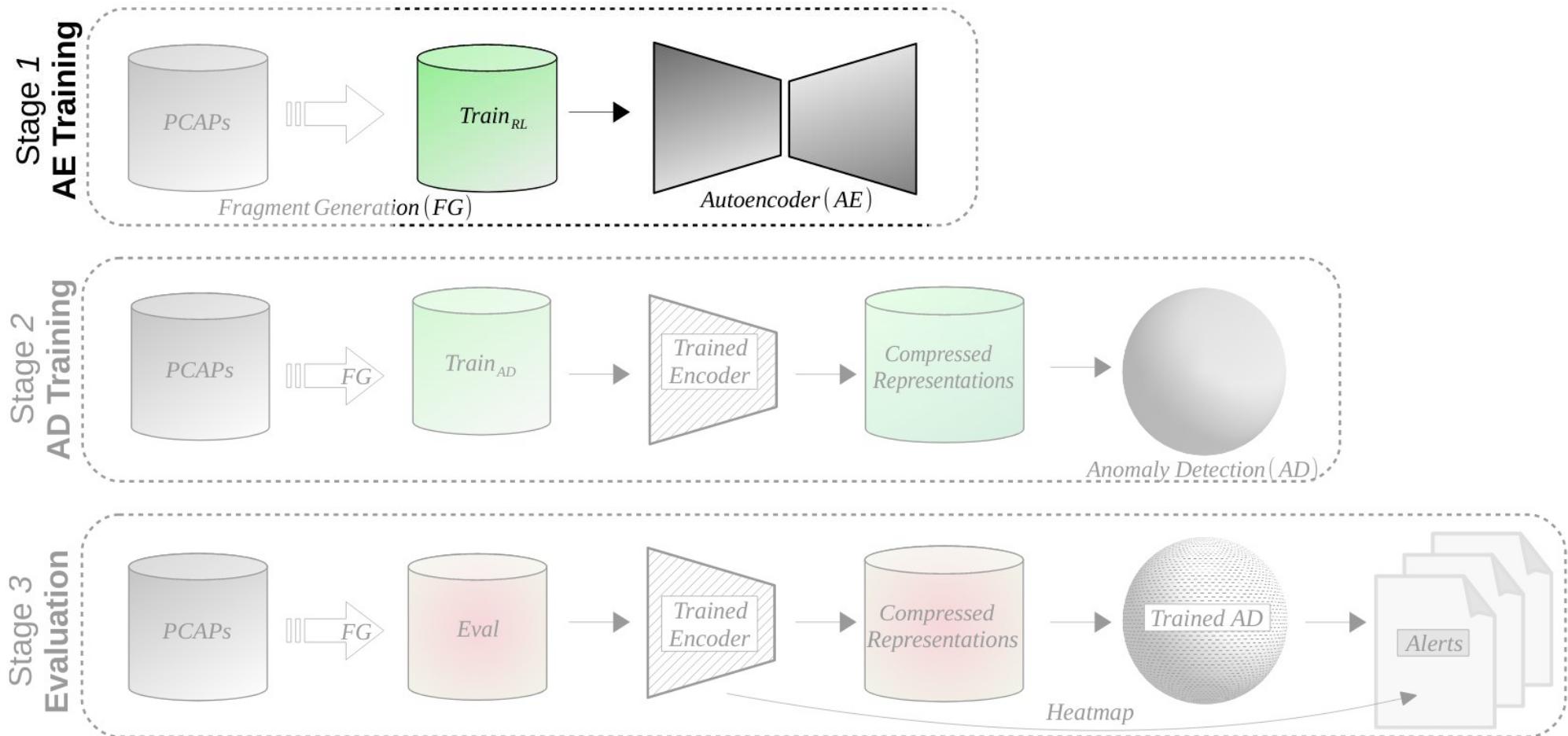


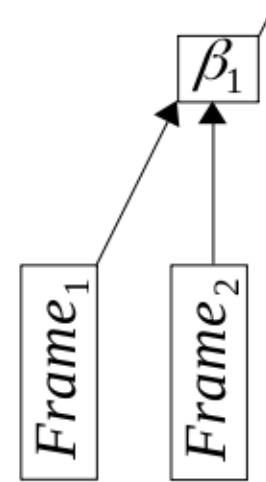
Figure 5.3: Four consecutive examples of flow-fragments shaded using a *viridis* color map

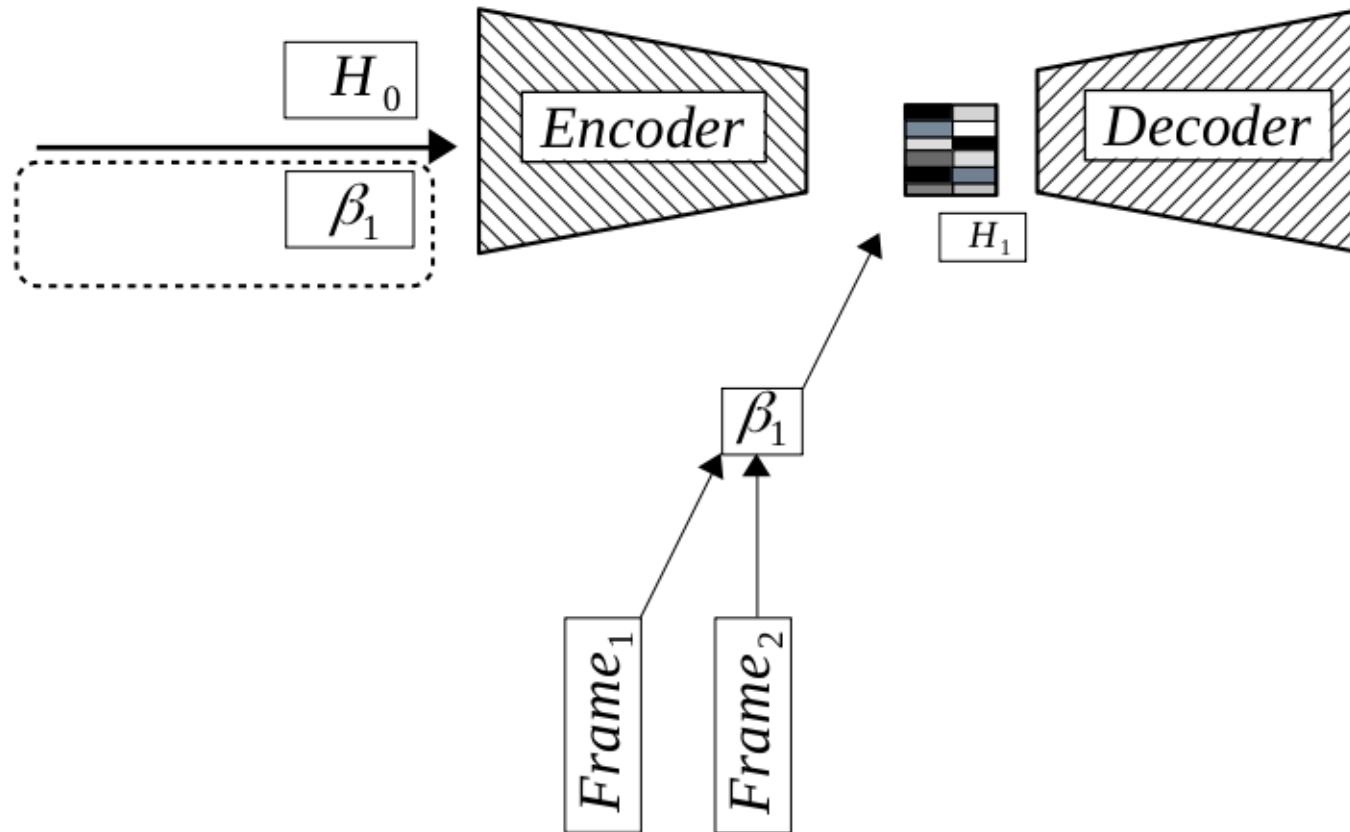
### 3 | Methodology: Overview

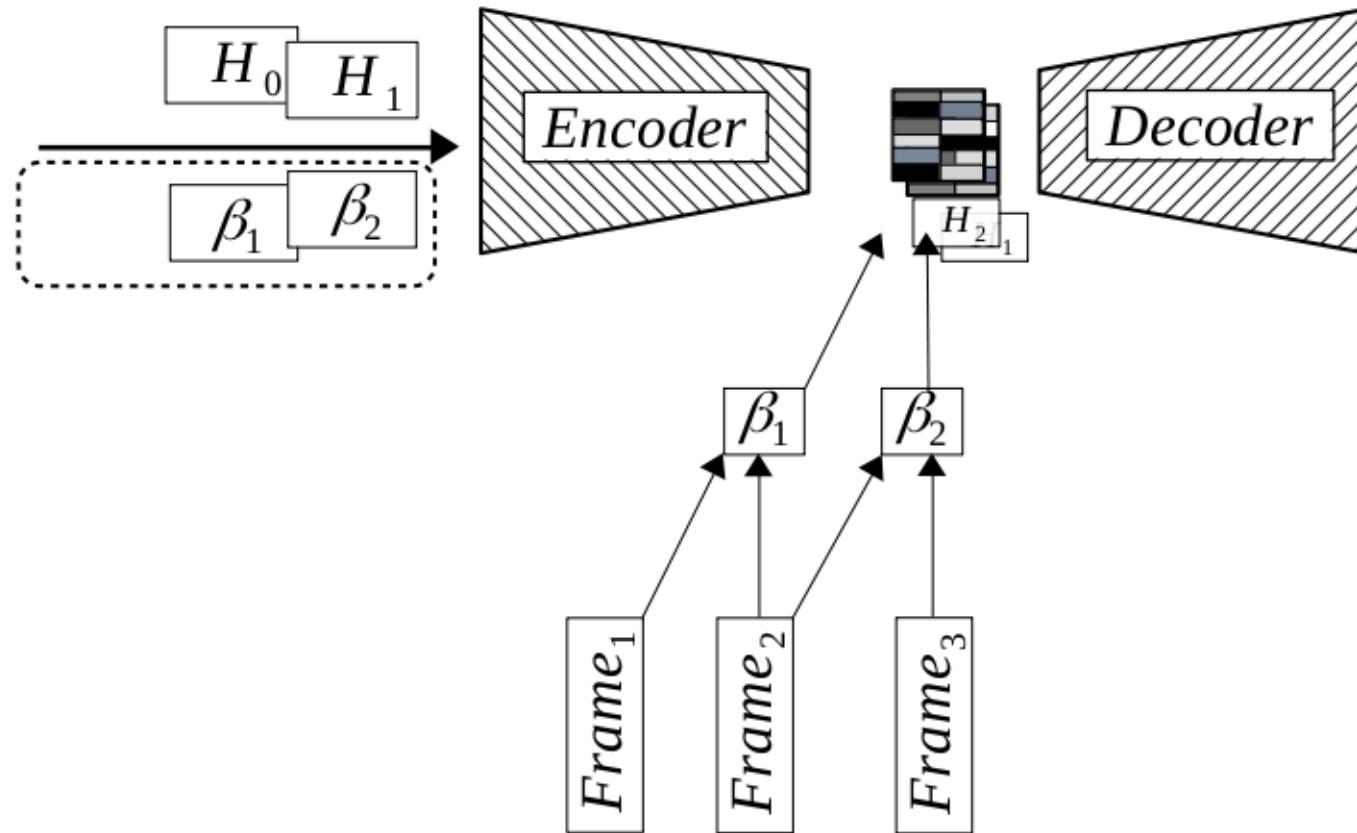


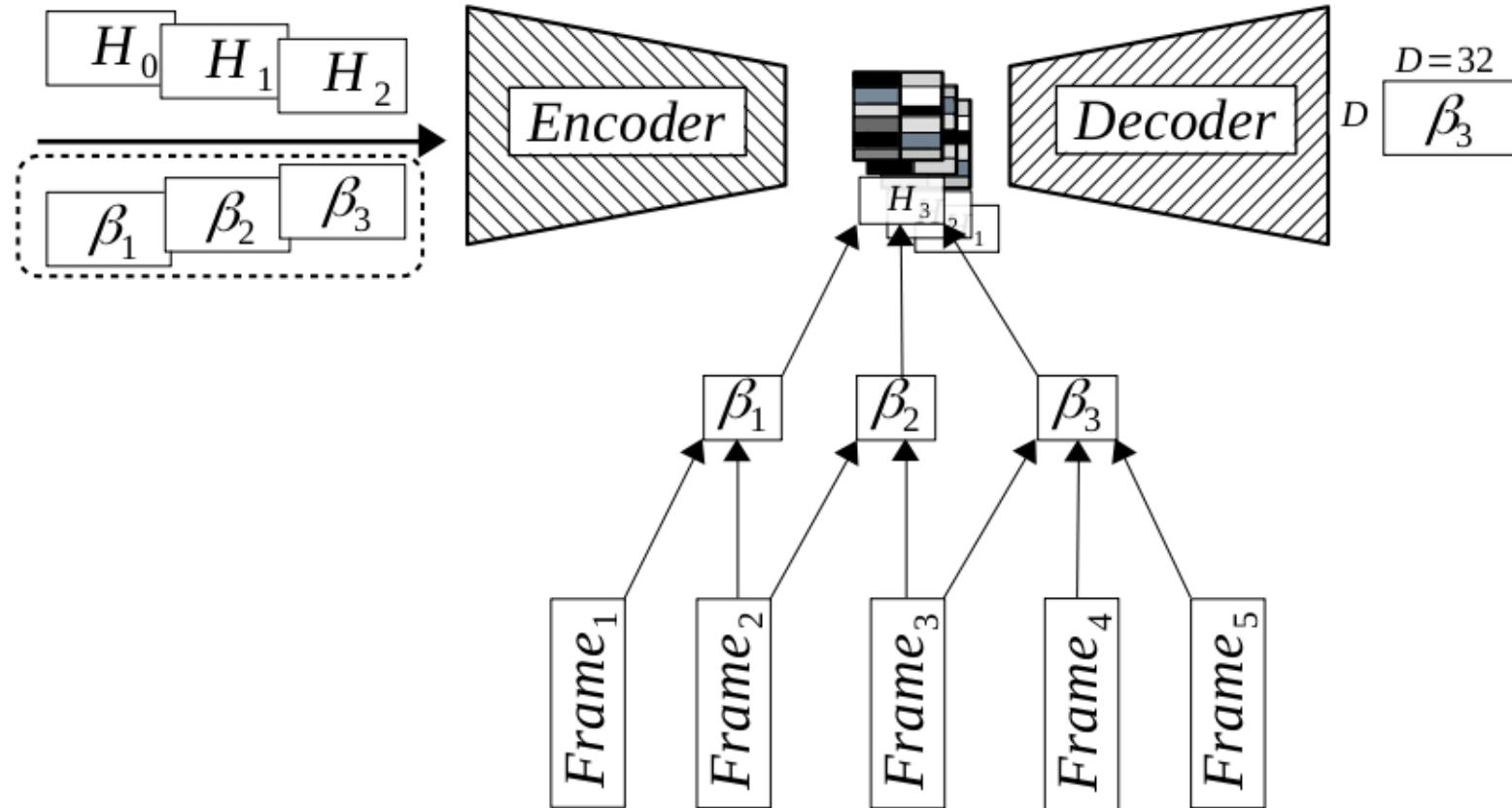
### 3 | Methodology: Overview

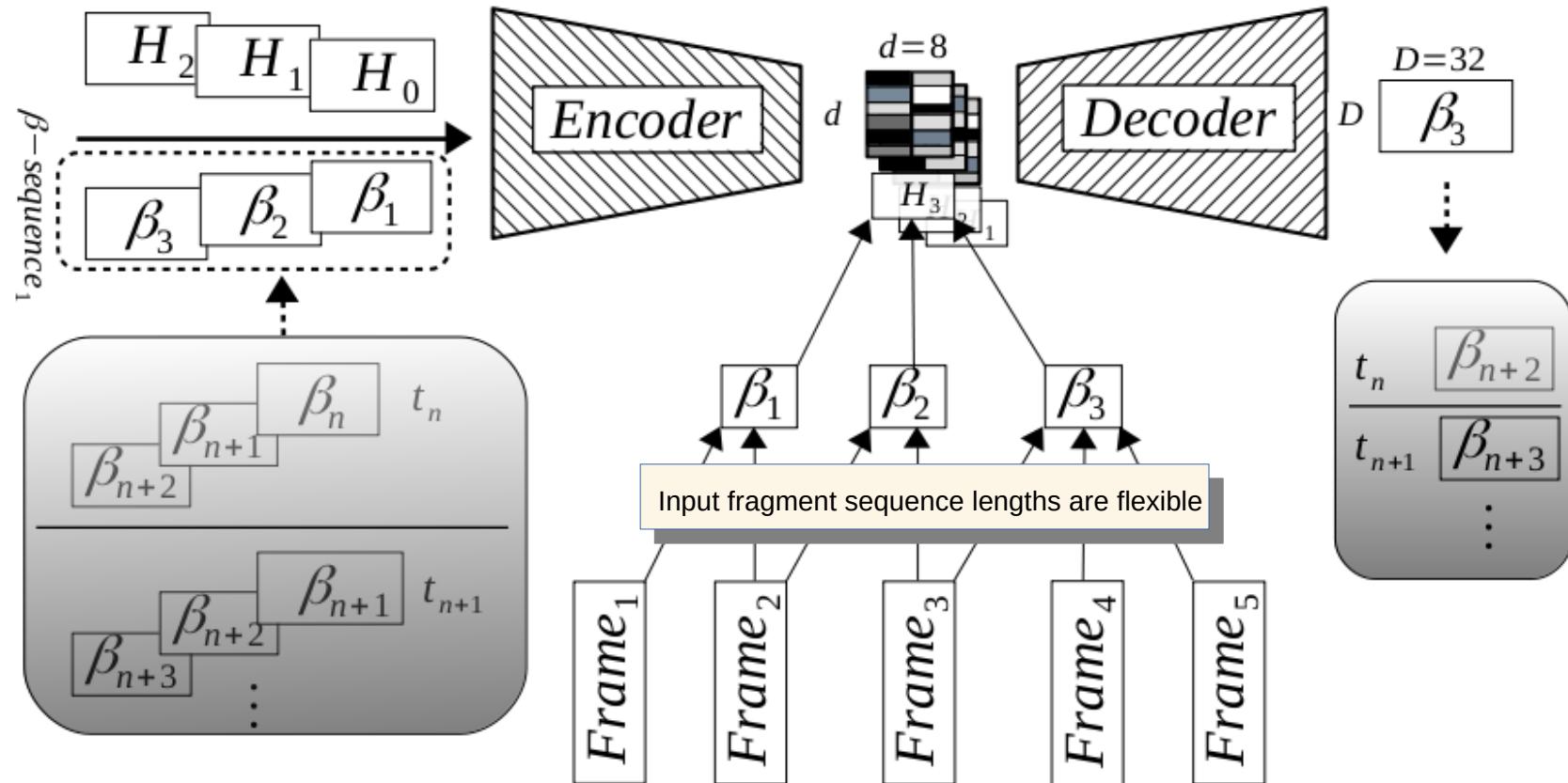








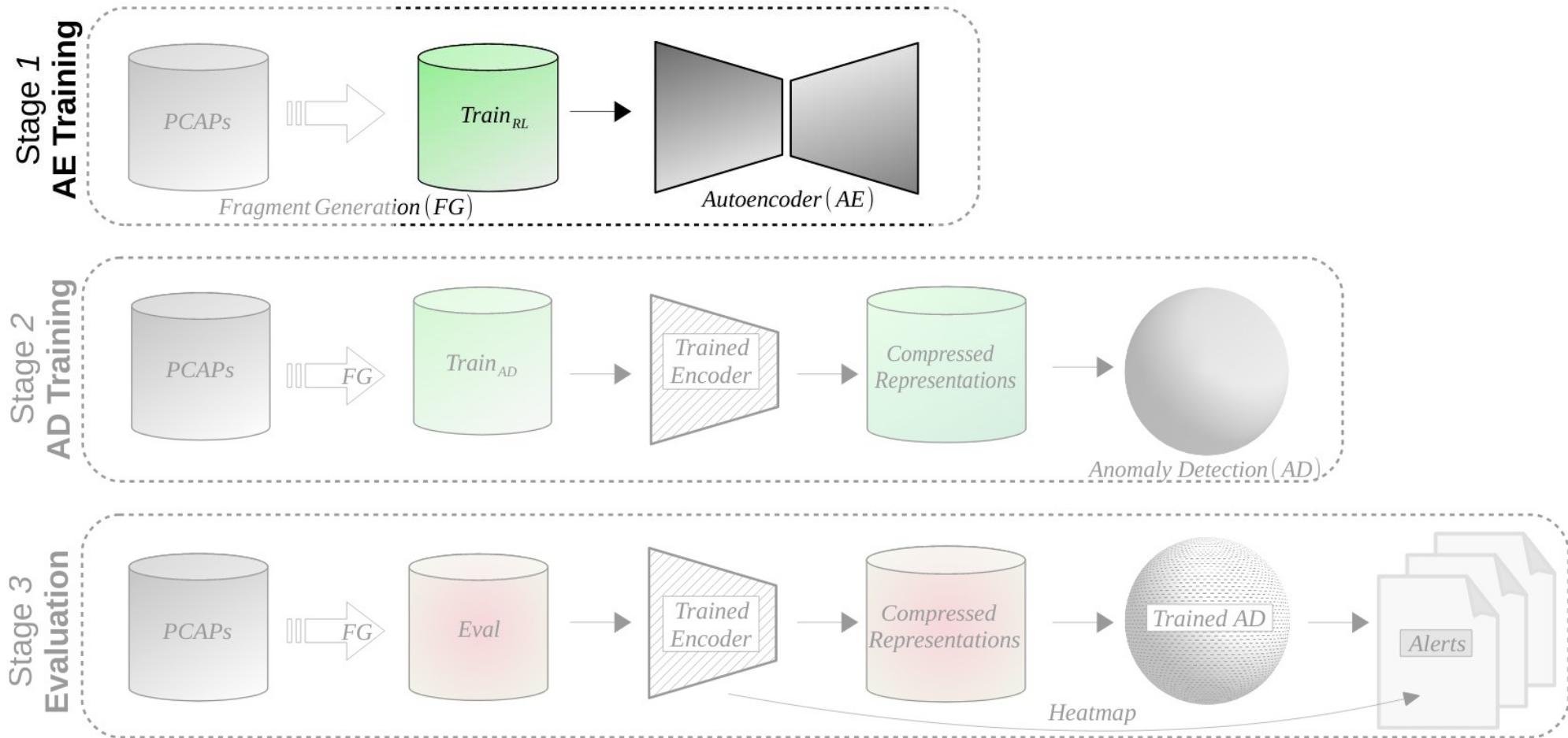




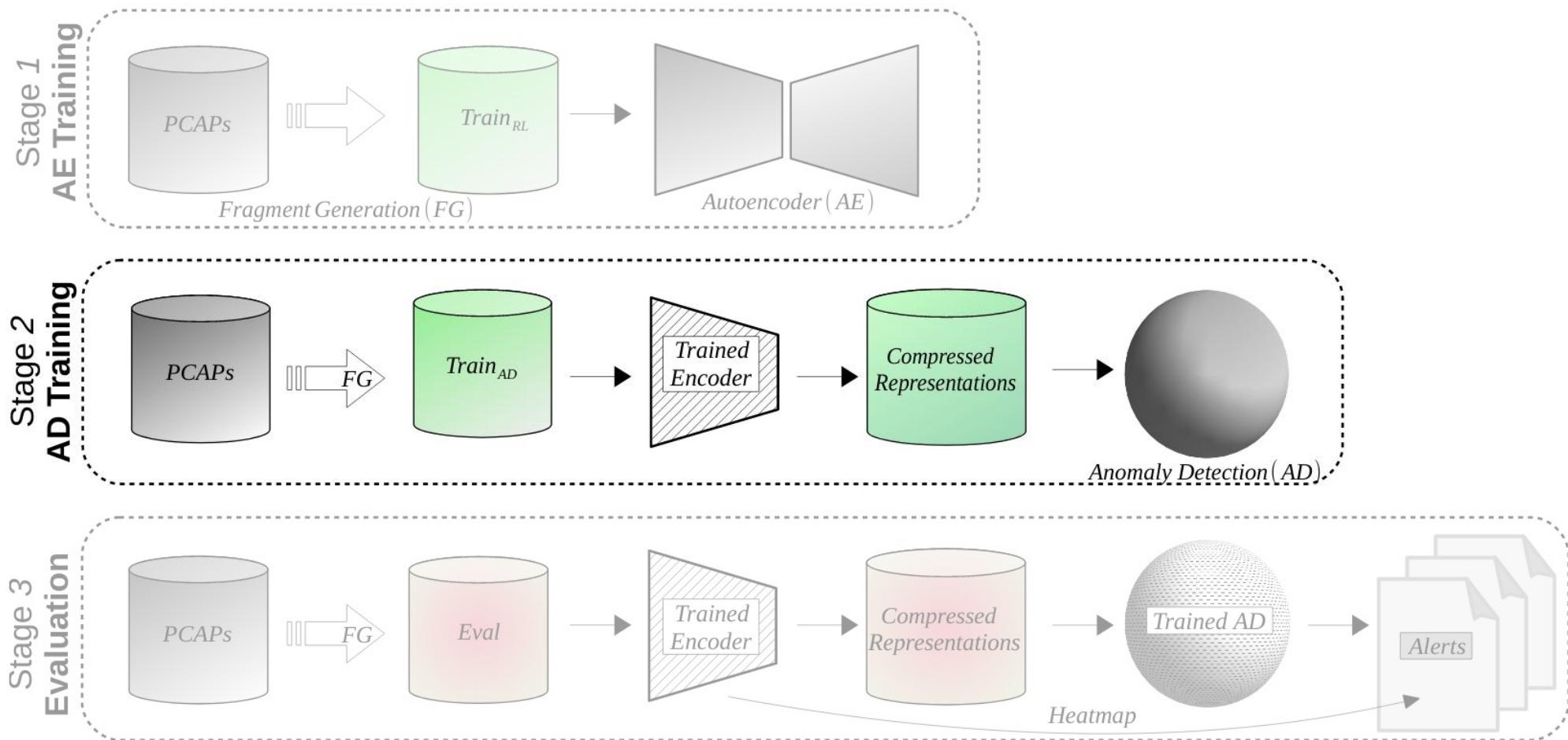
```
AutoEncoder(  
    (encoder): Encoder(  
        (stage1): Sequential(  
            (conv1_leaky_1): Conv2d(in=1, out=2, kernel_size=(3, 3),  
                stride=(1, 1),  
                padding=(1, 1))  
            (leaky_conv1_leaky_1): LeakyReLU(negative_slope=0.2, inplace=True)  
        )  
        (rnn1): CGRU_cell(  
            (dropout): Dropout2d(p=0.1, inplace=False)  
            (conv_gates): Sequential(  
                (0): Conv2d(in=6, out=8, kernel_size=(5, 5),  
                    stride=(1, 1),  
                    padding=(2, 2))  
                (1): GroupNorm(groups=8, channels=8, eps=1e-05, affine=True)  
            )  
            (conv_can): Sequential(  
                (0): Conv2d(in=6, out=4, kernel_size=(5, 5),  
                    stride=(1, 1),  
                    padding=(2, 2))  
                (1): GroupNorm(groups=4, channels=4, eps=1e-05, affine=True)  
            )  
        )  
    )
```

- Encoder compresses input data ( $32^2$  bytes) to **64 features**
- **convGRU** (X. Shi et al. – NIPS 2015) neurons for each layer
- Trained to **reconstruct the last input** of an input sequence
- **ReLU** (V. Nair et al. – ICML 2010) as an activation function
- **Mean square error (MSE)** as a loss function
- **Cycle LR-scheduler** (L. Smith et al. – SPIE 2019)
- Gradient descent with **small batch sizes** (J. Nocedal et al. – ICLR 2017)

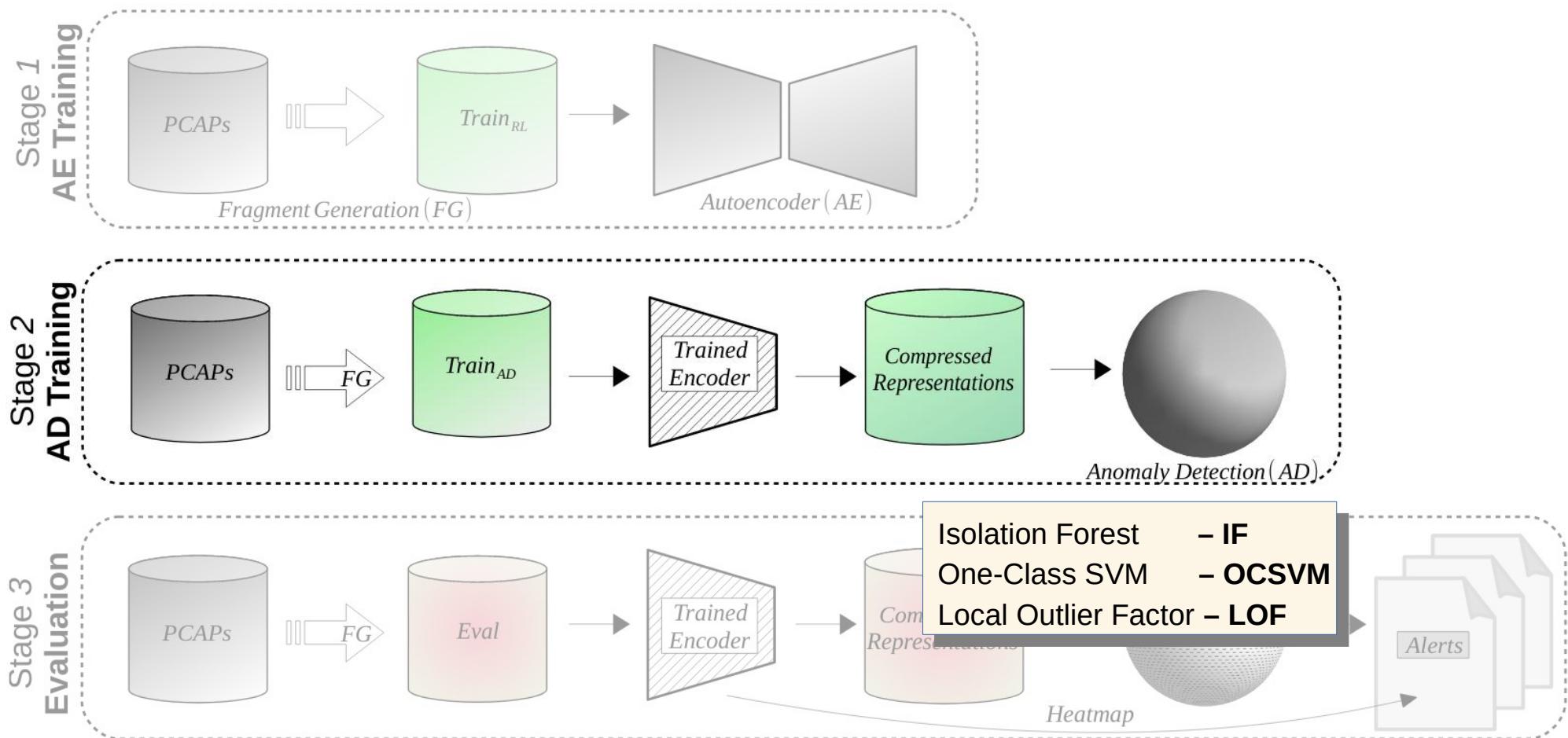
### 3 | Methodology: Overview

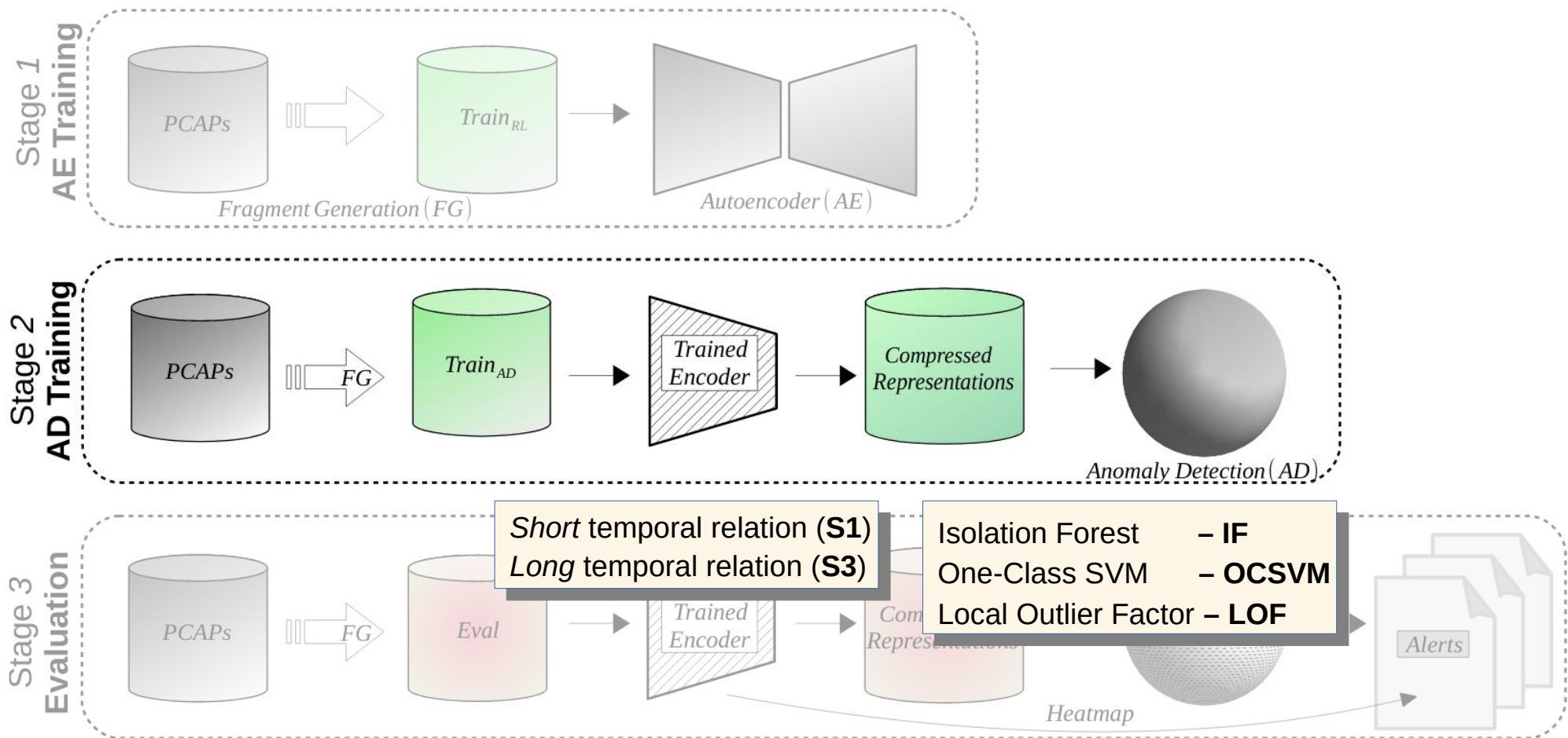


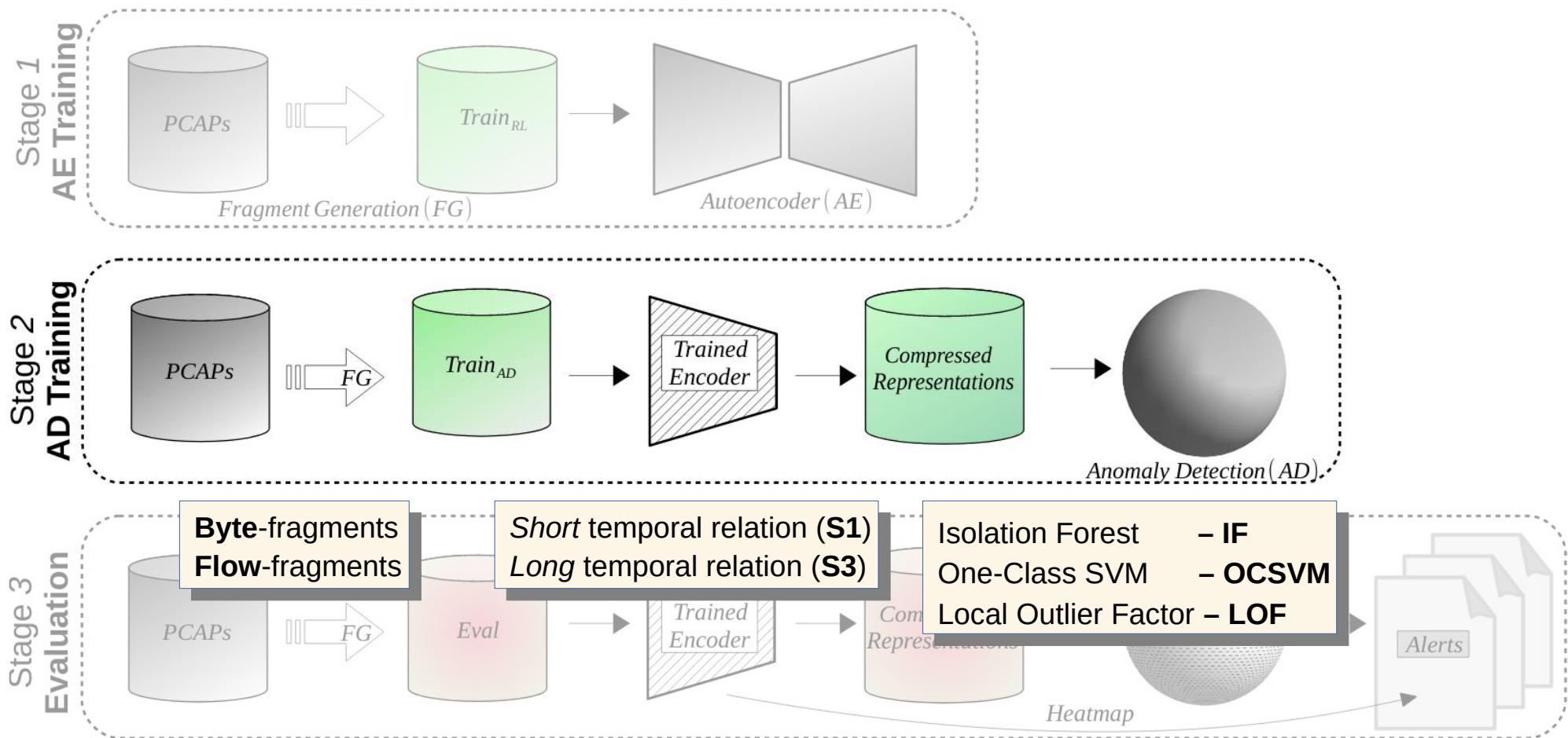
### 3 | Methodology: Overview



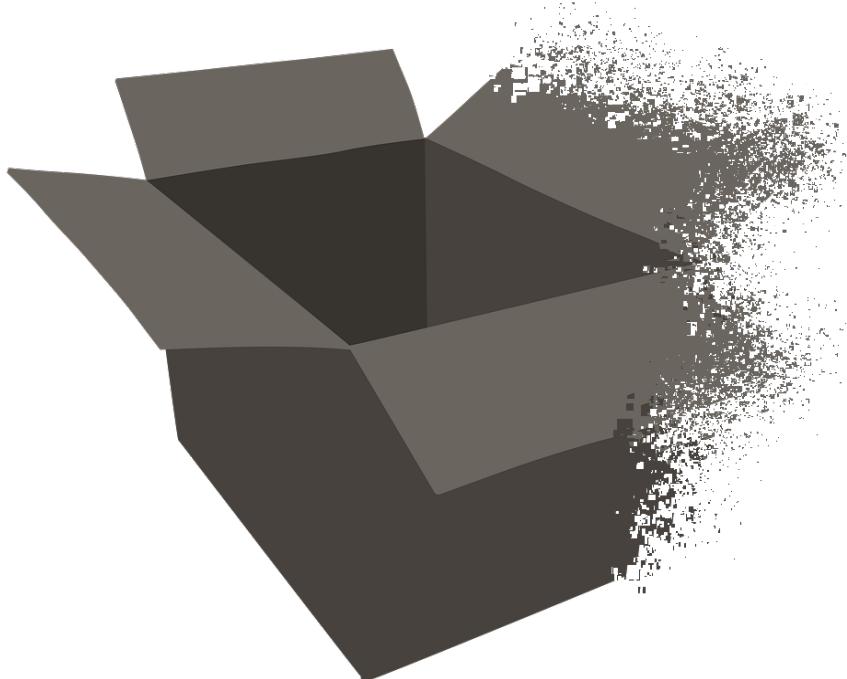
### 3 | Methodology: Overview



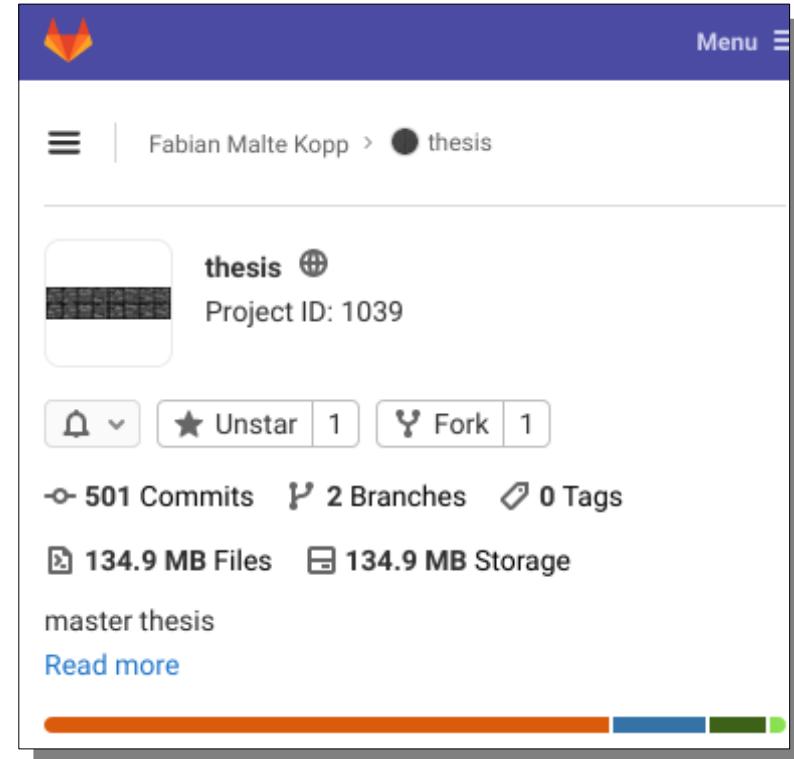




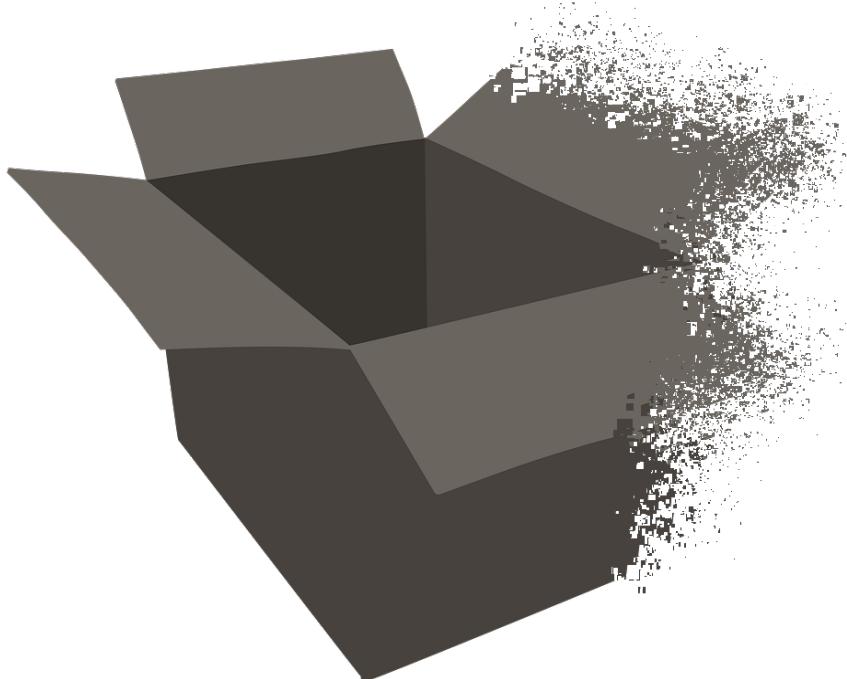
- 0 | Preamble
- 1 | Background
- 2 | Related Work
- 3 | Methodology
- 4 | Implementation**
- 5 | Evaluation
- 6 | Summary

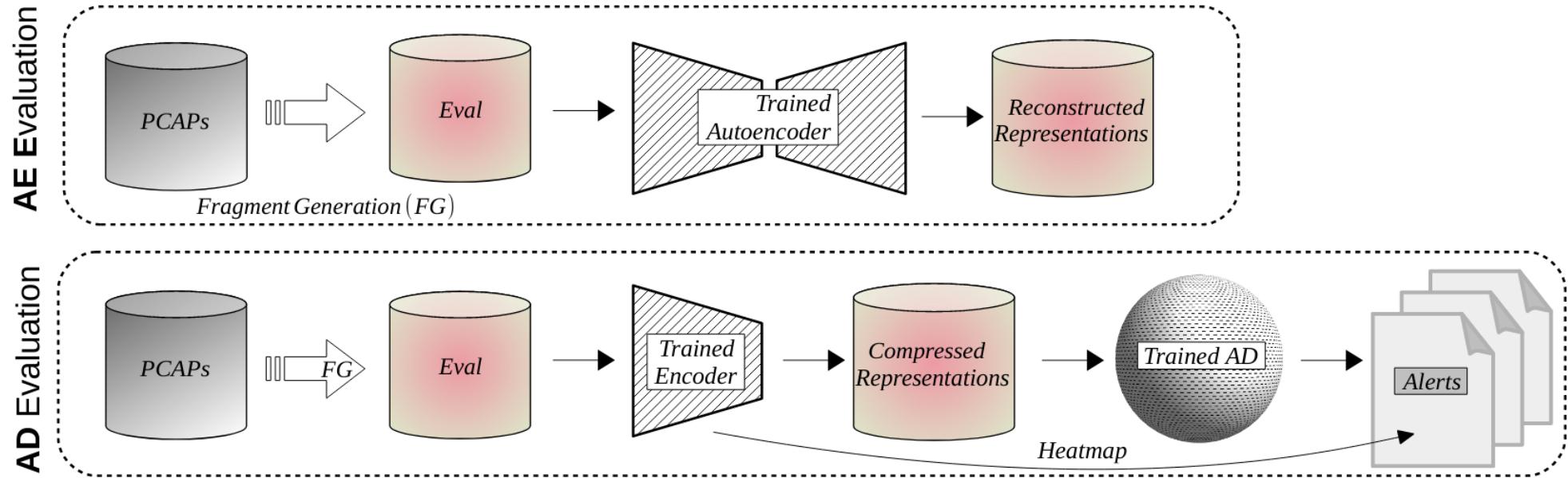


- Multi-threaded PCAP converter based on [dpkt](#)
  - Processes ~2,000 Ethernet frames/sec per computing core
  - Data is stored in [HDF5](#) containers (~10% overhead)
- [PyTorch](#)-based GPU-enabled model
  - Added integration for convGRU layers
- [Scikit-learn](#)-based anomaly detection integration
- Integrated [LRP implementation](#) for model interpretability
- Experiment tracking and logging via [Tensorboard](#)
  - [tensorboard.dev/experiment/OTD0qLBeQYK1DajlvkwoyQ/#graphs](https://tensorboard.dev/experiment/OTD0qLBeQYK1DajlvkwoyQ/#graphs)
- Developed and maintained via gitlab
  - [git.informatik.tu-cottbus.de/koppfab1/thesis](https://git.informatik.tu-cottbus.de/koppfab1/thesis)
  - Supplementary [Jupyter notebooks](#) for simple tutorials
  - ~3200 LOC / 500 commits



- 0 | Preamble
- 1 | Background
- 2 | Related Work
- 3 | Methodology
- 4 | Implementation
- 5 | Evaluation**
- 6 | Summary





- Industrial control system data
  - Power Plant traffic (BTU research data set, NYD)
  - Secure Water Treatment (*Tippenhauer – CySWater 2016*)

## 5 | Evaluation: Experiment Data

- Industrial control system data
  - Power Plant traffic (BTU research data set, NYD)
  - Secure Water Treatment (*Tippenhauer – CySWater 2016*)
- Ground truth information
  - Estimated ground truth
  - Synthetic data generation (*Hanusch – BTU 2018*)

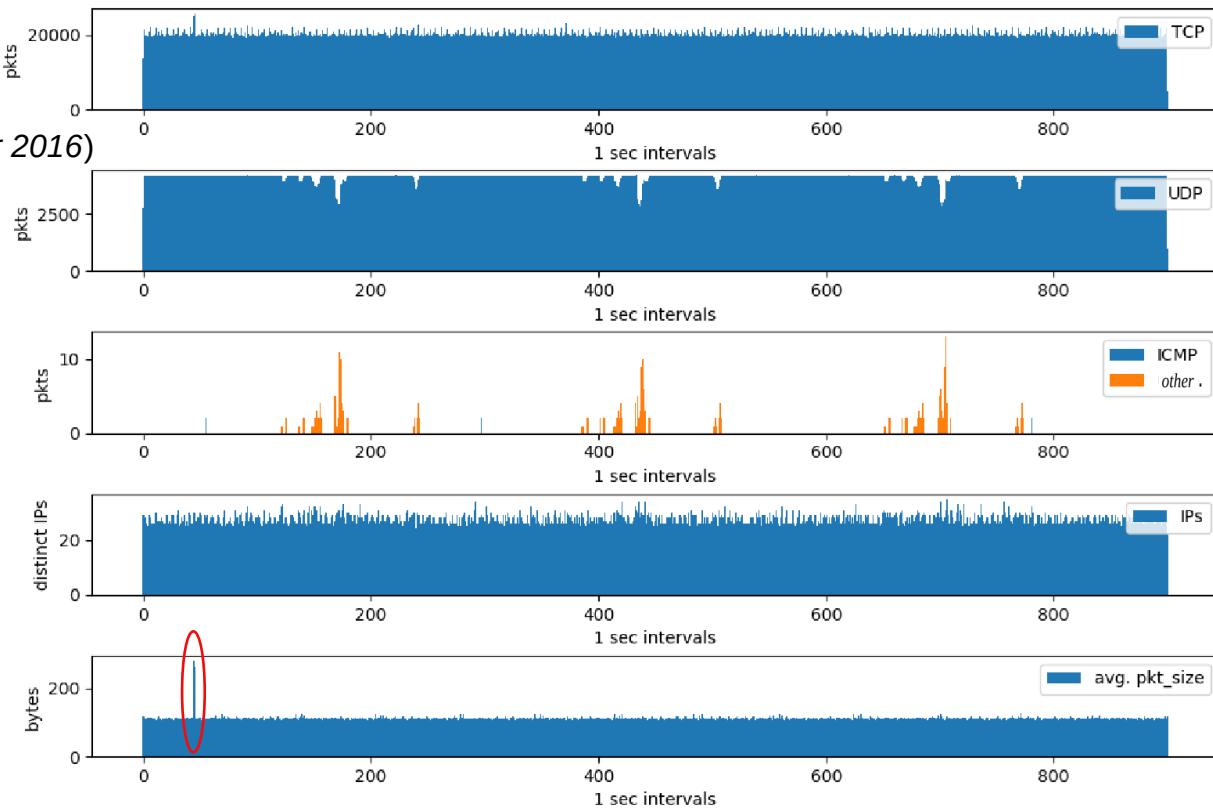


Figure 4.1: High-level traffic analysis of the SWaT data set

## 5 | Evaluation: Experiment Data

- Industrial control system data
  - Power Plant traffic (BTU research data set, NYD)
  - Secure Water Treatment (*Tippenhauer – CySWater 2016*)
- Ground truth information
  - Estimated ground truth
  - Synthetic data generation (*Hanusch – BTU 2018*)
- Data splitting
  - ~ 1.5 GB/3 min of training data (10M frames)
  - 10% for training validation

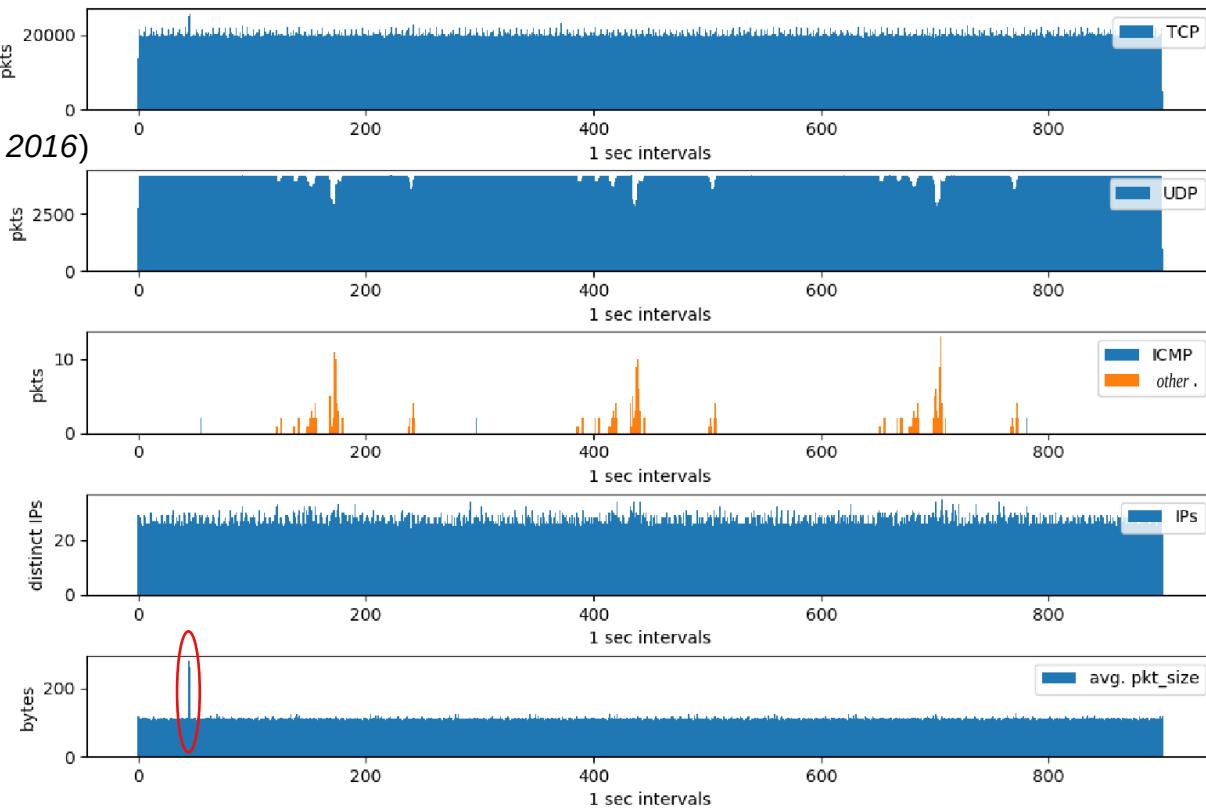


Figure 4.1: High-level traffic analysis of the SWaT data set

## 5 | Evaluation: Experiment Data

- Industrial control system data
  - Power Plant traffic (BTU research data set, NYD)
  - Secure Water Treatment (*Tippenhauer – CySWater 2016*)
- Ground truth information
  - Estimated ground truth
  - Synthetic data generation (*Hanusch – BTU 2018*)
- Data splitting
  - ~ 1.5 GB/3 min of training data (10M frames)
  - 10% for training validation
- Experiment trials
  - 2 input data orientations (*byte- & flow-frag.*)
  - 2 feature extraction models (*S1 & S3*)
  - 4 evaluation sets per data set (H, D, E, R)
  - 3+1 anomaly detection methods (IF, OCSVM, LOF, BL)

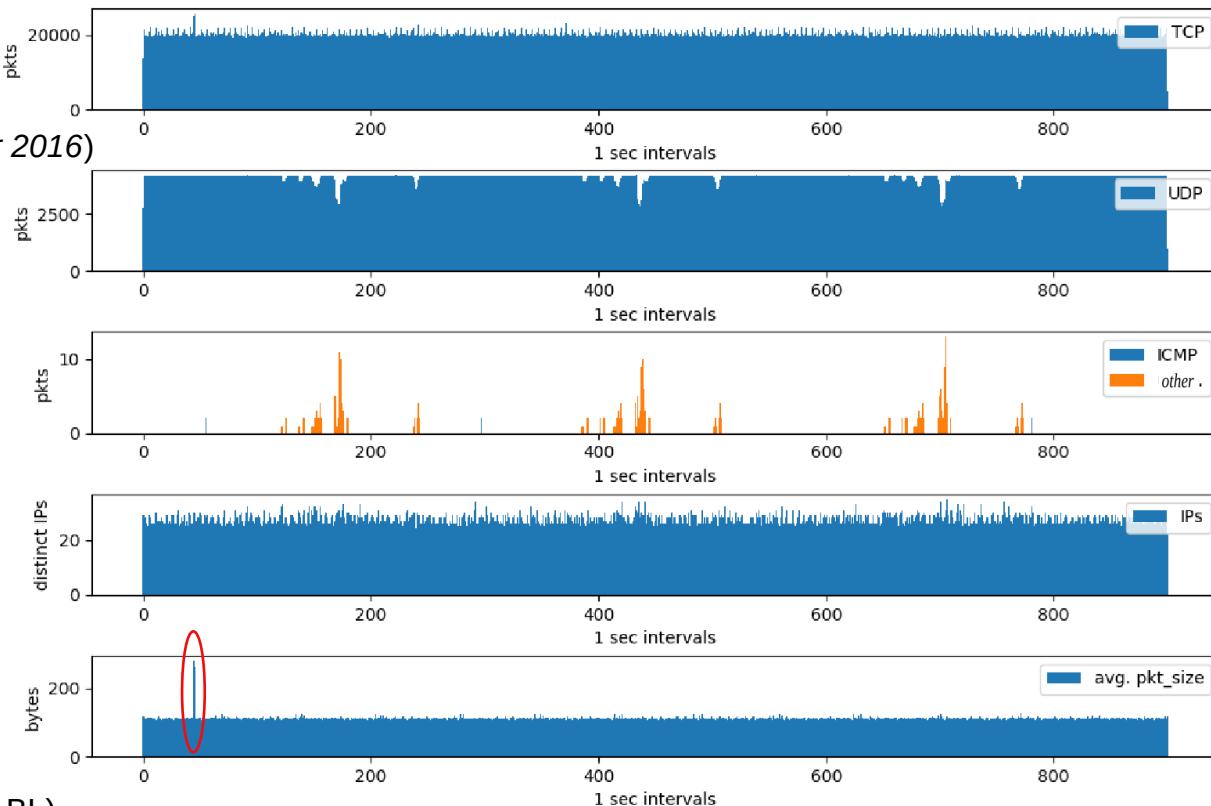
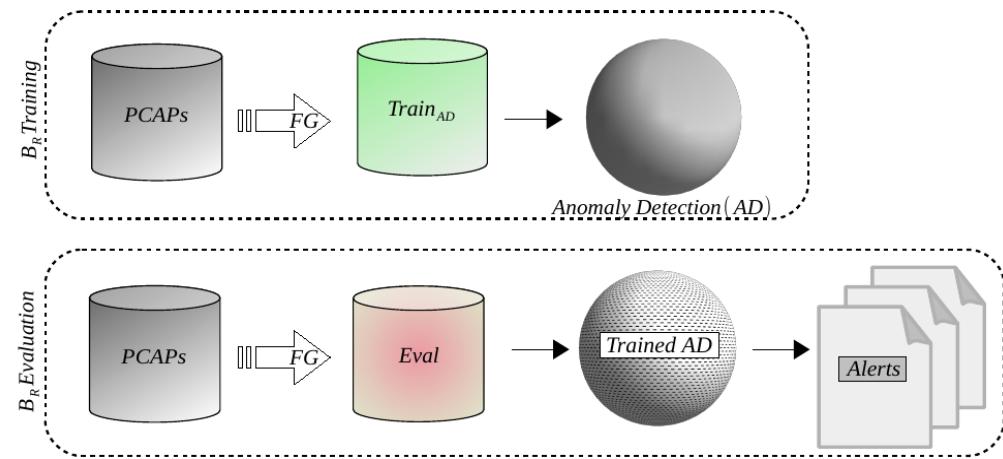
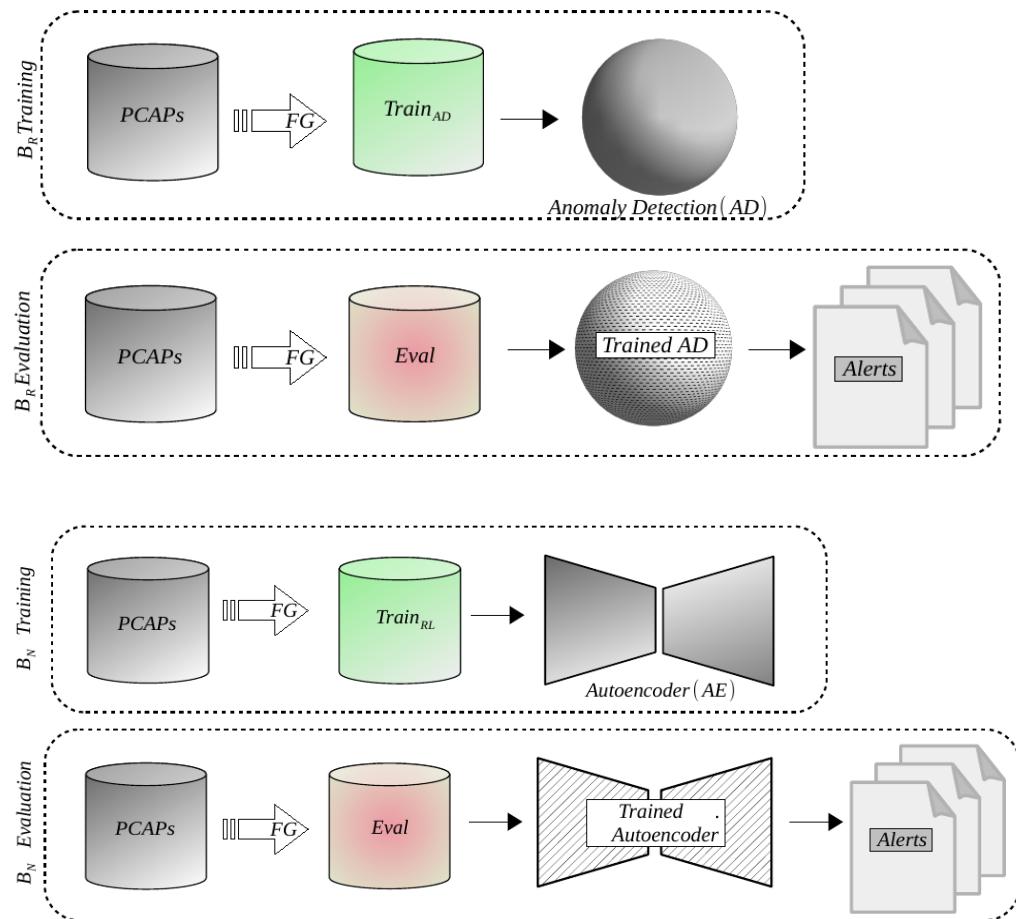


Figure 4.1: High-level traffic analysis of the SWaT data set

- Raw Baseline ( $B_r$ )
  - No additional feature extraction
  - Anomaly detection based on raw byte- or flow-frag.



- Raw Baseline ( $B_r$ )
  - No additional feature extraction
  - Anomaly detection based on raw byte- or flow-frag.
- Naive Baseline ( $B_n$ )
  - Feature extraction as an end-to-end model
  - Anomaly detection based on the model's loss values
  - Upper bound for normal behavior based on the average loss ( $AML$ ) plus a multiple ( $\nu$ ) of its standard deviation ( $\sigma$ )
    - $thr := AML + \nu\sigma$
    - $MSE(x, \hat{x}) > thr ? x \in A_{pred} : x \in N_{pred}$
    - For  $\nu = 2.5 \Rightarrow 98,9\%$  of normal examples are classified correctly (normal distribution)



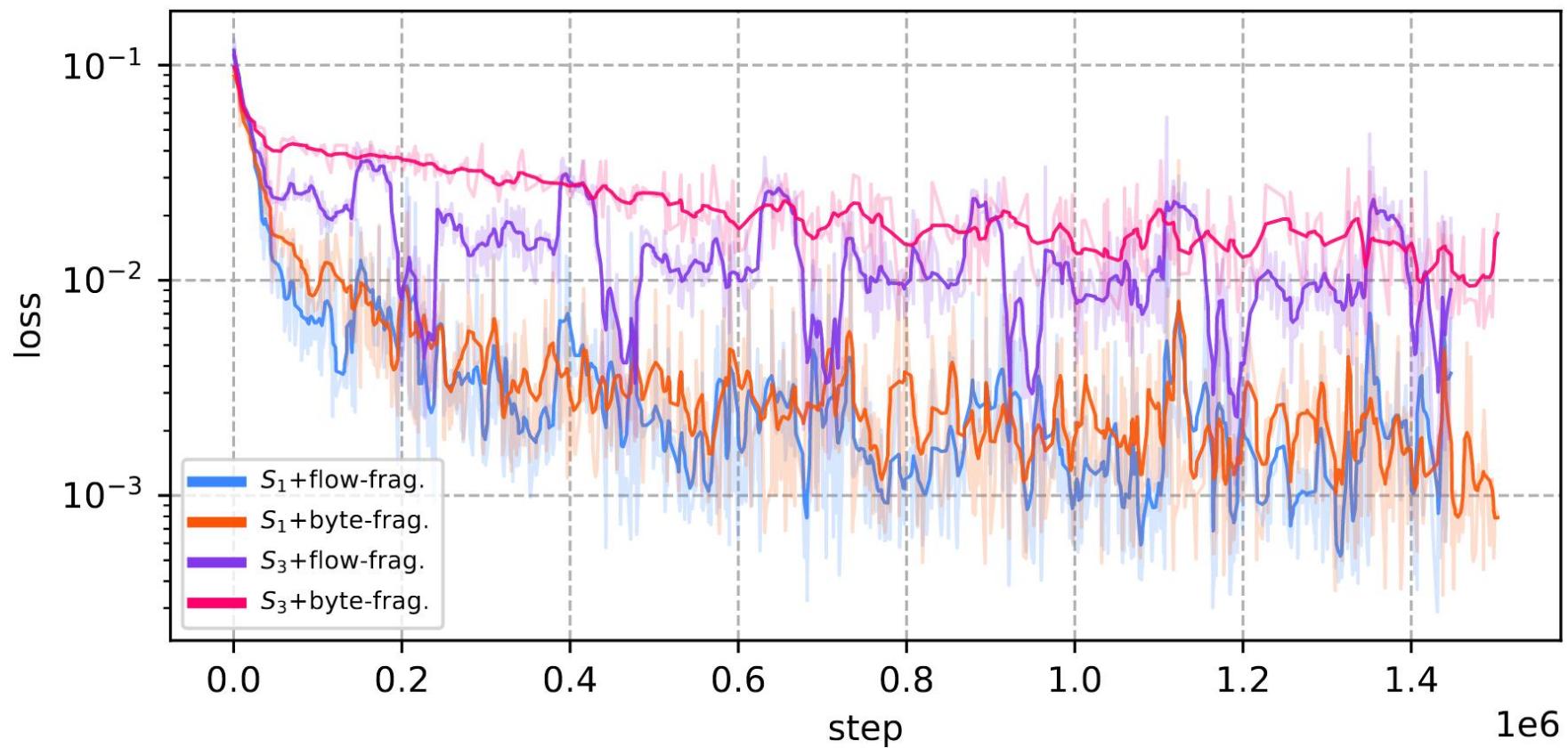


Figure 6.2: Training loss-curves for the 4 different trained models for the SWaT A6 data set

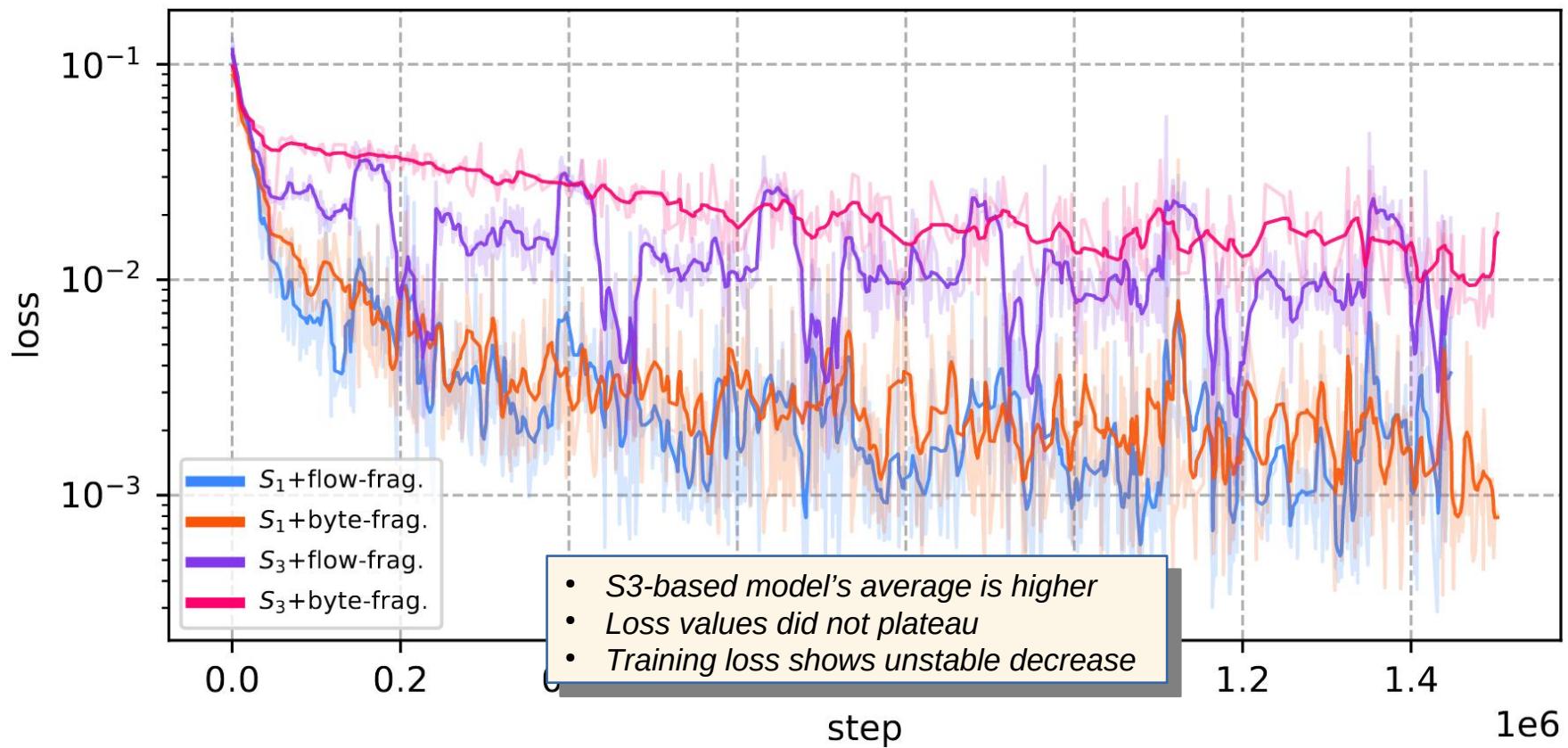


Figure 6.2: Training loss-curves for the 4 different trained models for the SWaT A6 data set

## 5 | Evaluation: (Flow-)Fragment Reconstruction

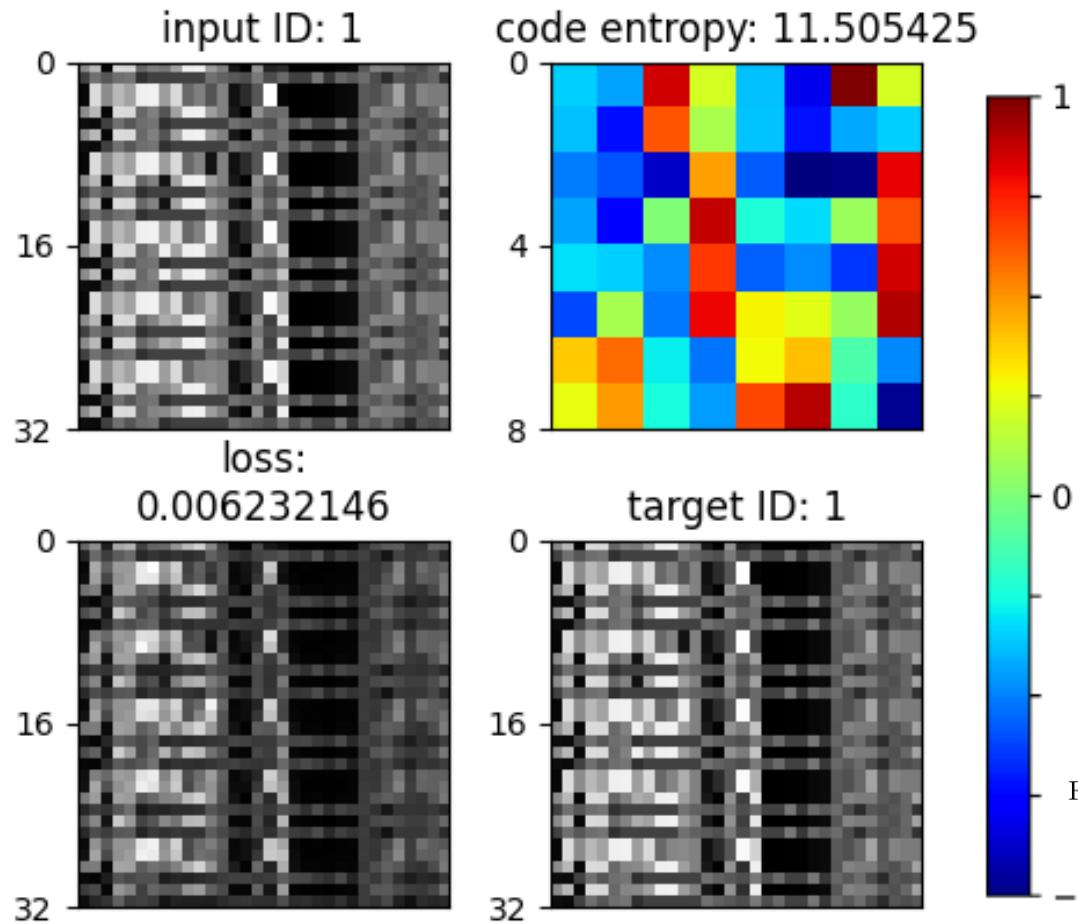


Figure 6.6: Input fragment (top left), extracted features (top right),  
reconstructed fragment (lower left), target fragment (lower right)

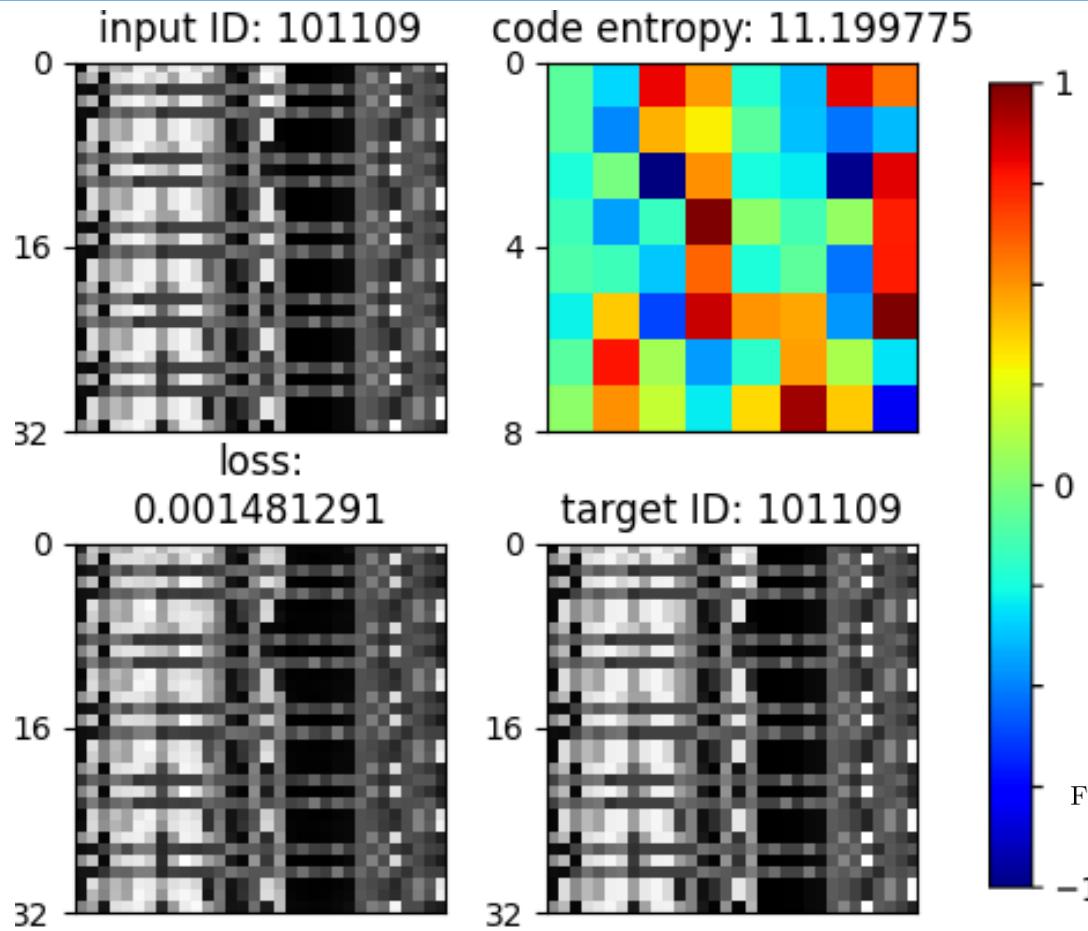
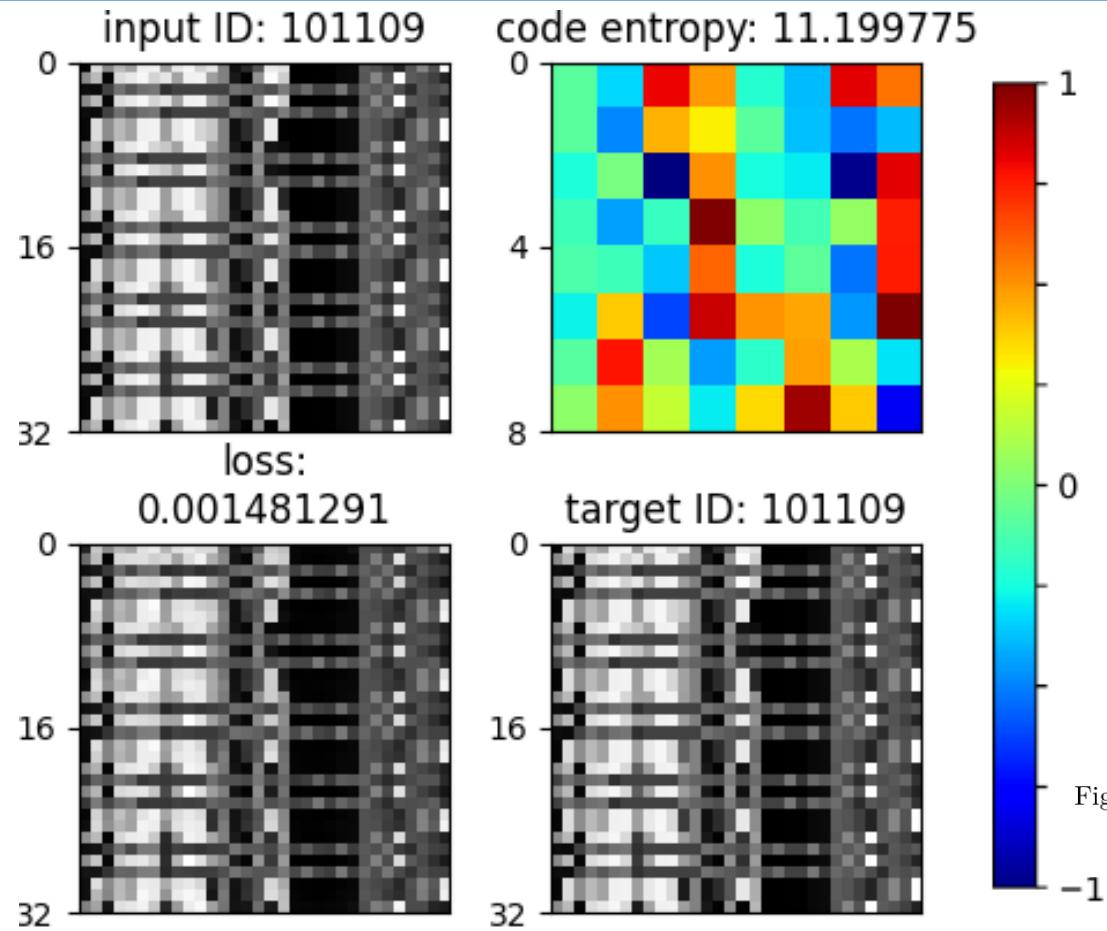


Figure 6.6: Input fragment (top left), extracted features (top right), reconstructed fragment (lower left), target fragment (lower right)

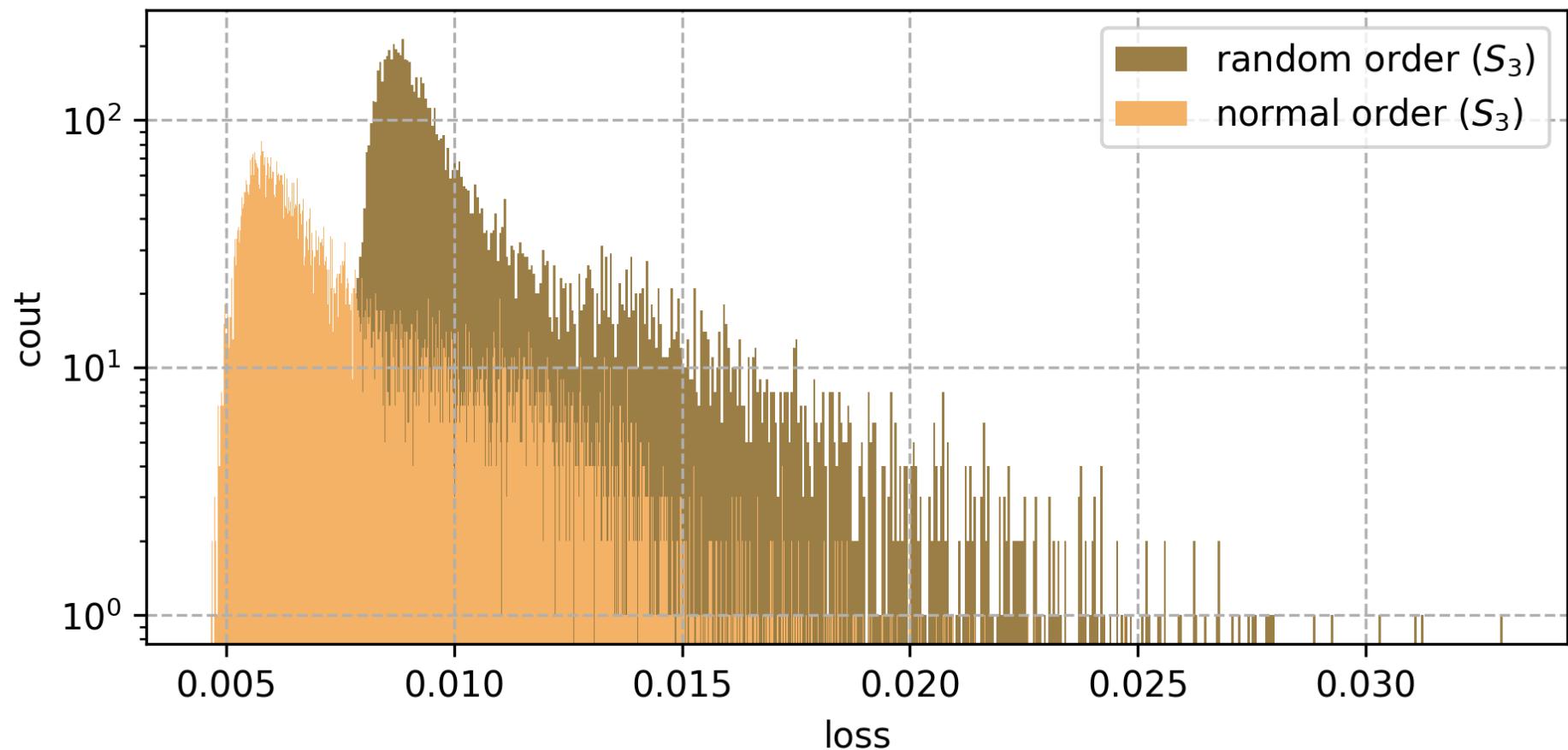
## 5 | Evaluation: (Flow-)Fragment Reconstruction

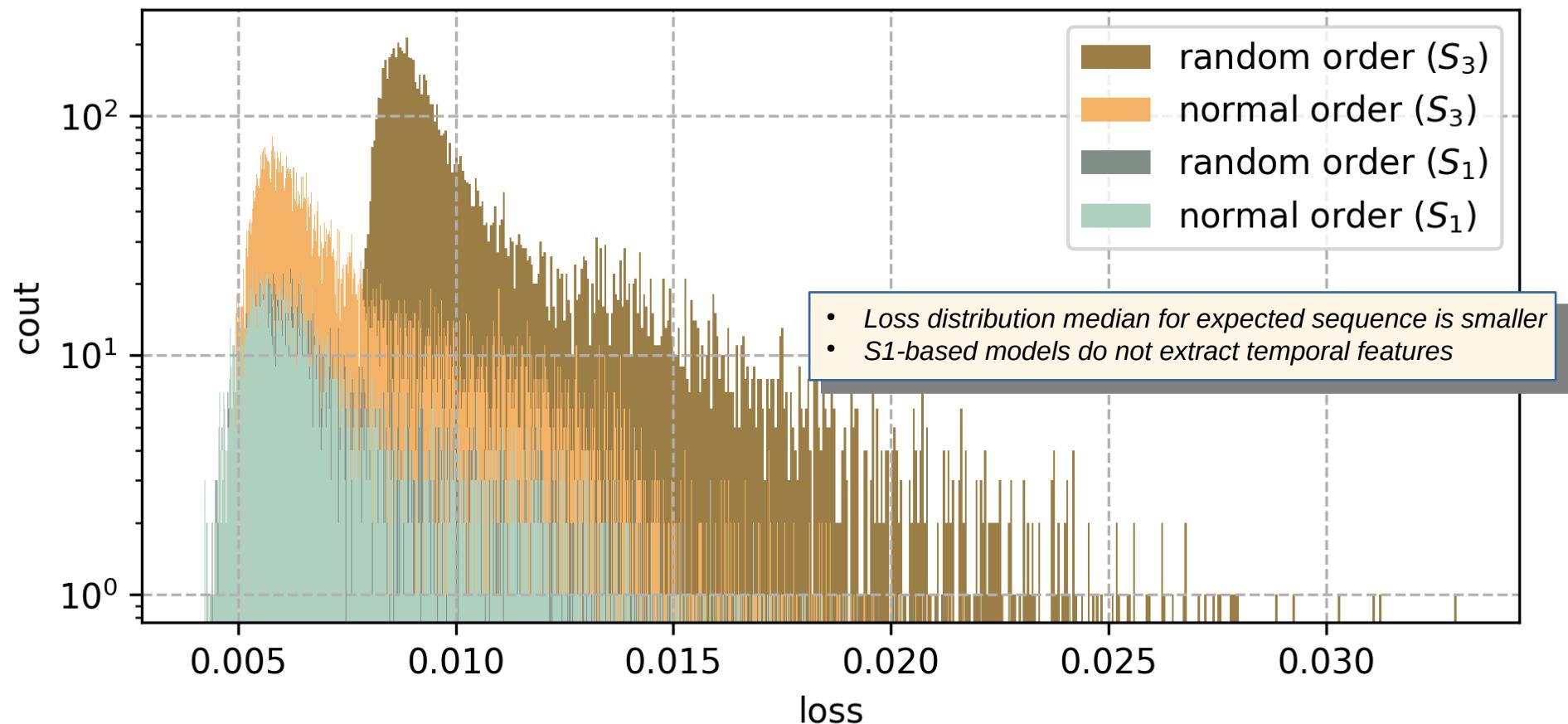


- Spatial features are reconstructed accurately
- Extracted features are structured & not intuitive

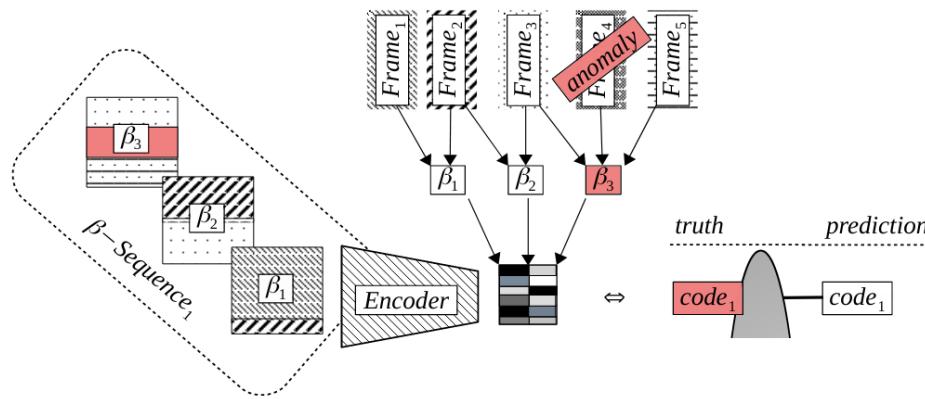
Figure 6.6: Input fragment (top left), extracted features (top right),  
reconstructed fragment (lower left), target fragment (lower right)

## 5 | Evaluation: Training Results (III)

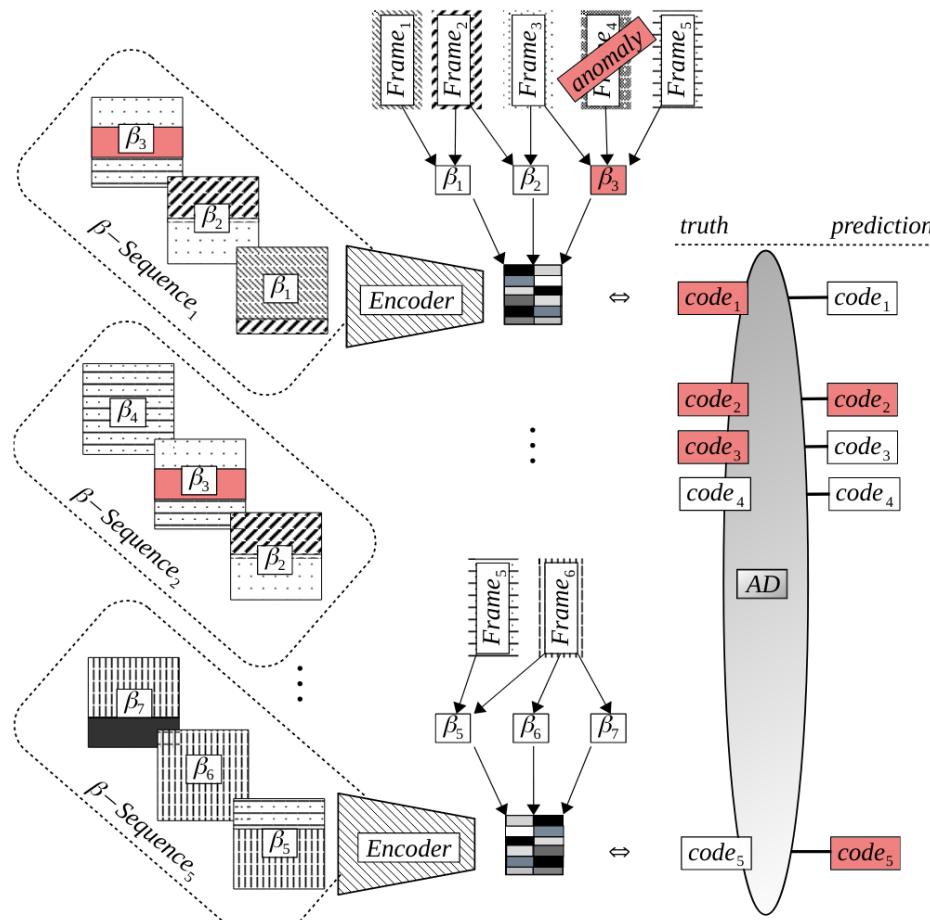




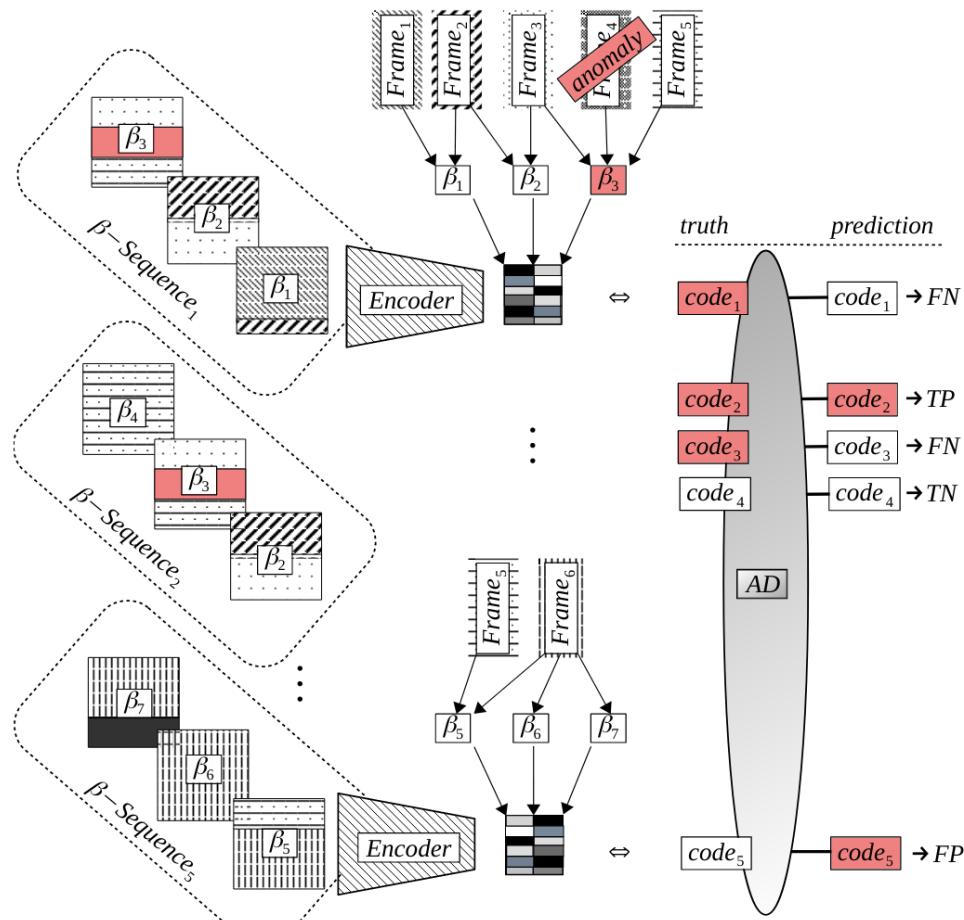
## 5 | Evaluation: Ground Truth Measurement



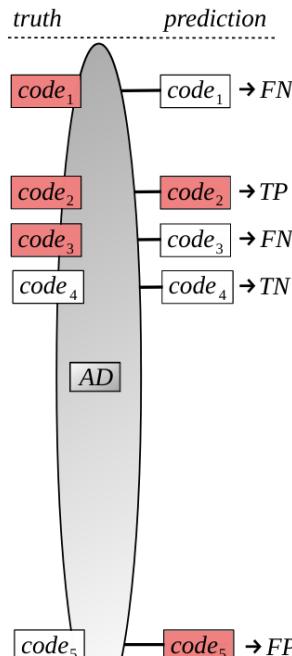
## 5 | Evaluation: Ground Truth Measurement



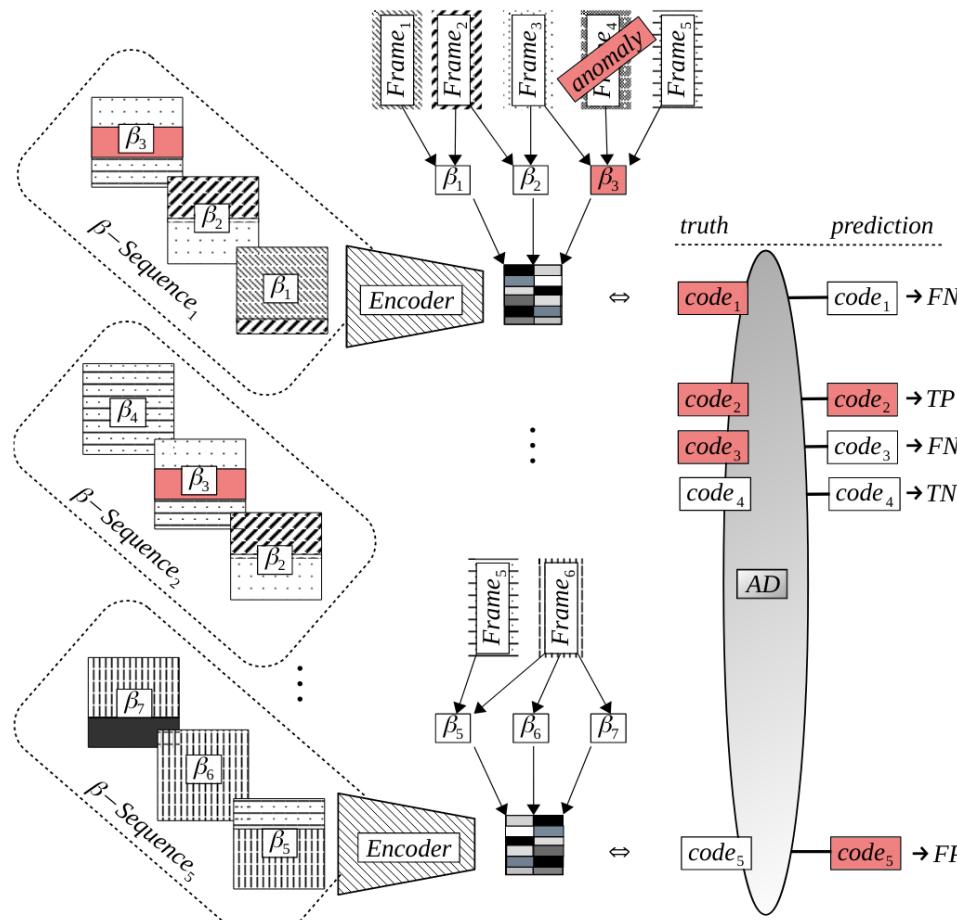
## 5 | Evaluation: Ground Truth Measurement



False Negatives – <b>FN</b>
True Positives – <b>TP</b>
True Negatives – <b>TN</b>
False Positives – <b>FP</b>



## 5 | Evaluation: Ground Truth Measurement



False Negatives – **FN**

True Positives – **TP**

True Negatives – **TN**

False Positives – **FP**

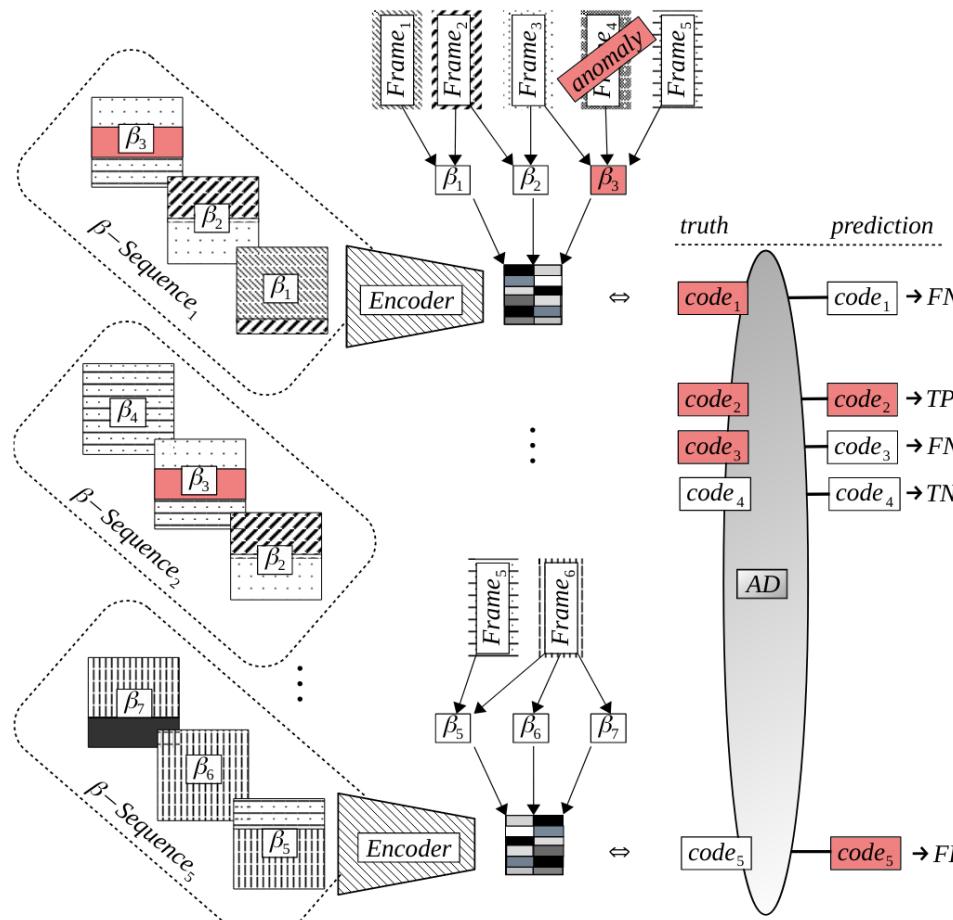
Recall – **RC**

$$RC = \frac{|TP|}{|TP|+|FN|}$$

Precision – **PR**

$$PR = \frac{|TP|}{|TP|+|FP|}$$

## 5 | Evaluation: Ground Truth Measurement



False Negatives – **FN**

True Positives – **TP**

True Negatives – **TN**

False Positives – **FP**

Recall – **RC**

$$RC = \frac{|TP|}{|TP|+|FN|}$$

Precision – **PR**

$$PR = \frac{|TP|}{|TP|+|FP|}$$

**F1-Score**

$$F_1 = 2 \times \left( \frac{|PR| * |RC|}{|PR| + |RC|} \right)$$

False Pos. Rate – **FPR**

$$FPR = \frac{|FP|}{|FP| + |TN|}$$

Fragment	$S_n$	Split	Time	#Samples	#Anomalies
Byte			SWaT		
$S_1$	<i>Train</i>	01:17:12	1,107,198	0	
	<i>Validation</i>	00:00:44	10,990	0	
	<i>Eval_H</i>	01:06:35	821,442	11,208	
	<i>Eval_E</i>	00:42:25	684,501	288	
	<i>Eval_D</i>	00:40:44	697,304	25,883	
	<i>Eval_R</i>	00:49:32	684,426	101	
$S_3$	<i>Train</i>	03:06:50	1,107,196	0	
	<i>Validation</i>	00:01:47	10,988	0	
	<i>Eval_H</i>	02:12:25	821,440	11,210	
	<i>Eval_E</i>	01:35:33	684,499	578	
	<i>Eval_D</i>	01:32:32	697,302	35,642	
	<i>Eval_R</i>	00:49:32	684,424	232	

Fragment	$S_n$	Split	Time	#Samples	#Anomalies
Flow			SWaT		
$S_1$	<i>Train</i>	00:27:18	482,280	0	
	<i>Validation</i>	00:00:16	4624	0	
	<i>Eval_H</i>	00:22:09	346,786	1,123	
	<i>Eval_E</i>	00:20:37	291,373	29	
	<i>Eval_D</i>	00:18:03	295,128	3,776	
	<i>Eval_R</i>	00:18:45	291,443	14	
$S_3$	<i>Train</i>	00:28:08	482,278	0	
	<i>Validation</i>	00:00:16	4,622	0	
	<i>Eval_H</i>	00:21:33	346,784	1,349	
	<i>Eval_E</i>	00:45:10	291,371	55	
	<i>Eval_D</i>	00:38:51	295,126	3,808	
	<i>Eval_R</i>	00:17:32	291,441	42	

- Fragment generation condense the total number of samples
- Fewer Flow-Fragments were extracted
- Feature extraction took longer on for models trained on longer sequences

## 5 | Evaluation: Test Results (I)

---

AD  
IF

OCSVM

LOF

$B_n$

—

## 5 | Evaluation: Test Results (I)

AD Fragments  
IF Byte

Flow

OCSVM  
Byte

Flow

LOF  
Byte

Flow

$B_n$   
Byte

Flow

## 5 | Evaluation: Test Results (I)

AD	Fragments	Model
IF		
Byte	$S_1$ $S_3$ $B_r$	
Flow	$S_1$ $S_3$ $B_r$	
OCSVM		
Byte	$S_1$ $S_3$ $B_r$	
Flow	$S_1$ $S_3$ $B_r$	
LOF		
Byte	$S_1$ $S_3$ $B_r$	
Flow	$S_1$ $S_3$ $B_r$	
$B_n$		
Byte	$S_1$ $S_3$	
Flow	$S_1$ $S_3$	

## 5 | Evaluation: Test Results (I)

AD	Fragments	Model	T2F↓	FPR↓
IF				
Byte				
	$S_1$	00:07:31	0	
	$S_3$	00:05:58	0	
	$B_r$	00:45:14	0.0001	
Flow				
	$S_1$	00:05:15	0	
	$S_3$	00:05:11	0	
	$B_r$	00:36:12	0.0089	
<hr/> OCSVM				
Byte				
	$S_1$	00:05:52	0	
	$S_3$	00:04:32	0.7869	
	$B_r$	19:51:57	0.5220	
Flow				
	$S_1$	00:01:02	0	
	$S_3$	00:01:07	0	
	$B_r$	10:13:46	0.5220	
<hr/> LOF				
Byte				
	$S_1$	00:01:19	0	
	$S_3$	00:05:50	0	
	$B_r$	40:13:40	0	
Flow				
	$S_1$	00:01:12	0	
	$S_3$	01:24:37	0	
	$B_r$	06:00:33	0	
<hr/> $B_n$				
Byte				
	$S_1$	05:31:37	0.0035	
	$S_3$	11:07:54	0.0234	
Flow				
	$S_1$	02:19:16	0.0702	
	$S_3$	04:41:41	0.0014	

Table B.3: SWaT train validation results

## 5 | Evaluation: Test Results (I)

AD	Fragments	Model	T2F↓	FPR↓
IF				
Byte				
	$S_1$	00:07:31	0	
	$S_3$	00:05:58	0	
	$B_r$	00:45:14	0.0001	
Flow				
	$S_1$	00:05:15	0	
	$S_3$	00:05:11	0	
	$B_r$	00:36:12	0.0089	
OCSVM				
Byte				
	$S_1$	00:05:52	0	
	$S_3$	00:04:32	0.7869	
	$B_r$	19:51:57	0.5220	
Flow				
	$S_1$	00:01:02	0	
	$S_3$	00:01:07	0	
	$B_r$	10:13:46	0.5220	
LOF				
Byte				
	$S_1$	00:01:19	0	
	$S_3$	00:05:50	0	
	$B_r$	40:13:40	0	
Flow				
	$S_1$	00:01:12	0	
	$S_3$	01:24:37	0	
	$B_r$	06:00:33	0	
$B_n$				
Byte				
	$S_1$	05:31:37	0.0035	
	$S_3$	11:07:54	0.0234	
Flow				
	$S_1$	02:19:16	0.0702	
	$S_3$	04:41:41	0.0014	

Table B.3: SWaT train validation results

AD	Fragments	Model	FPR↓	PR↑	RC↑	$F_1↑$	T2T↓
IF							
Byte							
	$S_1$	0	1	0	0	00:13:29	
	$S_3$	0	0	0	0	00:11:38	
	$B_r$	0	0.9783	0.0120	0.0238	01:15:11	
Flow							
	$S_1$	0	0	0	0	00:09:32	
	$S_3$	0	1	0	0	00:03:33	
	$B_r$	0	0	0	0	00:31:27	
OCSVM							
Byte							
	$S_1$	0	1	0	0	00:03:36	
	$S_3$	0.7915	0.0132	0.7656	0.0260	00:03:14	
	$B_r$	0.4316	0.0022	0.2876	0.0043	62:53:53	
Flow							
	$S_1$	0	1	0	0	00:02:45	
	$S_3$	0	1	0	0	00:01:14	
	$B_r$	0.4316	0.0022	0.2876	0.0043	62:53:53	
LOF							
Byte							
	$S_1$	0	1	0	0	01:32:54	
	$S_3$	0	0	0	0	06:03:28	
	$B_r$	0	0.5294	0.0008	0.0016	69:04:37	
Flow							
	$S_1$	0	0	0	0	01:28:06	
	$S_3$	0	0	0	0	00:21:47	
	$B_r$	0	0.0349	0.0001	0.0001	10:14:01	
$B_n$							
Byte							
	$S_1$	0.01	0.4959	0.7079	<u>0.5832</u>	05:31:39	
	$S_3$	0.0302	0.3059	0.9620	0.4642	11:07:56	
Flow							
	$S_1$	0.0732	0.0185	0.4239	0.0354	02:19:17	
	$S_3$	0.0040	0.2463	0.3306	<u>0.2823</u>	04:41:42	

Table B.4: SWaT anomaly detection results of  $\text{Eval}_H$

## 5 | Evaluation: Test Results (I)

AD	Fragments	Model	T2F↓	FPR↓
IF				
Byte				
$S_1$	00:07:31	0		
$S_3$	00:05:58	0		
$B_r$	00:45:14	0.0001		
Flow				
$S_1$	00:05:15	0		
$S_3$	00:05:11	0		
$B_r$	00:36:12	0.0089		
OCSVM				
Byte				
$S_1$	00:05:52	0		
$S_3$	00:04:32	0.7869		
$B_r$	19:51:57	0.5220		
Flow				
$S_1$	00:01:02	0		
$S_3$	00:01:07	0		
$B_r$	10:13:46	0.5220		
LOF				
Byte				
$S_1$	00:01:19	0		
$S_3$	00:05:50	0		
$B_r$	40:13:40	0		
Flow				
$S_1$	00:01:12	0		
$S_3$	01:24:37	0		
$B_r$	06:00:33	0		
$B_n$				
Byte				
$S_1$	05:31:37	0.0035		
$S_3$	11:07:54	0.0234		
Flow				
$S_1$	02:19:16	0.0702		
$S_3$	04:41:41	0.0014		

Table B.3: SWaT train validation results

AD	Fragments	Model	FPR↓	PR↑	RC↑	$F_1↑$	T2T↓
IF							
Byte							
$S_1$	0	1	0	0	0	00:13:29	
$S_3$	0	0	0	0	0	00:11:38	
$B_r$	0	0.9783	0.0120	0.0238	0.0238	01:15:11	
Flow							
$S_1$	0	0	0	0	0	00:09:32	
$S_3$	0	1	0	0	0	00:03:33	
$B_r$	0	0	0	0	0	00:31:27	
OCSVM							
Byte							
$S_1$	0	1	0	0	0	00:03:36	
$S_3$	0.7915	0.0132	0.7656	0.0260	0.0260	00:03:14	
$B_r$	0.4316	0.0022	0.2876	0.0043	0.0043	62:53:53	
Flow							
$S_1$	0	1	0	0	0	00:02:45	
$S_3$	0	1	0	0	0	00:01:14	
$B_r$	0.4316	0.0022	0.2876	0.0043	0.0043	62:53:53	
LOF							
Byte							
$S_1$	0	1	0	0	0	01:32:54	
$S_3$	0	0	0	0	0	06:03:28	
$B_r$	0	0.5294	0.0008	0.0016	0.0016	69:04:37	
Flow							
$S_1$	0	0	0	0	0	01:28:06	
$S_3$	0	0	0	0	0	00:21:47	
$B_r$	0	0.0349	0.0001	0.0001	0.0001	10:14:01	
$B_n$							
Byte							
$S_1$	0.01	0.4959	0.7079	<u>0.5832</u>	0.5832	05:31:39	
$S_3$	0.0302	0.3059	0.9620	<u>0.4642</u>	0.4642	11:07:56	
Flow							
$S_1$	0.0732	0.0185	0.4239	0.0354	0.0354	02:19:17	
$S_3$	0.0040	0.2463	0.3306	<u>0.2823</u>	0.2823	04:41:42	

Table B.4: SWaT anomaly detection results of  $\text{Eval}_H$

## 5 | Evaluation: Test Results (I)

AD	Fragments	Model	T2F↓	FPR↓
IF				
Byte	$S_1$	00:07:31	0	
	$S_3$	00:05:58	0	
	$B_r$	00:45:14	0.0001	
Flow	$S_1$	00:05:15	0	
	$S_3$	00:05:11	0	
	$B_r$	00:36:12	0.0089	
OCSVM				
Byte	$S_1$	00:05:52	0	
	$S_3$	00:04:32	0.7869	
	$B_r$	19:51:57	0.5220	
Flow	$S_1$	00:01:02	0	
	$S_3$	00:01:07	0	
	$B_r$	10:13:46	0.5220	
LOF				
Byte	$S_1$	00:01:19	0	
	$S_3$	00:05:50	0	
	$B_r$	40:13:40	0	
Flow	$S_1$	00:01:12	0	
	$S_3$	01:24:37	0	
	$B_r$	06:00:33	0	
$B_n$				
Byte	$S_1$	05:31:37	0.0035	
	$S_3$	11:07:54	0.0234	
Flow	$S_1$	02:19:16	0.0702	
	$S_3$	04:41:41	0.0014	

Table B.3: SWaT train validation results

Over-fitted to training data

AD	Fragments	Model	FPR↓	PR↑	RC↑	$F_1\uparrow$	T2T↓
IF							
Byte	$S_1$	0	1	0	0	00:13:29	
	$S_3$	0	0	0	0	00:11:38	
	$B_r$	0	0.9783	0.0120	0.0238	01:15:11	
Flow	$S_1$	0	0	0	0	00:09:32	
	$S_3$	0	1	0	0	00:03:33	
	$B_r$	0	0	0	0	00:31:27	
OCSVM							
Byte	$S_1$	0	1	0	0	00:03:36	
	$S_3$	0.7915	0.0132	0.7656	0.0260	00:03:14	
	$B_r$	0.4316	0.0022	0.2876	0.0043	62:53:53	
Flow	$S_1$	0	1	0	0	00:02:45	
	$S_3$	0	1	0	0	00:01:14	
	$B_r$	0.4316	0.0022	0.2876	0.0043	62:53:53	
LOF							
Byte	$S_1$	0	1	0	0	01:32:54	
	$S_3$	0	0	0	0	06:03:28	
	$B_r$	0	0.5294	0.0008	0.0016	69:04:37	
Flow	$S_1$	0	0	0	0	01:28:06	
	$S_3$	0	0	0	0	00:21:47	
	$B_r$	0	0.0349	0.0001	0.0001	10:14:01	
$B_n$							
Byte	$S_1$	0.01	0.4959	0.7079	0.5832	05:31:39	
	$S_3$	0.0302	0.3059	0.9620	0.4642	11:07:56	
Flow	$S_1$	0.0732	0.0185	0.4239	0.0354	02:19:17	
	$S_3$	0.0040	0.2463	0.3306	0.2823	04:41:42	

Table B.4: SWaT anomaly detection results of  $\text{Eval}_H$

## 5 | Evaluation: Test Results (I)

AD	Fragments	Model	T2F↓	FPR↓
IF				
Byte	$S_1$	00:07:31	0	
	$S_3$	00:05:58	0	
	$B_r$	00:45:14	0.0001	
Flow	$S_1$	00:05:15	0	
	$S_3$	00:05:11	0	
	$B_r$	00:36:12	0.0089	
OCSVM				
Byte	$S_1$	00:05:52	0	
	$S_3$	00:04:32	0.7869	
	$B_r$	19:51:57	0.5220	
Flow	$S_1$	00:01:02	0	
	$S_3$	00:01:07	0	
	$B_r$	10:13:46	0.5220	
LOF				
Byte	$S_1$	00:01:19	0	
	$S_3$	00:05:50	0	
	$B_r$	40:13:40	0	
Flow	$S_1$	00:01:12	0	
	$S_3$	01:24:37	0	
	$B_r$	06:00:38	0	
$B_n$				
Byte	$S_1$	05:31:37	0.0035	
	$S_3$	11:07:54	0.0234	
Flow	$S_1$	02:19:16	0.0702	
	$S_3$	04:41:41	0.0014	

Table B.3: SWaT train validation results

Over-fitted to training data

Long training time

AD	Fragments	Model	FPR↓	PR↑	RC↑	$F_1↑$	T2T↓
IF							
Byte	$S_1$	0	1	0	0	00:13:29	
	$S_3$	0	0	0	0	00:11:38	
	$B_r$	0	0.9783	0.0120	0.0238	01:15:11	
Flow	$S_1$	0	0	0	0	00:09:32	
	$S_3$	0	1	0	0	00:03:33	
	$B_r$	0	0	0	0	00:31:27	
OCSVM							
Byte	$S_1$	0	1	0	0	00:03:36	
	$S_3$	0.7915	0.0132	0.7656	0.0260	00:03:14	
	$B_r$	0.4316	0.0022	0.2876	0.0043	62:53:53	
Flow	$S_1$	0	1	0	0	00:02:45	
	$S_3$	0	1	0	0	00:01:14	
	$B_r$	0.4316	0.0022	0.2876	0.0043	62:53:53	
LOF							
Byte	$S_1$	0	1	0	0	01:32:54	
	$S_3$	0	0	0	0	06:03:28	
	$B_r$	0	0.5294	0.0008	0.0016	69:04:37	
Flow	$S_1$	0	0	0	0	01:28:06	
	$S_3$	0	0	0	0	00:21:47	
	$B_r$	0	0.0349	0.0001	0.0001	10:14:01	
$B_n$							
Byte	$S_1$	0.01	0.4959	0.7079	0.5832	05:31:39	
	$S_3$	0.0302	0.3059	0.9620	0.4642	11:07:56	
Flow	$S_1$	0.0732	0.0185	0.4239	0.0354	02:19:17	
	$S_3$	0.0040	0.2463	0.3306	0.2823	04:41:42	

Table B.4: SWaT anomaly detection results of *Eval<sub>H</sub>*

## 5 | Evaluation: Test Results (I)

AD	Fragments	Model	T2F↓	FPR↓
IF				
Byte	$S_1$	00:07:31	0	
	$S_3$	00:05:58	0	
	$B_r$	00:45:14	0.0001	
Flow	$S_1$	00:05:15	0	
	$S_3$	00:05:11	0	
	$B_r$	00:36:12	0.0089	
OCSVM				
Byte	$S_1$	00:05:52	0	
	$S_3$	00:04:32	0.7869	
	$B_r$	19:51:57	0.5220	
Flow	$S_1$	00:01:02	0	
	$S_3$	00:01:07	0	
	$B_r$	10:13:46	0.5220	
LOF				
Byte	$S_1$	00:01:19	0	
	$S_3$	00:05:50	0	
	$B_r$	40:13:40	0	
Flow	$S_1$	00:01:12	0	
	$S_3$	01:24:37	0	
	$B_r$	06:00:38	0	
$B_n$				
Byte	$S_1$	05:31:37	0.0035	
	$S_3$	11:07:54	0.0234	
Flow	$S_1$	02:19:16	0.0702	
	$S_3$	04:41:41	0.0014	

Over-fitted to training data

Long training time

Better than guessing

AD	Fragments	Model	FPR↓	PR↑	RC↑	$F_1↑$	T2T↓
IF							
Byte	$S_1$	0	1	0	0	0	00:13:29
	$S_3$	0	0	0	0	0	00:11:38
	$B_r$	0	0.9783	0.0120	0.0238	0	01:15:11
Flow	$S_1$	0	0	0	0	0	00:09:32
	$S_3$	0	1	0	0	0	00:03:33
	$B_r$	0	0	0	0	0	00:31:27
OCSVM							
Byte	$S_1$	0	1	0	0	0	00:03:36
	$S_3$	0.7915	0.0132	0.7656	0.0260	0	00:03:14
	$B_r$	0.4316	0.0022	0.2876	0.0043	0	62:53:53
Flow	$S_1$	0	1	0	0	0	00:02:45
	$S_3$	0	1	0	0	0	00:01:14
	$B_r$	0.4316	0.0022	0.2876	0.0043	0	62:53:53
LOF							
Byte	$S_1$	0	1	0	0	0	01:32:54
	$S_3$	0	0	0	0	0	06:03:28
	$B_r$	0	0.5294	0.0008	0.0016	0	69:04:37
Flow	$S_1$	0	0	0	0	0	01:28:06
	$S_3$	0	0	0	0	0	00:21:47
	$B_r$	0	0.0349	0.0001	0.0001	0	10:14:01
$B_n$							
Byte	$S_1$	0.01	0.4959	0.7079	0.5832	0	05:31:39
	$S_3$	0.0302	0.3059	0.9620	0.4642	0	11:07:56
Flow	$S_1$	0.0732	0.0185	0.4239	0.0354	0	02:19:17
	$S_3$	0.0040	0.2463	0.3306	0.2823	0	04:41:42

Table B.3: SWaT train validation results

Table B.4: SWaT anomaly detection results of *Eval<sub>H</sub>*

## 5 | Evaluation: Test Results (II)

AD	Fragments	Model	FPR↓	PR↑	RC↑	$F_1\uparrow$	T2T↓
IF							
Byte							
	$S_1$	0.9999	0.0371	0.9999	0.0716	00:01:08	
	$S_3$	0.9987	0.0511	0.9992	<u>0.0973</u>	00:01:09	
	$B_r$	1	0.0371	1	0.0716	00:18:27	
Flow							
	$S_1$	0.0705	0.0261	0.1457	0.0442	00:00:28	
	$S_3$	0	1	0	0	00:00:17	
	$B_r$	0.9346	0.0137	1	0.0270	00:08:07	
OCSVM							
Byte							
	$S_1$	1	0.0371	1	0.0716	00:00:36	
	$S_3$	0.0005	0.0384	0.0004	0.0008	00:00:38	
	$B_r$	1	0.0369	0.9950	0.0712	07:14:04	
Flow							
	$S_1$	1	0.0128	1	0.0253	00:00:16	
	$S_3$	1	0.0129	1	0.0255	00:00:18	
	$B_r$	1	0.0127	0.9902	0.0250	03:52:22	
LOF							
Byte							
	$S_1$	0.9999	0.0371	1	0.0716	06:33:01	
	$S_3$	1	0.0511	1	<u>0.0973</u>	06:14:50	
	$B_r$	1	0.0371	1	0.0716	11:13:28	
Flow							
	$S_1$	0.9387	0.0136	1	0.0269	01:05:51	
	$S_3$	1	0.0129	1	0.0255	00:48:12	
	$B_r$	1	0.0127	0.9952	0.0251	05:32:56	
$B_n$							
Byte							
	$S_1$	0.0091	0.0152	0.0036	0.0059	05:32:02	
	$S_3$	0.0281	0.0543	0.03	0.0386	11:17:33	
Flow							
	$S_1$	0.0085	0.2227	0.1870	<u>0.2033</u>	02:25:40	
	$S_3$	0.0030	0.0618	0.0150	0.0241	04:52:41	

Table B.5: SWaT anomaly detection results of  $Eval_D$

AD	Fragments	Model	FPR↓	PR↑	RC↑	$F_1\uparrow$	T2T↓
IF							
Byte							
	$S_1$	0.9999	0.0004	1		0.0008	00:00:45
	$S_3$	0.9987	0.0008	1		0.0017	00:00:46
	$B_r$	1	0.0004	1		0.0008	00:18:57
Flow							
	$S_1$	0.0706	0.0001	0.1034	0.0003	00:00:18	
	$S_3$	0	1	0	0	00:00:19	
	$B_r$	0.9346	0.0001	1	0.0002	00:07:56	
OCSVM							
Byte							
	$S_1$	1	0	1	0.0001	00:00:12	
	$S_3$	1	0.0001	1	0.0003	00:00:12	
	$B_r$	0.0169	0.0004	0.0169	0.0008	06:54:35	
Flow							
	$S_1$	1	0.0001	1	0.0002	00:00:18	
	$S_3$	1	0.0002	1	0.0004	00:00:12	
	$B_r$	0.0010	0.0001	0.0010	0.0002	02:48:17	
LOF							
Byte							
	$S_1$	0.9999	0.0004	1	0.0008	06:31:27	
	$S_3$	1	0.0008	1	0.0017	06:08:44	
	$B_r$	0.1050	0.0003	0.1050	0.0007	09:46:53	
Flow							
	$S_1$	0.9387	0.0001	1	0.0002	01:05:39	
	$S_3$	1	0.0002	1	0.0004	00:42:48	
	$B_r$	0.0968	0.0001	0.0968	0.0002	04:17:47	
$B_n$							
Byte							
	$S_1$	0.0090	0	0	0	0	05:28:49
	$S_3$	0.0278	0.0015	0.0484	<u>0.0029</u>	11:05:08	
Flow							
	$S_1$	0.0085	0.0020	0.1724	<u>0.0040</u>	02:19:47	
	$S_3$	0.0030	0.0012	0.0182	0.0022	04:50:35	

Table B.6: SWaT anomaly detection results of  $Eval_E$

## 5 | Evaluation: Test Results (II)

AD	Fragments	Model	FPR↓	PR↑	RC↑	$F_1\uparrow$	T2T↓
IF							
Byte							
$S_1$	0.9999	0.0371	0.9999	0.0716	00:01:08		
$S_3$	0.9987	0.0511	0.9992	<u>0.0973</u>	00:01:09		
$B_r$	1	0.0371	1	0.0716	00:18:27		
Flow							
$S_1$	0.0705	0.0261	0.1457	0.0442	00:00:28		
$S_3$	0	1	0	0	00:00:17		
$B_r$	0.9346	0.0137	1	0.0270	00:08:07		
OCSVM							
Byte							
$S_1$	1	0.0371	1	0.0716	00:00:36		
$S_3$	0.0005	0.0384	0.0004	0.0008	00:00:38		
$B_r$	1	0.0369	0.9950	0.0712	07:14:04		
Flow							
$S_1$	1	0.0128	1	0.0253	00:00:16		
$S_3$	1	0.0129	1	0.0255	00:00:18		
$B_r$	1	0.0127	0.9902	0.0250	03:52:22		
LOF							
Byte							
$S_1$	0.9999	0.0371	1	0.0716	06:33:01		
$S_3$	1	0.0511	1	<u>0.0973</u>	06:14:50		
$B_r$	1	0.0371	1	0.0716	11:13:28		
Flow							
$S_1$	0.9387	0.0136	1	0.0269	01:05:51		
$S_3$	1	0.0129	1	0.0255	00:48:12		
$B_r$	1	0.0127	0.9952	0.0251	05:32:56		
$B_n$							
Byte							
$S_1$	0.0091	0.0152	0.0036	0.0059	05:32:02		
$S_3$	0.0281	0.0543	0.03	0.0386	11:17:33		
Flow							
$S_1$	0.0085	0.2227	0.1870	<u>0.2033</u>	02:25:40		
$S_3$	0.0030	0.0618	0.0150	0.0241	04:52:41		

Table B.5: SWaT anomaly detection results of  $Eval_D$

AD	Fragments	Model	FPR↓	PR↑	RC↑	$F_1\uparrow$	T2T↓
IF							
Byte							
$S_1$	0.9999	0.0004	1		0.0008	00:00:45	
$S_3$	0.9987	0.0008	1		0.0017	00:00:46	
$B_r$	1	0.0004	1		0.0008	00:18:57	
Flow							
$S_1$	0.0706	0.0001	0.1034	0.0003	00:00:18		
$S_3$	0	1	0	0	00:00:19		
$B_r$	0.9346	0.0001	1	0.0002	00:07:56		
OCSVM							
Byte							
$S_1$	1	0	1	0.0001	00:00:12		
$S_3$	1	0.0001	1	0.0003	00:00:12		
$B_r$	0.0169	0.0004	0.0169	0.0008	06:54:35		
Flow							
$S_1$	1	0.0001	1	0.0002	00:00:18		
$S_3$	1	0.0002	1	0.0004	00:00:12		
$B_r$	0.0010	0.0001	0.0010	0.0002	02:48:17		
LOF							
Byte							
$S_1$	0.9999	0.0004	1	0.0008	06:31:27		
$S_3$	1	0.0008	1	0.0017	06:08:44		
$B_r$	0.1050	0.0003	0.1050	0.0007	09:46:53		
Flow							
$S_1$	0.9387	0.0001	1	0.0002	01:05:39		
$S_3$	1	0.0002	1	0.0004	00:42:48		
$B_r$	0.0968	0.0001	0.0968	0.0002	04:17:47		
$B_n$							
Byte							
$S_1$	0.0090	0	0	0	0	05:28:49	
$S_3$	0.0278	0.0015	0.0484	<u>0.0029</u>	11:05:08		
Flow							
$S_1$	0.0085	0.0020	0.1724	<u>0.0040</u>	02:19:47		
$S_3$	0.0030	0.0012	0.0182	0.0022	04:50:35		

Table B.6: SWaT anomaly detection results of  $Eval_E$

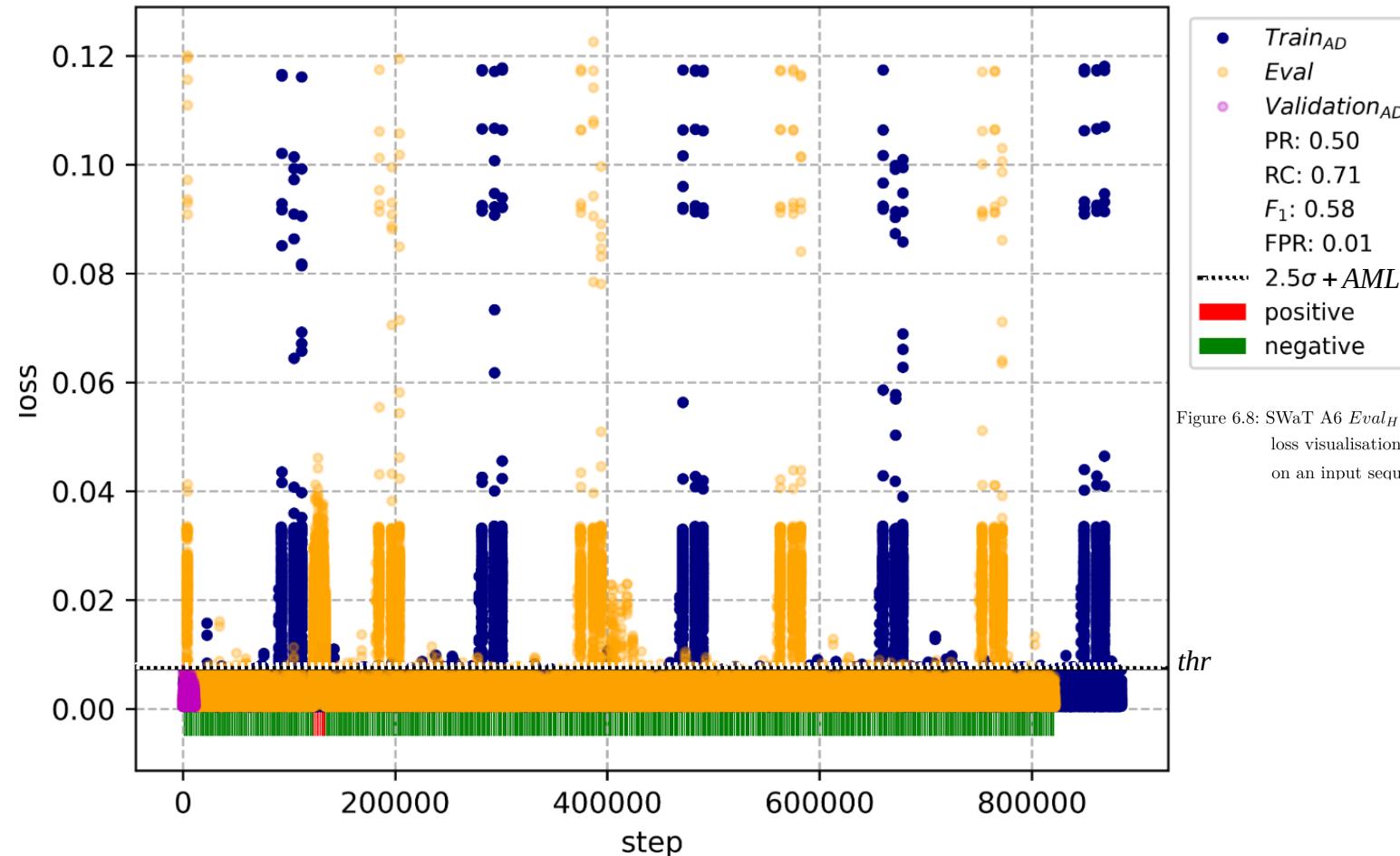
## 5 | Evaluation: Test Results (II)

AD	Fragments	Model	FPR↓	PR↑	RC↑	$F_1\uparrow$	T2T↓
IF							
Byte							
$S_1$	0.9999	0.0371	0.9999	0.0716	00:01:08		
$S_3$	0.9987	0.0511	0.9992	<u>0.0973</u>	00:01:09		
$B_r$	1	0.0371	1	0.0716	00:18:27		
Flow							
$S_1$	0.0705	0.0261	0.1457	0.0442	00:00:28		
$S_3$	0	1	0	0	00:00:17		
$B_r$	0.9346	0.0137	1	0.0270	00:08:07		
OCSVM							
Byte							
$S_1$	1	0.0371	1	0.0716	00:00:36		
$S_3$	0.0005	0.0384	0.0004	0.0008	00:00:38		
$B_r$	1	0.0369	0.9950	<u>0.0712</u>	<u>07:14:04</u>		
Flow							
$S_1$	1	0.0128	1	0.0128			
$S_3$	1	0.0129	1	0.0255	00:00:18		
$B_r$	1	0.0127	0.9902	0.0250	03:52:22		
LOF							
Byte							
$S_1$	0.9999	0.0371	1	0.0716	06:33:01		
$S_3$	1	0.0511	1	0.0973	06:14:50		
$B_r$	1	0.0371	1	0.0716	11:13:28		
Flow							
$S_1$	0.9387	0.0136	1	0.0269	01:05:51		
$S_3$	1	0.0129	1	0.0255	00:48:12		
$B_r$	1	0.0127	0.9952	0.0251	05:32:56		
$B_n$							
Byte							
$S_1$	0.0091	0.0152	0.0036	0.0059	05:32:02		
$S_3$	0.0281	0.0543	0.03	0.0386	11:17:33		
Flow							
$S_1$	0.0085	0.2227	0.1870	<u>0.2033</u>	02:25:40		
$S_3$	0.0030	0.0618	0.0150	0.0241	04:52:41		

Table B.5: SWaT anomaly detection results of *Eval<sub>D</sub>*

AD	Fragments	Model	FPR↓	PR↑	RC↑	$F_1\uparrow$	T2T↓
IF							
Byte							
$S_1$	0.9999	0.0004	1		0.0008	00:00:45	
$S_3$	0.9987	0.0008	1		0.0017	00:00:46	
$B_r$	1	0.0004	1		0.0008	00:18:57	
Flow							
$S_1$	0.0706	0.0001	0.1034	0.0003	00:00:18		
$S_3$	0	1	0	0	00:00:19		
$B_r$	0.9346	0.0001	1	0.0002	00:07:56		
OCSVM							
Byte							
$S_1$	1	0	1	0.0001	00:00:12		
$S_3$	1	0.0001	1	0.0003	00:00:12		
$B_r$	0.0169	0.0004	0.0169	0.0008	06:54:35		
Flow							
$S_1$	0.0001	1	0.0002	00:00:18			
$S_3$	0.0002	1	0.0004	00:00:12			
$B_r$	0.0010	0.0001	0.0010	0.0002	02:48:17		
LOF							
Byte							
$S_1$	0.9999	0.0004	1	0.0008	06:31:27		
$S_3$	1	0.0008	1	0.0017	06:08:44		
$B_r$	0.1050	0.0003	0.1050	0.0007	09:46:53		
Flow							
$S_1$	0.9387	0.0001	1	0.0002	01:05:39		
$S_3$	1	0.0002	1	0.0004	00:42:48		
$B_r$	0.0968	0.0001	0.0968	0.0002	04:17:47		
$B_n$							
Byte							
$S_1$	0.0090	0	0	0	0	05:28:49	
$S_3$	0.0278	0.0015	0.0484	<u>0.0029</u>	11:05:08		
Flow							
$S_1$	0.0085	0.0020	0.1724	<u>0.0040</u>	02:19:47		
$S_3$	0.0030	0.0012	0.0182	0.0022	04:50:35		

Table B.6: SWaT anomaly detection results of *Eval<sub>E</sub>*



## 5 | Evaluation: Naive Baseline

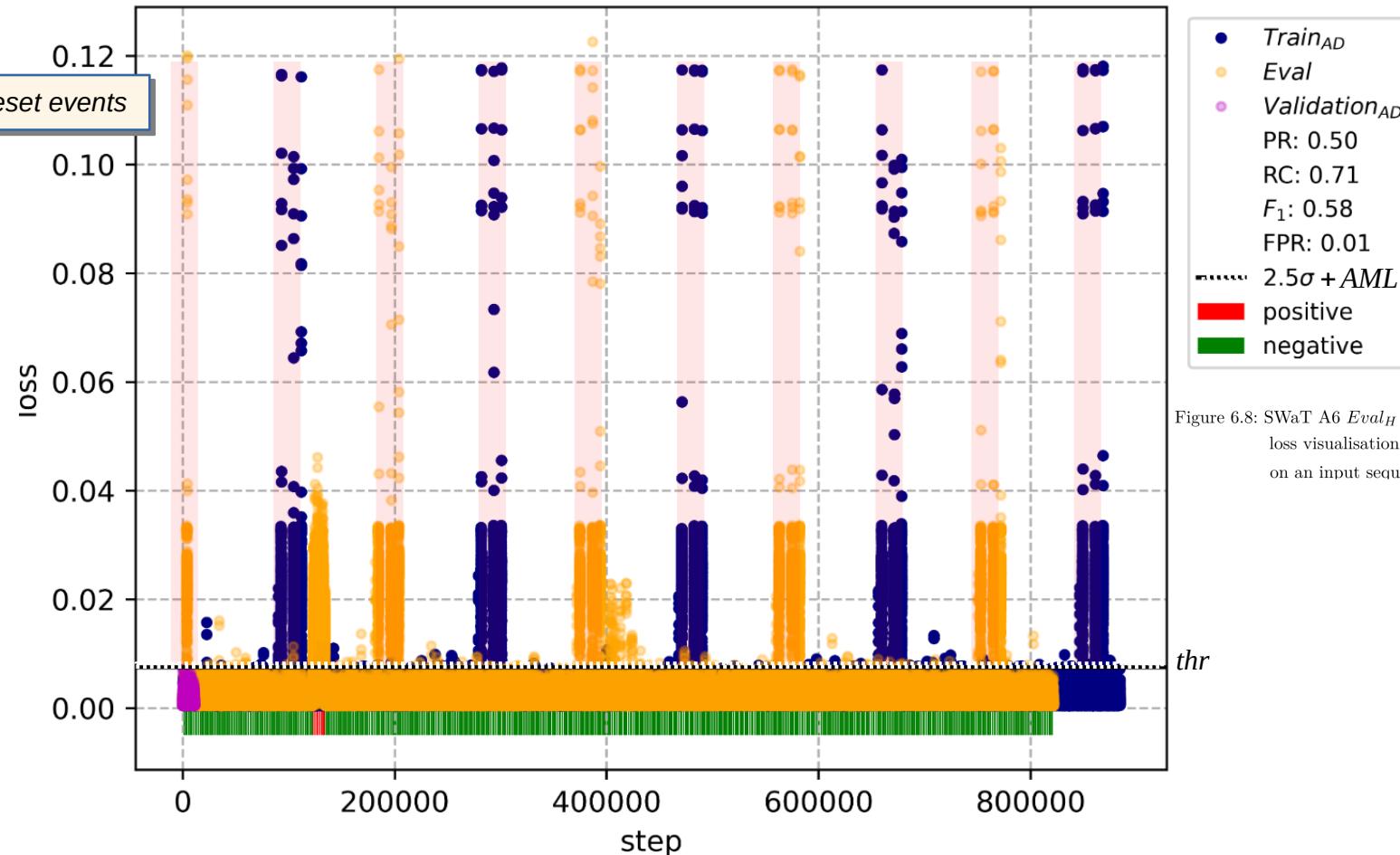


Figure 6.8: SWaT A6  $Eval_H$  naive baseline evaluation loss visualisation for byte-fragments trained on an input sequence length of  $n=1$

## 5 | Evaluation: Naive Baseline

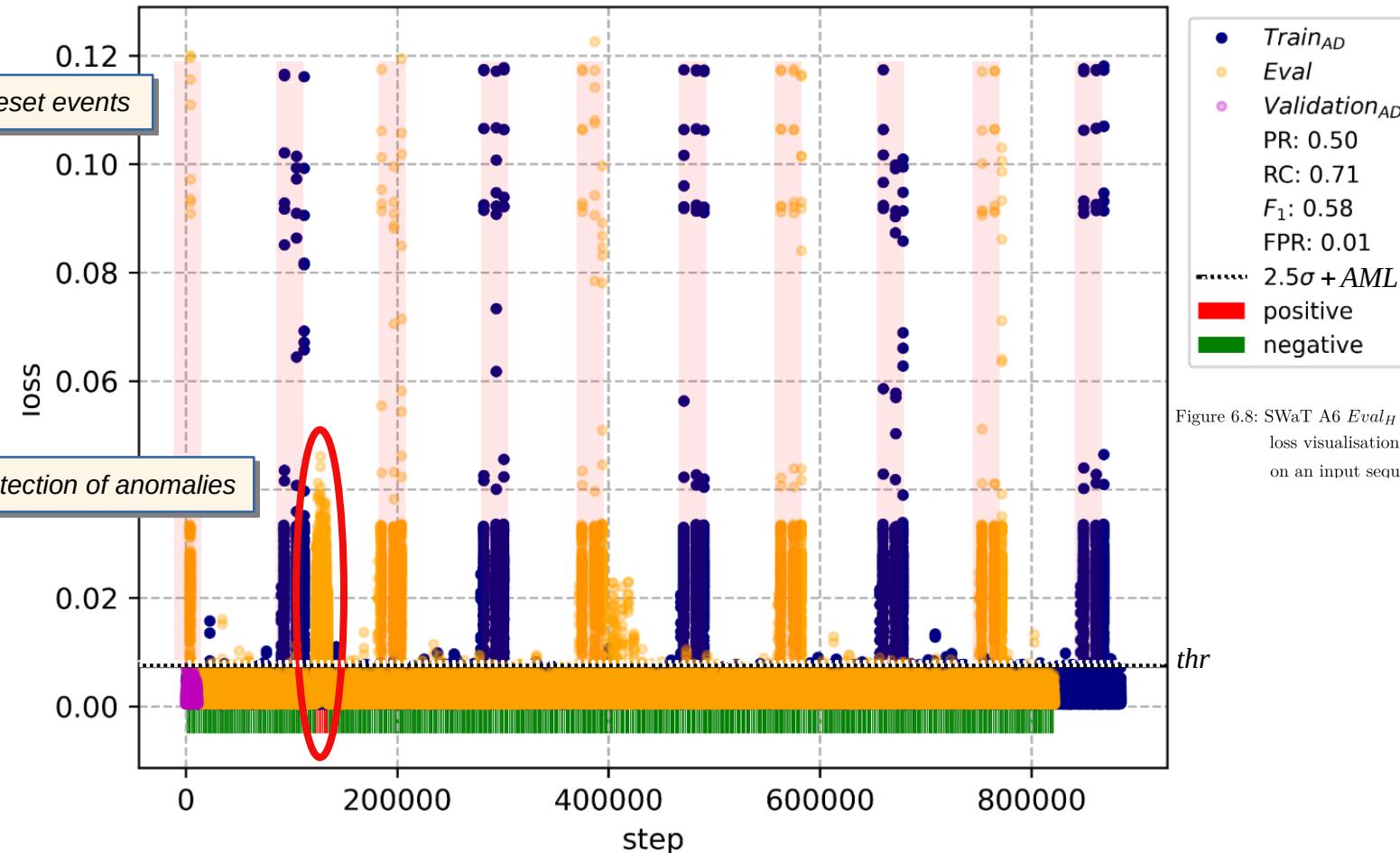


Figure 6.8: SWaT A6  $\text{Eval}_H$  naive baseline evaluation loss visualisation for byte-fragments trained on an input sequence length of  $n=1$

## 5 | Evaluation: Naive Baseline

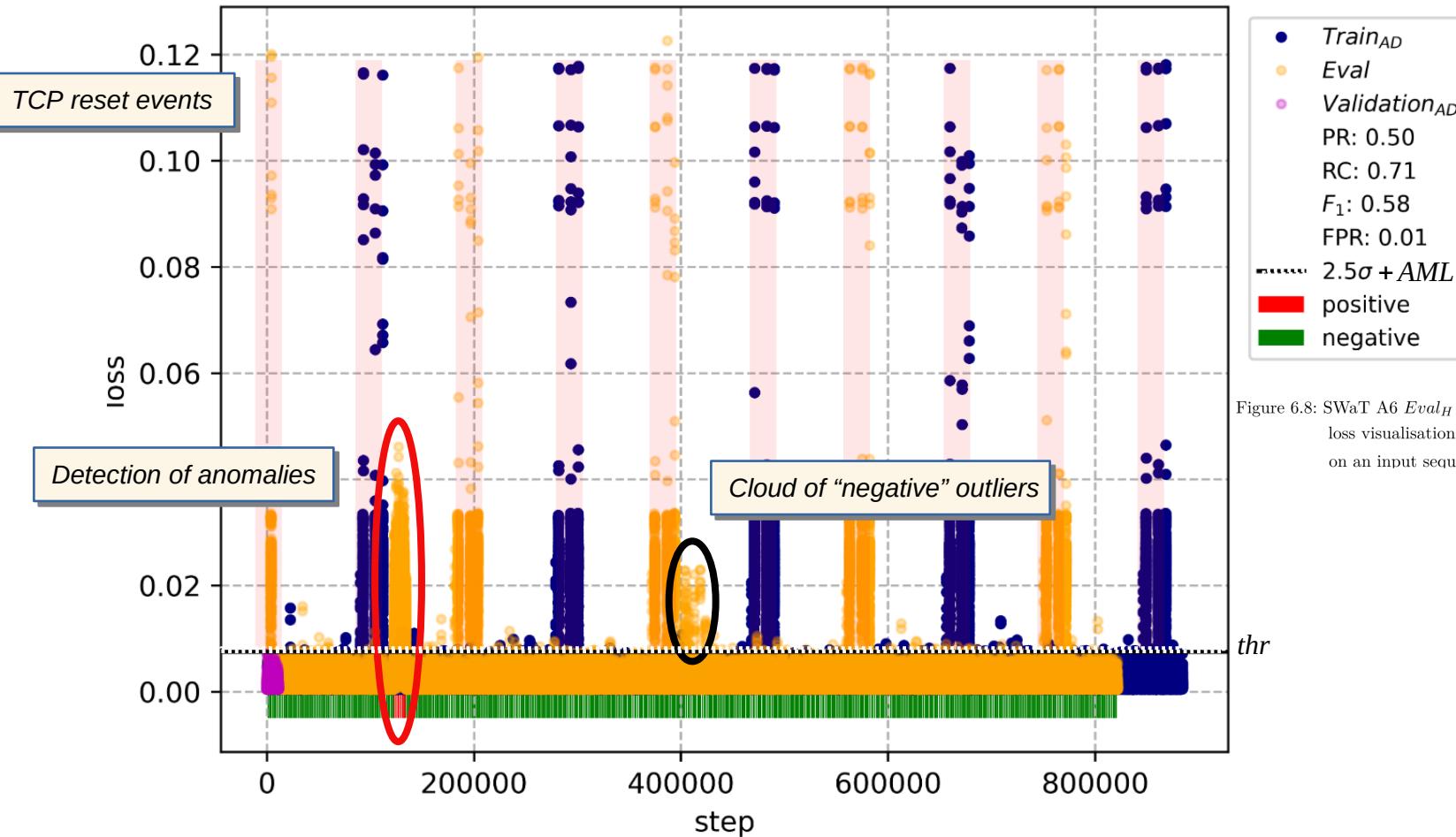


Figure 6.8: SWaT A6  $Eval_H$  naive baseline evaluation loss visualisation for byte-fragments trained on an input sequence length of  $n=1$

## 5 | Evaluation: Decision Analysis

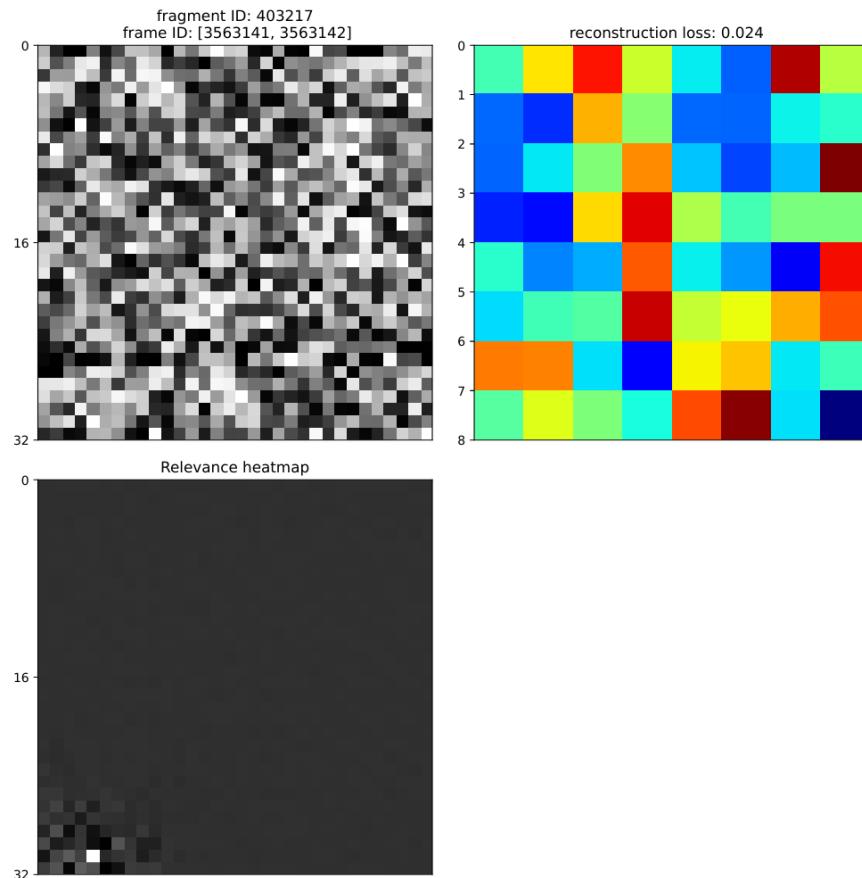


Figure B.9: Byte-fragments relevance heatmaps of the unusual UDP packet

# 5 | Evaluation: Decision Analysis

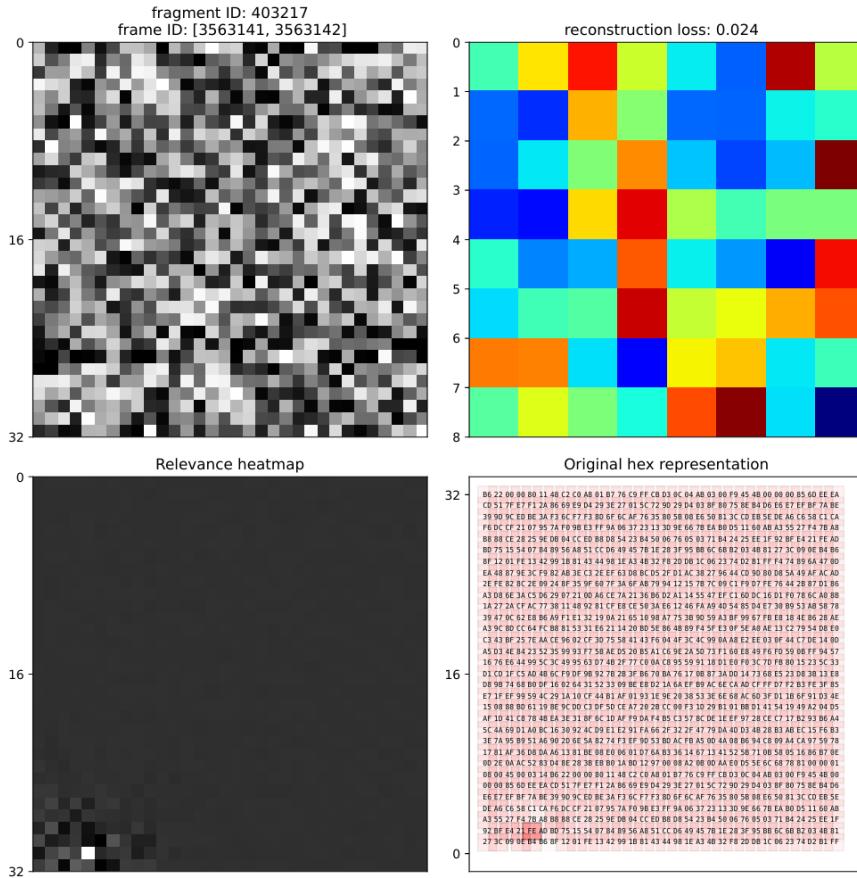


Figure B.9: Byte-fragments relevance heatmaps of the unusual UDP packet

# 5 | Evaluation: Decision Analysis

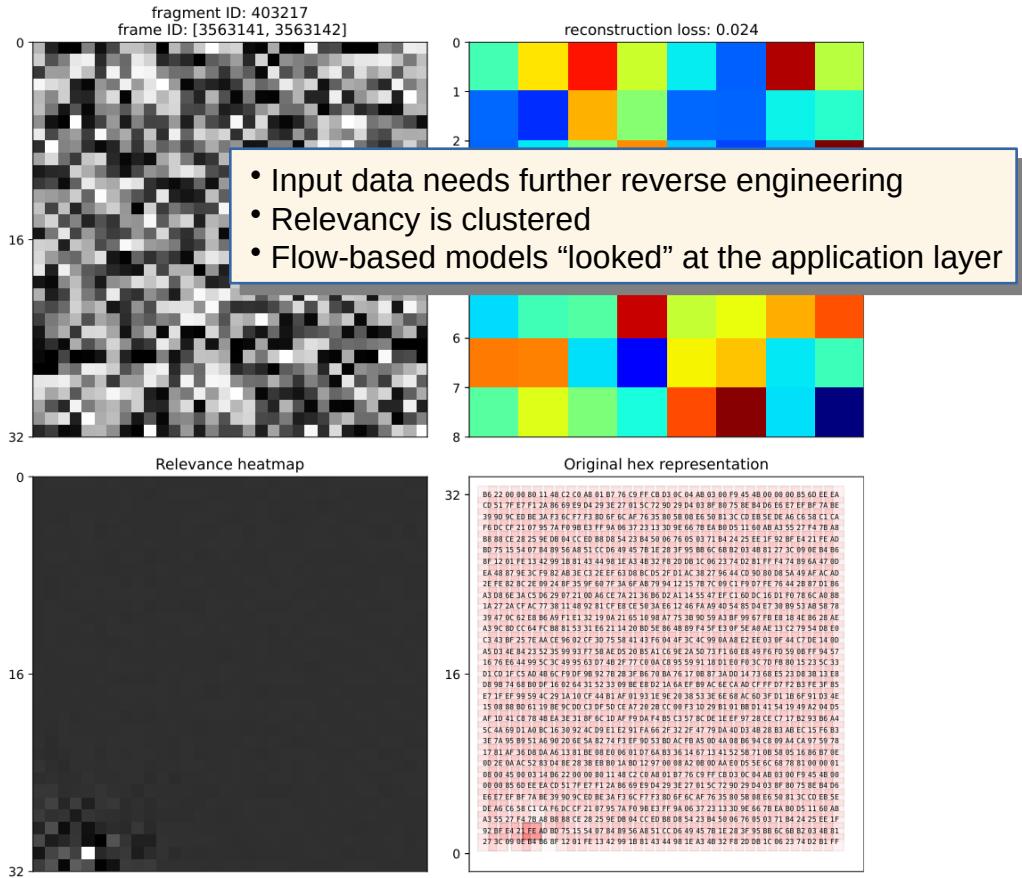


Figure B.9: Byte-fragments relevance heatmaps of the unusual UDP packet

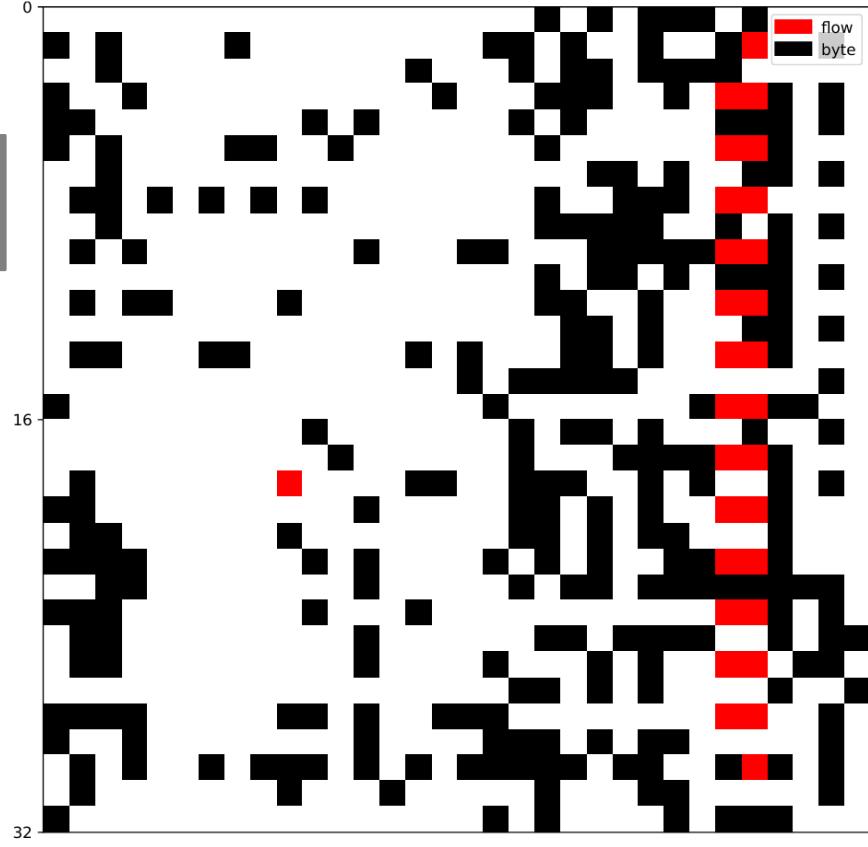


Figure 6.11: Aggregated relevance heatmaps for 1% of randomly selected evaluation fragments

- Hybrid model produces poor results
  - Under-fitted feature extraction model
    - Model loss fluctuates
    - Premature training stop (@6 epochs)
    - Small/fixed validation set
  - Over-fitted anomaly detection model
    - Biased hyper-parameter search

- Hybrid model produces poor results
  - Under-fitted feature extraction model
    - Model loss fluctuates
    - Premature training stop (@6 epochs)
    - Small/fixed validation set
  - Over-fitted anomaly detection model
    - Biased hyper-parameter search
- Detection models took longer to train on raw features
  - Raw baseline model not advantageous

- Hybrid model produces poor results
  - Under-fitted feature extraction model
    - Model loss fluctuates
    - Premature training stop (@6 epochs)
    - Small/fixed validation set
  - Over-fitted anomaly detection model
    - Biased hyper-parameter search
- Detection models took longer to train on raw features
  - Raw baseline model not advantageous
- Naive baseline model identified anomalies
  - Extracted features differed from normal samples

- Hybrid model produces poor results
  - Under-fitted feature extraction model
    - Model loss fluctuates
    - Premature training stop (@6 epochs)
    - Small/fixed validation set
  - Over-fitted anomaly detection model
    - Biased hyper-parameter search
- Detection models took longer to train on raw features
  - Raw baseline model not advantageous
- Naive baseline model identified anomalies
  - Extracted features differed from normal samples

- Byte-fragments outperformed flow-fragments
  - Byte-Fragments contain more information
  - Fewer flow-fragments were generated

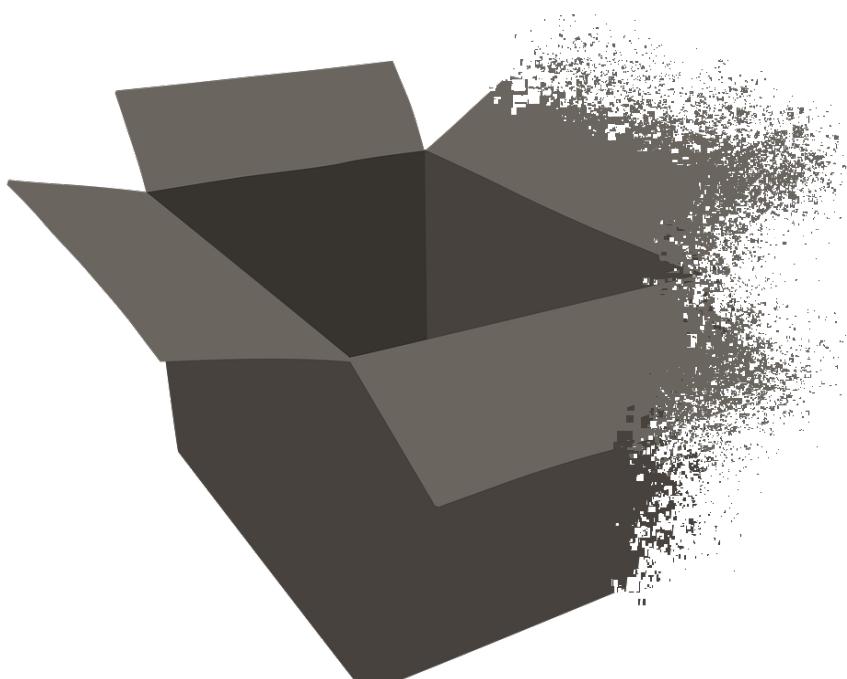
- Hybrid model produces poor results
  - Under-fitted feature extraction model
    - Model loss fluctuates
    - Premature training stop (@6 epochs)
    - Small/fixed validation set
  - Over-fitted anomaly detection model
    - Biased hyper-parameter search
- Detection models took longer to train on raw features
  - Raw baseline model not advantageous
- Naive baseline model identified anomalies
  - Extracted features differed from normal samples

- Byte-fragments outperformed flow-fragments
  - Byte-Fragments contain more information
  - Fewer flow-fragments were generated
- Inter-packet anomalies were not identified
  - Feature extraction learned temporal aspects
  - Detecting of temporal deviations were not demonstrated

- Hybrid model produces poor results
  - Under-fitted feature extraction model
    - Model loss fluctuates
    - Premature training stop (@6 epochs)
    - Small/fixed validation set
  - Over-fitted anomaly detection model
    - Biased hyper-parameter search
- Detection models took longer to train on raw features
  - Raw baseline model not advantageous
- Naive baseline model identified anomalies
  - Extracted features differed from normal samples

- Byte-fragments outperformed flow-fragments
  - Byte-Fragments contain more information
  - Fewer flow-fragments were generated
- Inter-packet anomalies were not identified
  - Feature extraction learned temporal aspects
  - Detecting of temporal deviations were not demonstrated
- Heatmaps reveal insights about extracted features
  - LRP does explain *why* certain areas were important
  - Manual reverse engineering is needed for analysis

- 0 | Preamble
- 1 | Background
- 2 | Related Work
- 3 | Methodology
- 4 | Implementation
- 5 | Evaluation
- 6 | Summary**



- Introduction of an ANN-based feature extraction framework
  - *Unsupervised, protocol-agnostic and content sensitive* feature extraction
  - Extraction of spatial- & temporal-features from a stream of frames
  - Offer capabilities to detect anomalies with out the need for additional models (end-to-end detection)
  - Helps to bridges the *semantic gap* by offering LRP-generated heatmaps for alerts

- Introduction of an ANN-based feature extraction framework
  - *Unsupervised, protocol-agnostic and content sensitive* feature extraction
  - Extraction of spatial- & temporal-features from a stream of frames
  - Offer capabilities to detect anomalies with out the need for additional models (end-to-end detection)
  - Helps to bridges the *semantic gap* by offering LRP-generated heatmaps for alerts
- Comparison between different detection configurations over 2 data sets with 2 baselines
  - Input data representations (byte- & flow-fragments)
  - Small & longer input sequences (S1 & S3)
  - Isolation Forest & OC-SVM & Local outlier factor

- Introduction of an ANN-based feature extraction framework
  - *Unsupervised, protocol-agnostic and content sensitive* feature extraction
  - Extraction of spatial- & temporal-features from a stream of frames
  - Offer capabilities to detect anomalies with out the need for additional models (end-to-end detection)
  - Helps to bridges the *semantic gap* by offering LRP-generated heatmaps for alerts
- Comparison between different detection configurations over 2 data sets with 2 baselines
  - Input data representations (byte- & flow-fragments)
  - Small & longer input sequences (S1 & S3)
  - Isolation Forest & OC-SVM & Local outlier factor
- Bad evaluation results
  - Feature extraction model training was stopped prematurely
  - Validation data was not representative
  - Limited hyper-parameter grid-search (AE & AD)
  - Unbalanced evaluation (more byte-fragments then flow-fragments)
  - Estimated ground truth labels

I. Are computing systems able to perform for a spatial-temporal extraction of feature from raw data ?

- Detection of intra-packet anomalies demonstrated the extraction of spatial features
- Temporal extraction was demonstrated through constructed experiments

I. Are computing systems able to perform for a spatial-temporal extraction of feature from raw data ?

- Detection of intra-packet anomalies demonstrated the extraction of spatial features
- Temporal extraction was demonstrated through constructed experiments

II. Can anomaly detection be improved when trained on extracted features ?

- Conducted evaluation does not give evidence to support this claim

I. Are computing systems able to perform for a spatial-temporal extraction of feature from raw data ?

- Detection of intra-packet anomalies demonstrated the extraction of spatial features
- Temporal extraction was demonstrated through constructed experiments

II. Can anomaly detection be improved when trained on extracted features ?

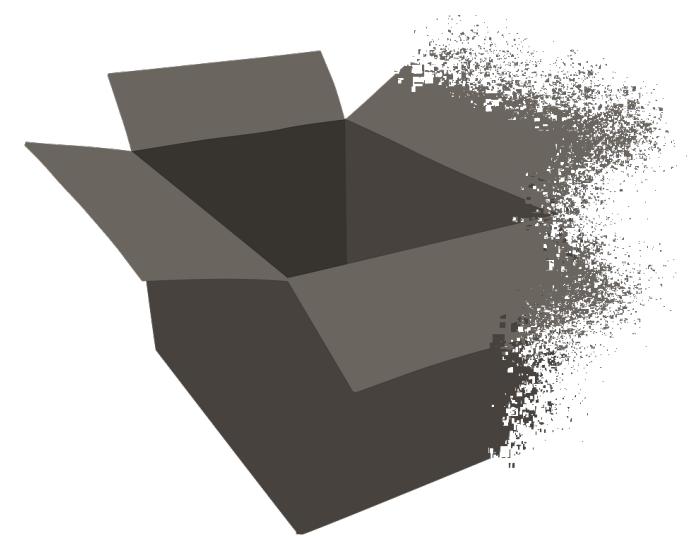
- Conducted evaluation does not give evidence to support this claim

III. What type of input data representations yields better results for different anomalies ?

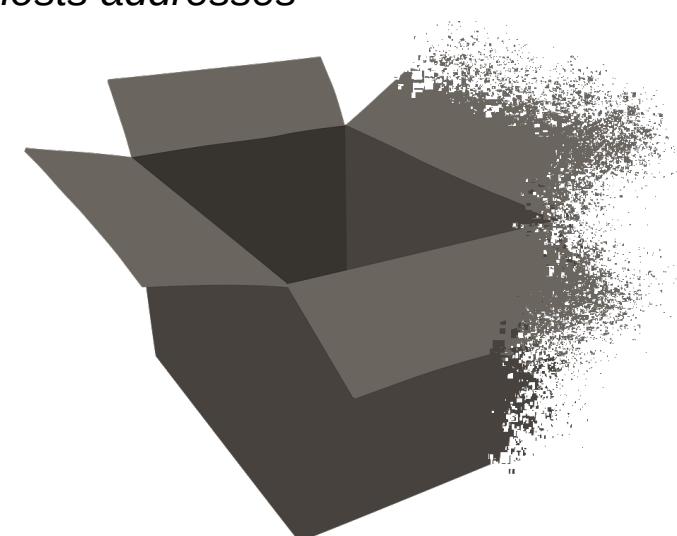
- Byte-fragments were advantageous (higher F1-Score) for intra-packet anomalies
- Inter-packet anomalies were not detected overall

- I. Are computing systems able to perform for a spatial-temporal extraction of feature from raw data ?
  - Detection of intra-packet anomalies demonstrated the extraction of spatial features
  - Temporal extraction was demonstrated through constructed experiments
- II. Can anomaly detection be improved when trained on extracted features ?
  - Conducted evaluation does not give evidence to support this claim
- III. What type of input data representations yields better results for different anomalies ?
  - Byte-fragments were advantageous (higher F1-Score) for intra-packet anomalies
  - Inter-packet anomalies were not detected overall
- IV. How can the semantic gap be closed with the help of explainable models ?
  - Auxiliary techniques (LRP) offer a natural investigation for ANN

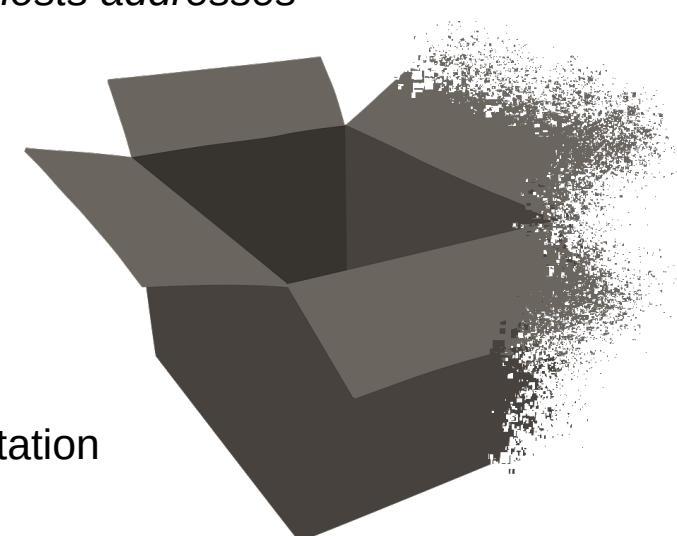
- Simplify problem space
  - *Reassembly* of packets before features are extracted
  - *Restrict* model to train on only certain *network protocols / hosts addresses*



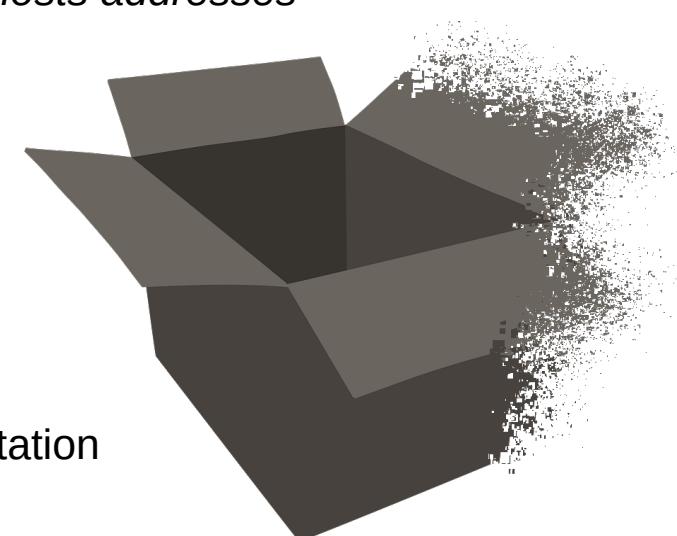
- Simplify problem space
  - *Reassembly* of packets before features are extracted
  - *Restrict* model to train on only certain *network protocols / hosts addresses*
- Extend feature extraction model
  - Introduce more layers per sub-network / longer training
  - Train to reconstruct future fragments
  - Process *1d representations* of network data



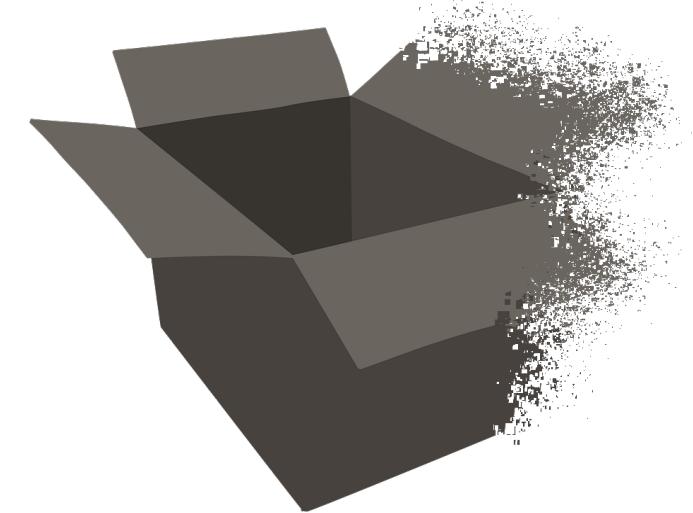
- Simplify problem space
  - Reassembly of packets before features are extracted
  - Restrict model to train on only certain *network protocols / hosts addresses*
- Extend feature extraction model
  - Introduce more layers per sub-network / longer training
  - Train to reconstruct future fragments
  - Process *1d representations* of network data
- Extend model interpretability
  - Use *negative & positive relevance propagation*
  - Introduce protocol decoders and UI's for semantic interpretation



- Simplify problem space
  - Reassembly of packets before features are extracted
  - Restrict model to train on only certain *network protocols / hosts addresses*
- Extend feature extraction model
  - Introduce more layers per sub-network / longer training
  - Train to reconstruct future fragments
  - Process *1d representations* of network data
- Extend model interpretability
  - Use *negative & positive relevance propagation*
  - Introduce protocol decoders and UI's for semantic interpretation
- Utilize synthetic anomalies
  - Aid the feature extraction training
  - Aid the anomaly detection training



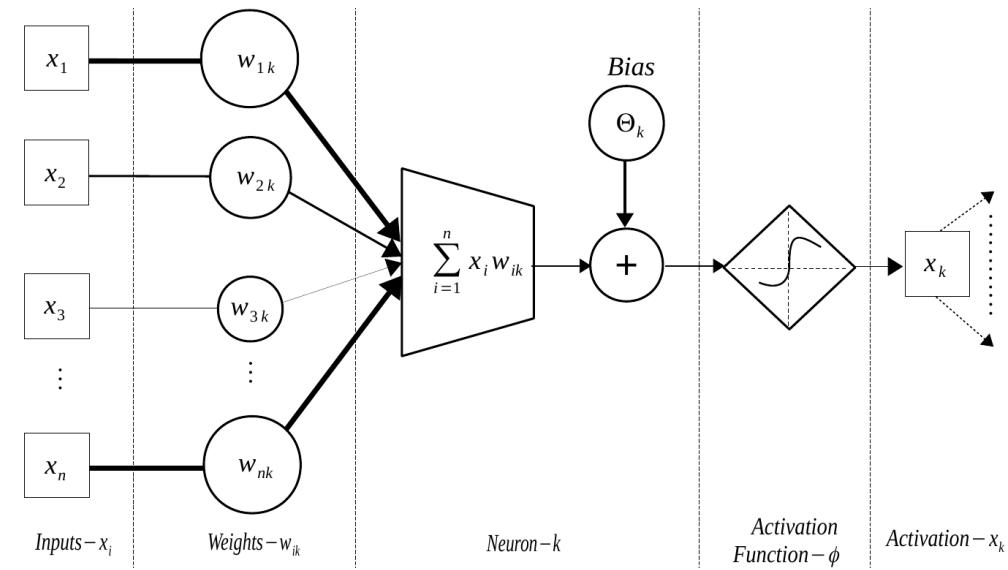
**Thanks for the attention!**  
Questions?



<git.informatik.tu-cottbus.de/koppfabi/thesis>

- Y. Bengio et al. Contractive Auto-Encoders: Explicit Invariance During Feature Extraction – **ICML 2008**
- J. Schmidhuber et al. Stacked Convolutional Auto-Encoders for Hierarchical Feature Extraction – **ICANN 2011**
- B. Frey et al. Winner-Take-All Autoencoders – **NIPS 2015**
- K. Kim et al. RaPP: Novelty Detection with Reconstruction along Projection Pathway – **ICLR 2020**
- A. Binder et al. Layer-Wise Relevance Propagation for Neural Networks with Local Renormalization Layers – **ICANN 2016**
- P. Schneider et al. High-Performance Unsupervised Anomaly Detection for Cyber-Physical System Networks – **CPS-SPC 2018**
- P. Casas et al. DeepMAL – Deep Learning Models for Malware Traffic Detection and Classification – **iDSC 2020**
- ZQ. Qin et al. Attentional Payload Anomaly Detector for Web Applications – **ICONIP 2018**
- X. Shi et al. Convolutional LSTM Network: A Machine Learning Approach for Precipitation Nowcasting – **NIPS 2015**
- J. Nocedal et al. On Large-Batch Training for Deep Learning: Generalization Gap and Sharp Minima – **ICLR 2017**
- V. Nair et al. Rectified Linear Units Improve Restricted Boltzmann Machines – **ICML 2010**
- L. Smith et al. Super-Convergence: Very Fast Training of Neural Networks Using Large Learning Rates – **SPIE 2019**
- N. Tippenhauer et al. SWaT: A water treatment testbed for research and training on ICS security – **CySWater 2016**
- R. Hanusch Synthetische Generierung von netzbasierten Angriffsdaten – **BTU 2018**

- Hierarchical network of layers of neurons
  - Extracts information through the hierarchy
- Loss function is used to judge how wrong the model's decisions are
- ANN trained by minimizing an optimization problem
- Advantages:
  - (+) No feature extraction need
  - (+) Scales with more data & more compute
- Disadvantage:
  - (-) Many hyper-parameters & model architectures
  - (-) Complex models are hard to train



- Transformation of input data into a set of derived values (features)
- Features intended to be informative, non-redundant and describe the data with sufficient accuracy

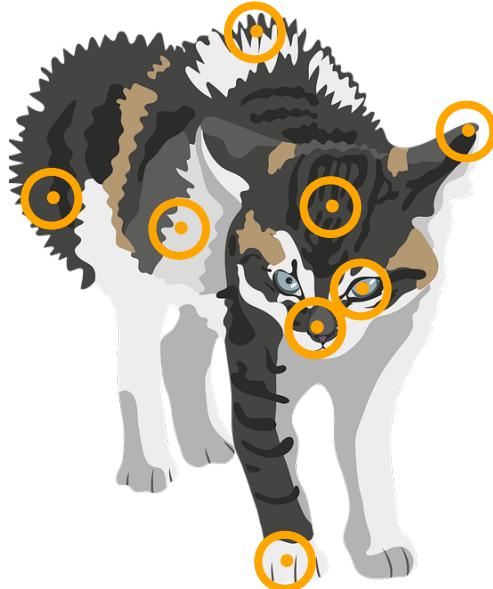
- Transformation of input data into a set of derived values (features)
- Features intended to be informative, non-redundant and describe the data with sufficient accuracy
- Special form of dimensional reduction

- Transformation of input data into a set of derived values (features)
- Features intended to be informative, non-redundant and describe the data with sufficient accuracy
- Special form of dimensional reduction



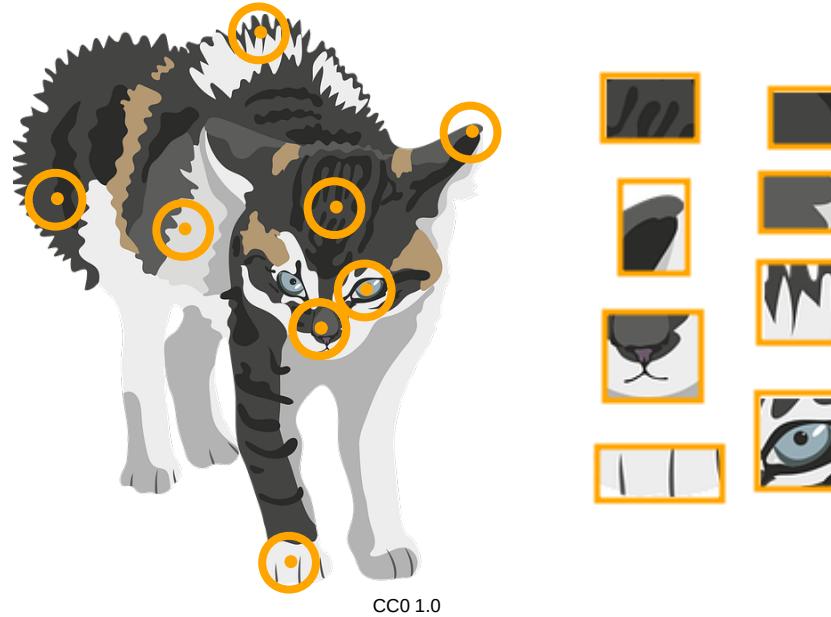
CC0 1.0

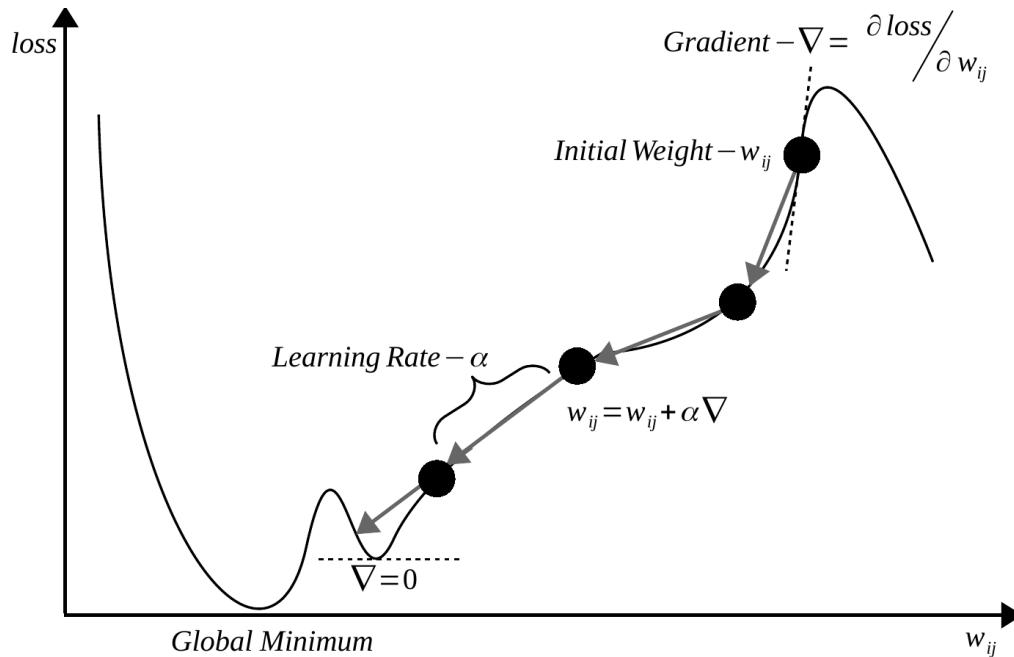
- Transformation of input data into a set of derived values (features)
- Features intended to be informative, non-redundant and describe the data with sufficient accuracy
- Special form of dimensional reduction



CC0 1.0

- Transformation of input data into a set of derived values (features)
- Features intended to be informative, non-redundant and describe the data with sufficient accuracy
- Special form of dimensional reduction
- Representation learning (feature learning) - automatically discovery of suitable representations





- *Loss function* is used to measure how *good* the network understood the problem
- Error is propagated back after each *batch* of samples → *backpropagation*
- *Partial derivatives for each weight* (gradient) shows direction where weight will minimize
- *Learn rate* dictates amount of change in each step

## Network Data

- Series of *packets*
- Weak order of packets
- Packets consist of *layers*
- Different types of protocols
- Layers have *header* and *payload* fields

00000260:	54	50	2f	31	2e	31	0d	0a	48	6f	73	74	3a	20	78	6e	TP/1.1 · Host: xn		
00000270:	2d	2d	6d	62	69	75	73	2d	6a	75	61	2c	62	61	6e	64	--mbius-jua.band		
00000280:	0d	07	55	73	65	72	2d	41	67	65	6e	74	3a	20	4d	6f	·User-Agent: Mozilla/5.0 (X11; Ubuntu; Linux x8		
00000290:	7a	69	6c	6c	61	21	35	2e	30	20	28	58	31	31	3b	20	6_64; rv:81.0) Gecko/20100101 Firefox/81.0 · Acce		
000002a0:	55	62	75	6e	74	75	3b	20	4d	69	6e	75	78	20	78	38	pt: text/html,application/xhtml+xml,application/xml;q=0.9,image/webp,*/*;q=0.8 · Accept-Language: en-GB,en;q=0.5 ·Accept-Encoding: gzip, deflate ·DNT: 1 ·Connect ion: keep-alive ·Upgrade-Insecure-Requests: 1 ·If-Modified-Since: Fri, 11 Sep 2020 08:42:27 GMT ·If-None-Match: W/"5ca-5af05aa276499" ·Cache-Control: max-age=0 ·-----_-----B ·B		
000002b0:	36	51	36	34	3b	20	72	76	3a	38	31	2e	30	29	20	47	·User-Agent: Mozilla/5.0 (X11; Ubuntu; Linux x8_64; rv:81.0) Gecko/20100101 Firefox/81.0 · Acce		
000002c0:	65	63	6b	6f	21	32	30	31	30	30	31	30	31	20	46	69	·User-Agent: Mozilla/5.0 (X11; Ubuntu; Linux x8_6_64; rv:81.0) Gecko/20100101 Firefox/81.0 · Acce		
000002d0:	72	65	66	6f	78	2f	38	31	2e	30	0d	0a	41	63	63	65	·User-Agent: Mozilla/5.0 (X11; Ubuntu; Linux x8_6_64; rv:81.0) Gecko/20100101 Firefox/81.0 · Acce		
000002e0:	70	74	3a	20	74	65	78	74	2f	68	74	6d	6c	2c	61	70	·User-Agent: Mozilla/5.0 (X11; Ubuntu; Linux x8_6_64; rv:81.0) Gecko/20100101 Firefox/81.0 · Acce		
000002f0:	70	6c	69	63	61	74	69	6f	6e	21	78	68	74	6d	6c	2b	·User-Agent: Mozilla/5.0 (X11; Ubuntu; Linux x8_6_64; rv:81.0) Gecko/20100101 Firefox/81.0 · Acce		
00000300:	78	6d	6c	2c	61	70	70	6c	69	63	01	74	69	6f	6e	21	·User-Agent: Mozilla/5.0 (X11; Ubuntu; Linux x8_6_64; rv:81.0) Gecko/20100101 Firefox/81.0 · Acce		
00000310:	78	6d	6c	3b	71	50	30	2e	39	2c	69	6d	61	67	65	20	·User-Agent: Mozilla/5.0 (X11; Ubuntu; Linux x8_6_64; rv:81.0) Gecko/20100101 Firefox/81.0 · Acce		
00000320:	77	65	62	70	2c	2a	2f	2a	3b	71	30	30	2e	38	0d	02	·User-Agent: Mozilla/5.0 (X11; Ubuntu; Linux x8_6_64; rv:81.0) Gecko/20100101 Firefox/81.0 · Acce		
00000330:	4	63	63	65	70	74	2d	4c	61	6e	67	75	61	67	65	3a	·User-Agent: Mozilla/5.0 (X11; Ubuntu; Linux x8_6_64; rv:81.0) Gecko/20100101 Firefox/81.0 · Acce		
00000340:	20	65	6e	2d	47	42	2c	65	6e	3b	71	30	30	2e	35	0d	06	·User-Agent: Mozilla/5.0 (X11; Ubuntu; Linux x8_6_64; rv:81.0) Gecko/20100101 Firefox/81.0 · Acce	
00000350:	0a	4	63	63	63	65	70	74	2d	45	6e	63	6f	64	69	6e	67	·User-Agent: Mozilla/5.0 (X11; Ubuntu; Linux x8_6_64; rv:81.0) Gecko/20100101 Firefox/81.0 · Acce	
00000360:	3a	20	67	7a	69	70	2c	20	64	65	66	6c	61	74	65	0d	06	·User-Agent: Mozilla/5.0 (X11; Ubuntu; Linux x8_6_64; rv:81.0) Gecko/20100101 Firefox/81.0 · Acce	
00000370:	0a	44	4e	54	3a	20	31	0d	0a	43	6f	6e	6e	65	63	74	0d	·User-Agent: Mozilla/5.0 (X11; Ubuntu; Linux x8_6_64; rv:81.0) Gecko/20100101 Firefox/81.0 · Acce	
00000380:	69	6f	6e	3a	20	6b	65	65	70	2d	01	6c	69	76	65	0d	06	·User-Agent: Mozilla/5.0 (X11; Ubuntu; Linux x8_6_64; rv:81.0) Gecko/20100101 Firefox/81.0 · Acce	
00000390:	0a	55	70	67	72	61	64	65	2d	49	6e	73	63	63	75	72	0d	06	·User-Agent: Mozilla/5.0 (X11; Ubuntu; Linux x8_6_64; rv:81.0) Gecko/20100101 Firefox/81.0 · Acce
000003a0:	65	2d	52	65	71	75	65	73	74	73	3a	20	31	0d	0a	49	0d	06	·User-Agent: Mozilla/5.0 (X11; Ubuntu; Linux x8_6_64; rv:81.0) Gecko/20100101 Firefox/81.0 · Acce
000003b0:	66	2d	4d	6f	64	69	66	69	65	64	2d	53	69	6e	63	65	0d	06	·User-Agent: Mozilla/5.0 (X11; Ubuntu; Linux x8_6_64; rv:81.0) Gecko/20100101 Firefox/81.0 · Acce
000003c0:	3a	20	46	72	69	2c	20	31	31	20	53	65	70	20	32	30	0d	06	·User-Agent: Mozilla/5.0 (X11; Ubuntu; Linux x8_6_64; rv:81.0) Gecko/20100101 Firefox/81.0 · Acce
000003d0:	32	30	20	30	38	3a	34	32	3a	32	37	20	47	4d	54	0d	06	·User-Agent: Mozilla/5.0 (X11; Ubuntu; Linux x8_6_64; rv:81.0) Gecko/20100101 Firefox/81.0 · Acce	
000003e0:	0a	49	60	2d	4e	6f	6e	65	2d	4d	61	74	63	68	3a	20	0d	06	·If-None-Match: W/"5ca-5af05aa276499" ·Cache-Control: max-age=0 ·-----_-----B ·B
000003f0:	57	2f	22	35	63	61	2d	35	61	66	30	35	01	61	32	37	0d	06	·If-None-Match: W/"5ca-5af05aa276499" ·Cache-Control: max-age=0 ·-----_-----B ·B
00000400:	36	34	39	39	22	0d	0a	43	61	63	68	65	2d	43	6f	6e	0d	06	·If-None-Match: W/"5ca-5af05aa276499" ·Cache-Control: max-age=0 ·-----_-----B ·B
00000410:	74	72	6f	6c	3a	20	6d	61	78	2d	61	67	65	30	30	0d	06	·If-None-Match: W/"5ca-5af05aa276499" ·Cache-Control: max-age=0 ·-----_-----B ·B	
00000420:	0a	0d	07	cf	7f	90	5f	10	1b	08	00	41	00	00	00	00	0d	06	·If-None-Match: W/"5ca-5af05aa276499" ·Cache-Control: max-age=0 ·-----_-----B ·B

## Network Data

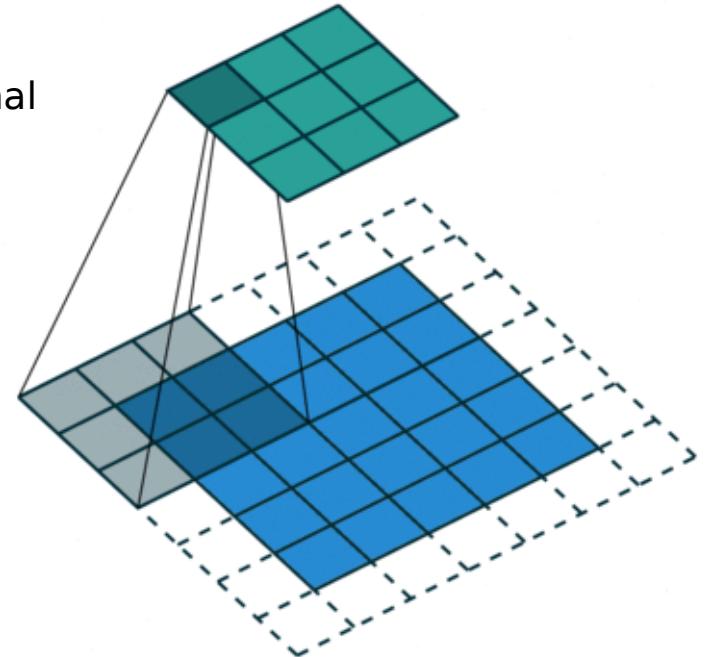
- Series of *packets*
- Weak order of packets
- Packets consist of *layers*
- Different types of protocols
- Layers have *header* and *payload* fields

## Industrial Control Systems (ICS)

- Monitoring & Control functions
- Fixed network topology
- Proprietary Infrastructure
- Cycle information exchange

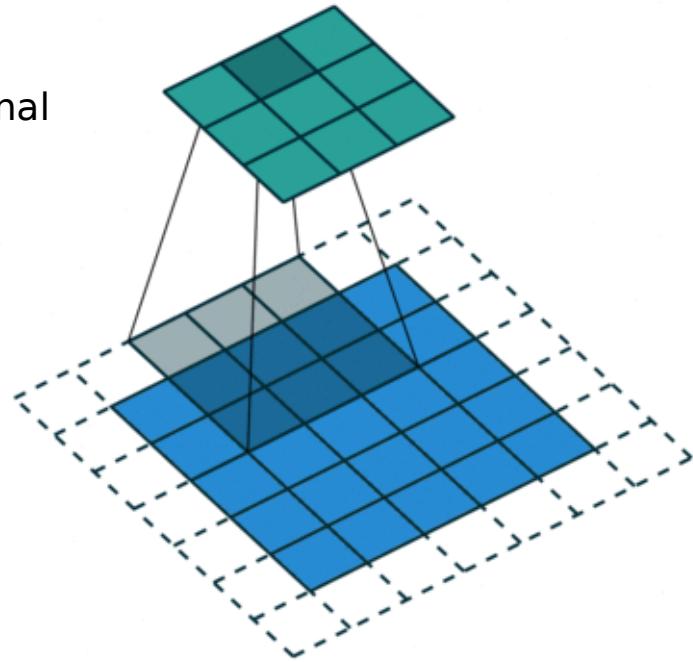
00000260:	54	50	2f	31	2e	31	0d	0a	48	6f	73	74	3a	20	78	6e	TP/1.1 · Host: xn	
00000270:	2d	2d	6d	62	69	75	73	2d	6a	75	61	20	62	61	6e	64	--mbius-jua.band	
00000280:	0d	07	55	73	65	72	2d	41	67	65	6e	74	3a	20	4d	6f	···User-Agent: Mo	
00000290:	7a	69	6c	6c	61	21	35	2e	30	20	28	58	31	31	3b	20	zilla/5.0 (X11;	
000002a0:	55	62	75	6e	74	75	3b	20	4d	69	6e	75	78	20	78	38	Ubuntu; Linux x8	
000002b0:	36	51	36	34	3b	20	72	76	3a	38	31	2e	30	20	20	47	6_64; rv:81.0) G	
000002c0:	65	63	6b	6f	21	32	30	31	30	30	31	30	31	20	46	69	ecko/20100101 Fi	
000002d0:	72	65	66	6f	78	2f	38	31	2e	30	0d	0a	41	63	63	65	refox/81.0 · Acce	
000002e0:	70	74	3a	20	74	65	78	74	2f	68	74	6d	6c	2c	61	70	pt: text/html,ap	
000002f0:	70	6c	69	63	61	74	69	6f	6e	21	78	68	74	6d	6c	2b	plication/xhtml+xml,application/xml+xml;q=0.9,image/webp,*/*;q=0.8 ·	
00000300:	78	6d	6c	2c	61	70	70	6c	69	63	01	74	69	6f	6e	21	Accept-Language: en-GB,en;q=0.5 ·	
00000310:	78	6d	6c	3b	71	50	30	2e	39	2c	69	6d	61	67	65	20	·Accept-Encoding: gzip, deflate ·	
00000320:	77	65	62	70	2c	2a	2f	2a	3b	71	3d	30	2e	38	0d	02	·DNT: 1 · Connect ion: keep-alive ·	
00000330:	4	63	63	65	70	74	2d	4c	61	6e	67	75	61	67	65	3a	Upgrade-Insecur e-Requests: 1 · If-Modified-Since : Fri, 11 Sep 20	
00000340:	20	65	6e	2d	47	42	2c	65	6e	3b	71	3d	30	2e	35	0d	20 08:42:27 GMT ·	
00000350:	0a	4	63	63	63	65	70	74	2d	45	6e	63	6f	64	69	6e	67	·If-None-Match: W/"5ca-5af05aa27
00000360:	3a	20	67	7a	69	70	2c	20	64	65	66	6c	61	74	65	0d	6499 · Cache-Con trol: max-age=0 ·	
00000370:	0a	44	4e	54	3a	20	31	0d	0a	43	6f	6e	66	65	63	74	....._.....B · B	
00000380:	69	6f	6e	3a	20	6b	65	65	70	2d	61	6c	69	76	65	0d		
00000390:	0a	55	70	67	72	61	64	65	2d	49	6e	73	63	63	75	72		
000003a0:	65	2d	52	65	71	75	65	73	74	73	3a	20	31	0d	0a	49		
000003b0:	66	2d	4d	6f	64	69	66	69	65	64	2d	53	69	6e	63	65		
000003c0:	3a	20	46	72	69	2c	20	31	31	20	53	65	70	20	32	30		
000003d0:	32	30	20	30	38	3a	34	32	3a	32	37	20	47	4d	54	0d		
000003e0:	0a	49	60	2d	4e	6f	6e	65	2d	4d	61	74	63	68	3a	20		
000003f0:	57	2f	22	35	63	61	2d	35	61	66	30	35	61	62	32	37		
00000400:	36	34	39	39	22	0d	0a	43	61	63	68	65	2d	43	6f	6e		
00000410:	74	72	6f	6c	3a	20	6d	61	78	2d	61	67	65	30	30	0d		
00000420:	0a	0d	07	cf	7f	90	5f	10	1b	08	00	41	00	00	00	00		

- Neurons that *convolve (dot product)* a filter with the input signal
- Filters/Kernels are adjustable through learning
- Hierarchical layers of neurons learn spatial patterns
- Often used for images based signal processing



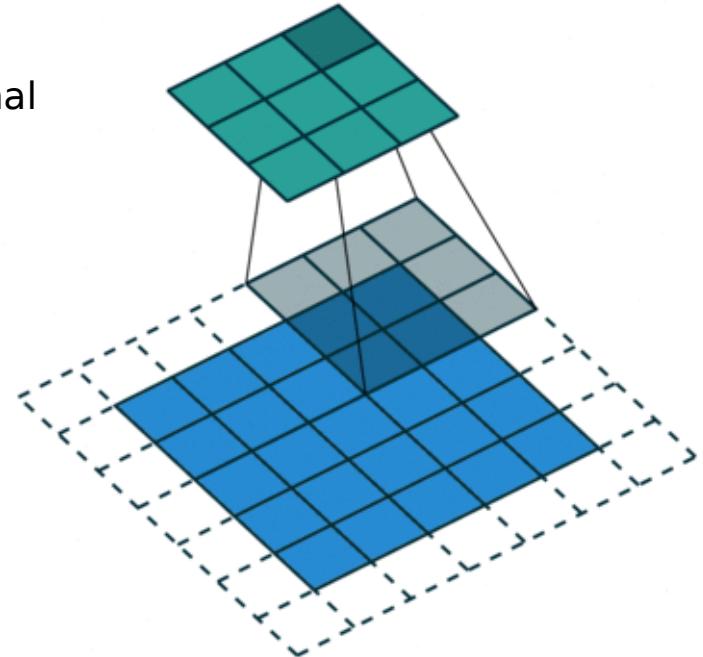
[github.com/vdumoulin/conv\\_arithmetic](https://github.com/vdumoulin/conv_arithmetic)

- Neurons that *convolve (dot product)* a filter with the input signal
- Filters/Kernels are adjustable through learning
- Hierarchical layers of neurons learn spatial patterns
- Often used for images based signal processing



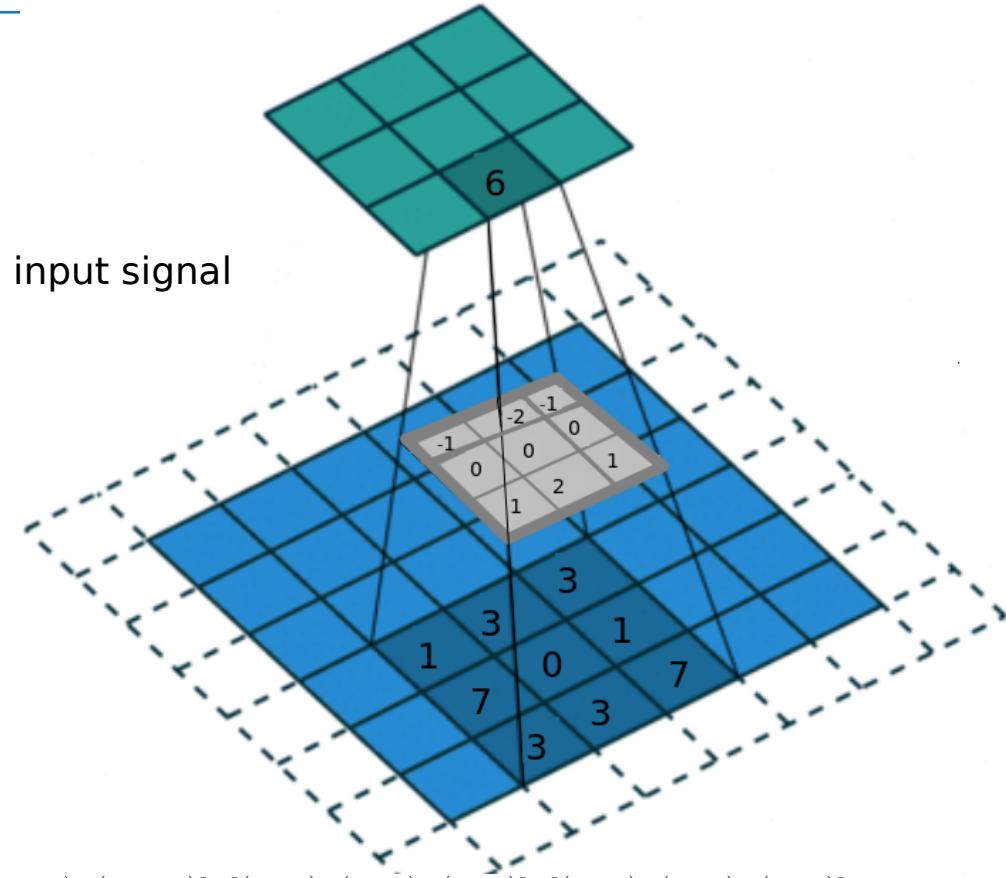
[github.com/vdumoulin/conv\\_arithmetic](https://github.com/vdumoulin/conv_arithmetic)

- Neurons that *convolve (dot product)* a filter with the input signal
- Filters/Kernels are adjustable through learning
- Hierarchical layers of neurons learn spatial patterns
- Often used for images based signal processing



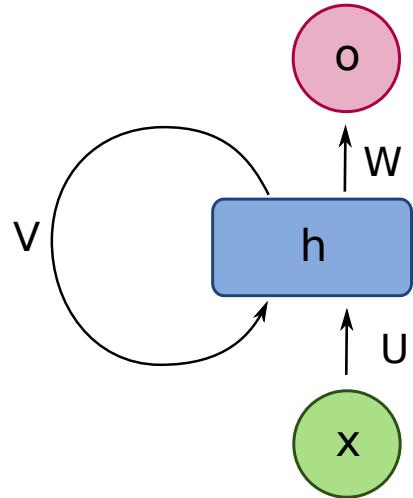
[github.com/vdumoulin/conv\\_arithmetic](https://github.com/vdumoulin/conv_arithmetic)

- Neurons that *convolve (dot product)* a filter with the input signal
- Filters/Kernels are adjustable through learning
- Hierarchical layers of neurons learn spatial patterns
- Often used for images based signal processing

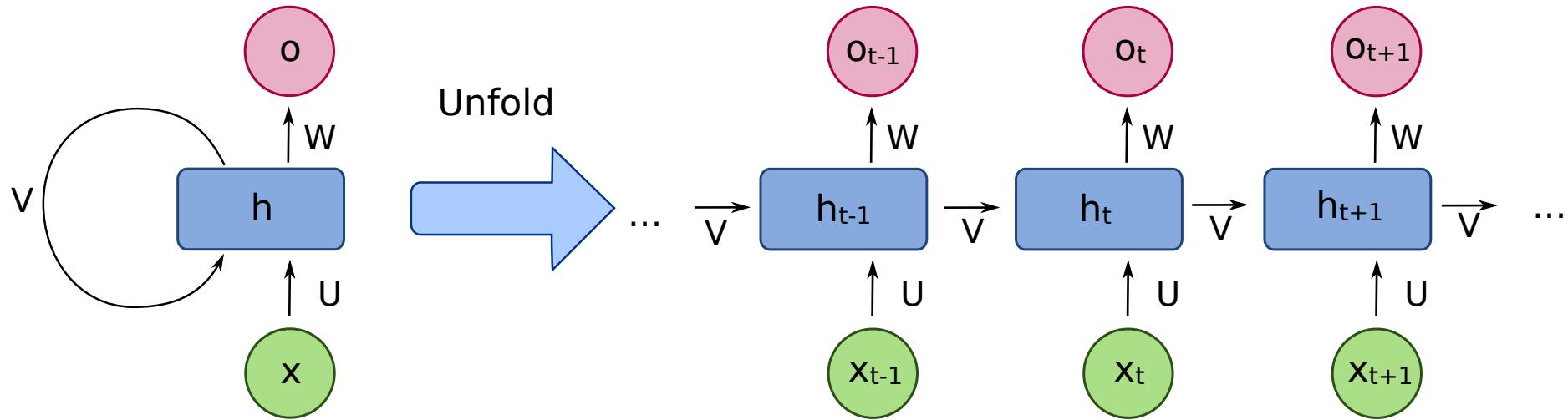


$$[(1 * -1) + (3 * -2) + (3 * -1)] + [(7 * 0) + (0 * 0) + (1 * 0)] + [(3 * 1) + (3 * 2) + (7 * 1)] = 6$$

[github.com/vdumoulin/conv\\_arithmetic](https://github.com/vdumoulin/conv_arithmetic)

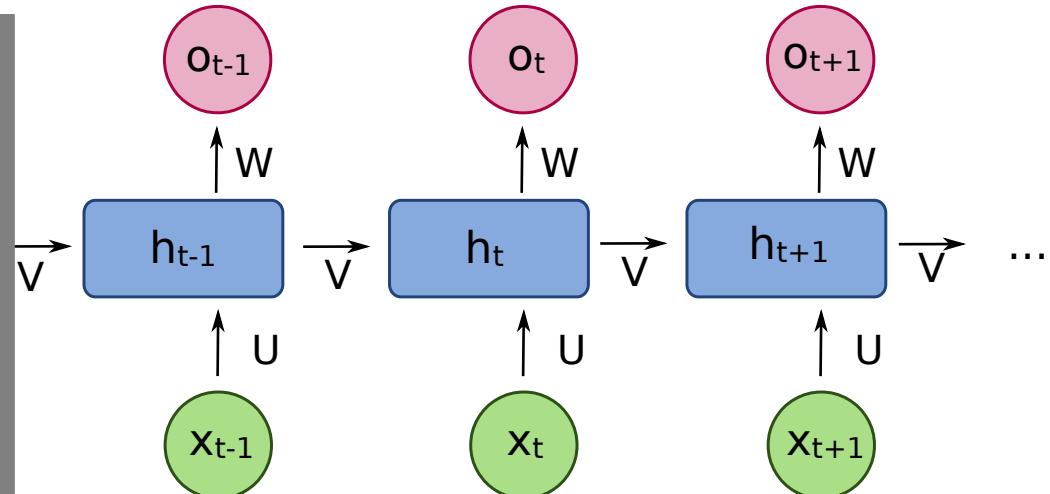


unchanged, Licence: <https://creativecommons.org/licenses/by-sa/4.0/deed.en>



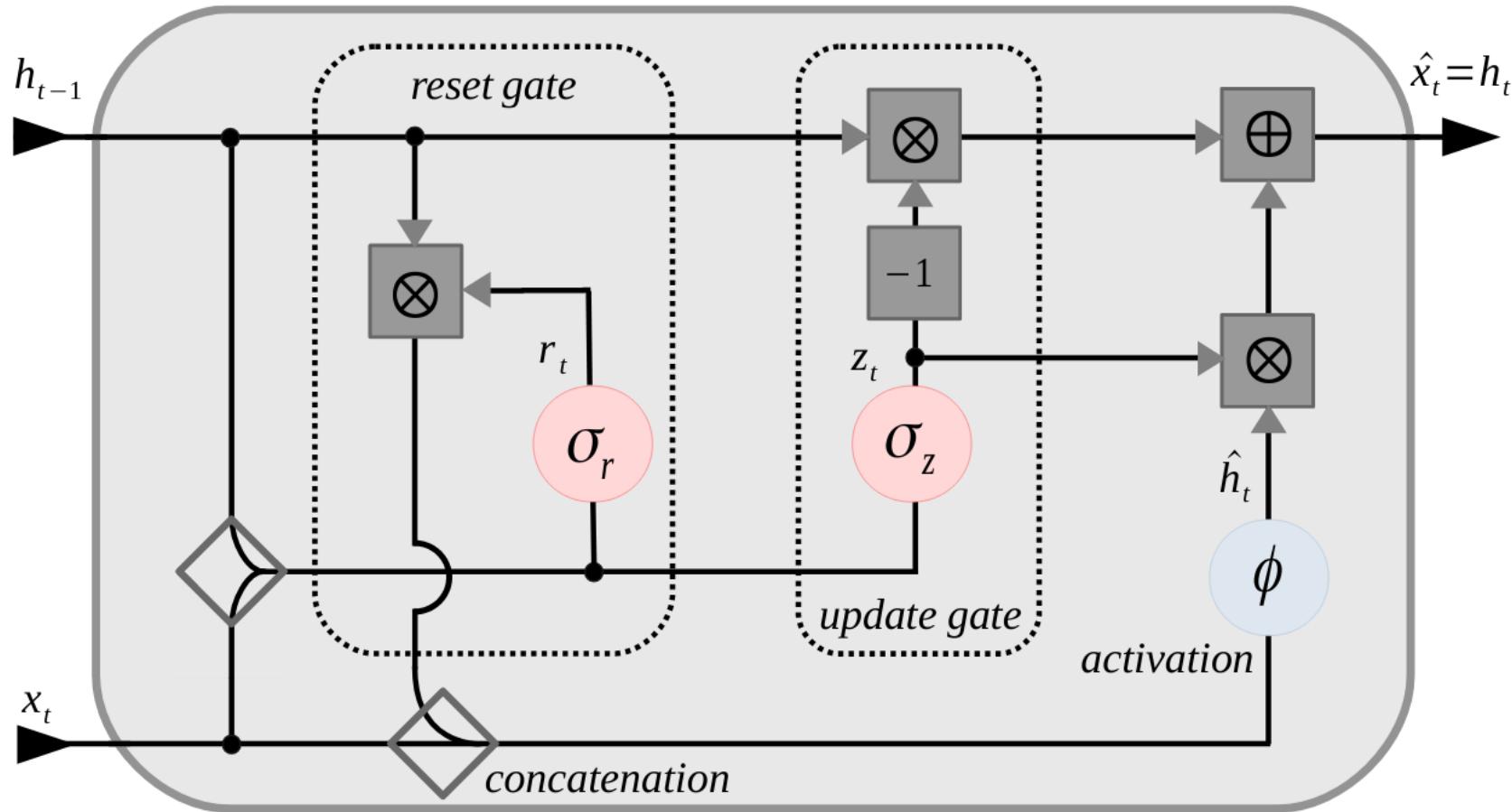
unchanged, Licence: <https://creativecommons.org/licenses/by-sa/4.0/deed.en>

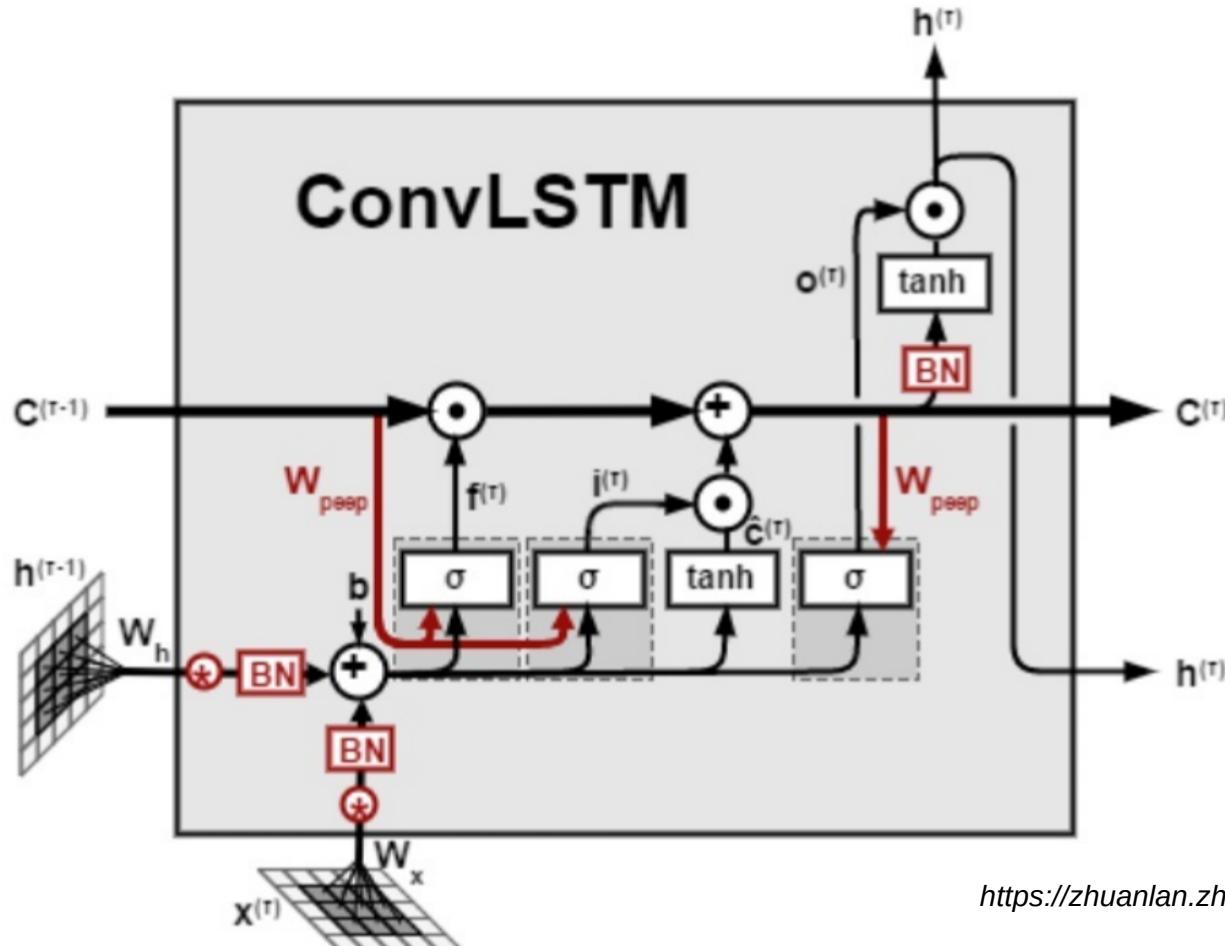
- Incorporate historical data
- Output signal is used for the next step
- Designed for sequential problems
- Used in generative models
- Complex units can host a number of gates



unchanged, Licence: <https://creativecommons.org/licenses/by-sa/4.0/deed.en>

## A | Appendix: Gated Rectifier Unit





<https://zhuanlan.zhihu.com/p/148122328>

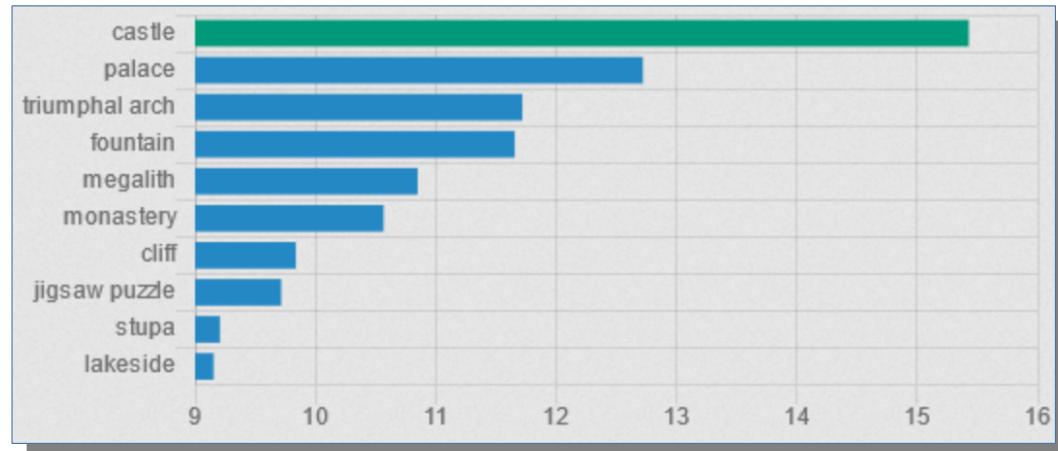


© 2021 Ulf Büschleb

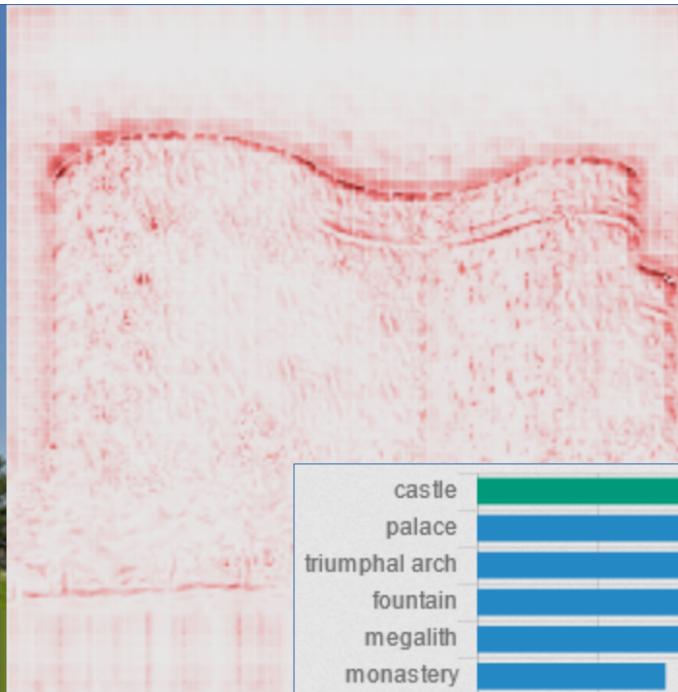
## A | Appendix: LPR - Example



© 2021 Ulf Büschleb

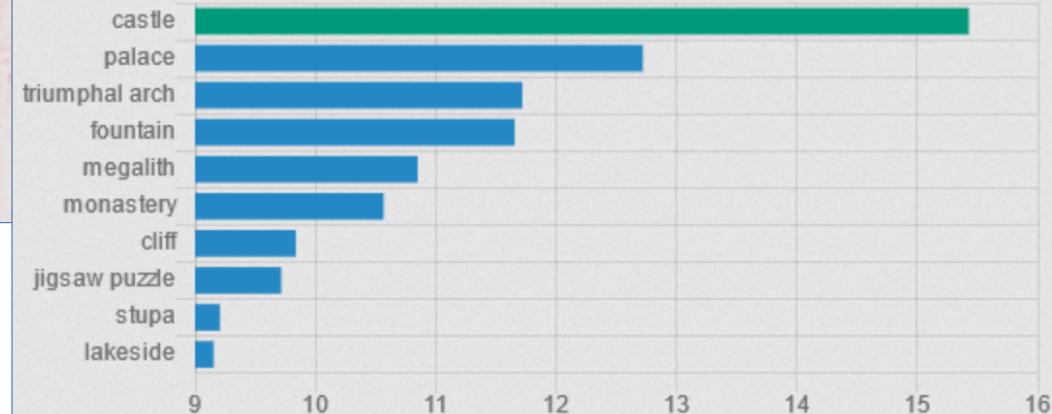


# A | Appendix: LPR - Example



lprserver.hhi.fraunhofer.de

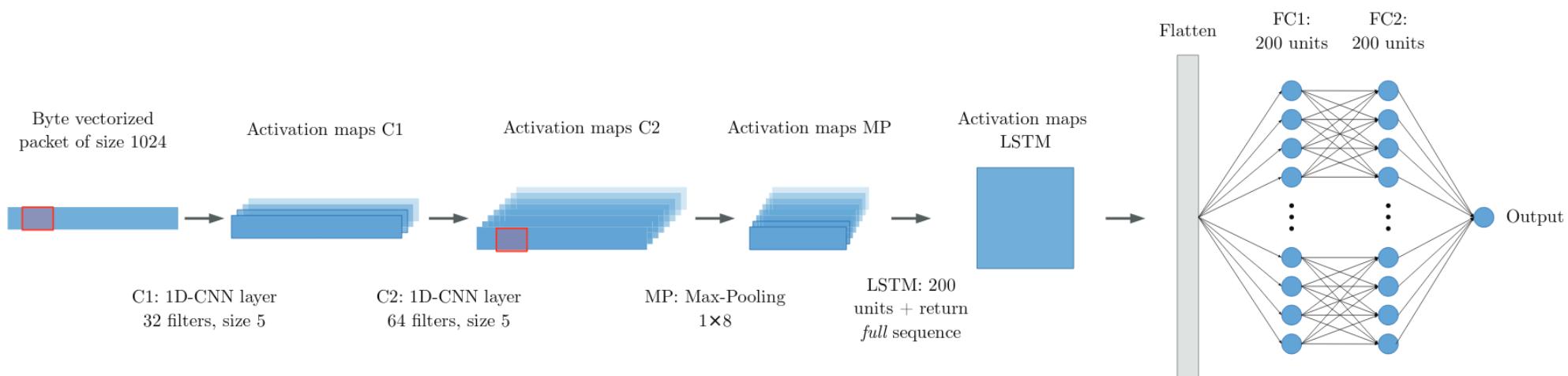
© 2021 Ulf Büschleb



## DeepMAL - Deep Learning Models for Malware Traffic Detection and Classification

(P. Casas et al. - iDSC 2020)

- Feature extraction on first 1000 bytes of every packet / flow
- Spatial-Temporal representation learning on raw traffic (+)
- Supervised softmax-based classification (-)
- Comparison to shallow methods & handcrafted features extraction
- Evaluated on ICS unrelated data



## An Unsupervised Deep Learning Model for Early Network Traffic Anomaly Detection

- *IEEE Access*, 2020

## An LSTM-Based Deep Learning Approach for Classifying Malicious Traffic at the Packet Level

- *Applied Sciences*, 2019

## ZOE: Content-based Anomaly Detection for Industrial Control Systems

- *IEEE/IFIP DSN*, 2018

## A Deep Learning Approach for Network Intrusion Detection System

- *ACM BIONETICS*, 2016

## Deep packet: a novel approach for encrypted traffic classification using deep learning

- *Soft Computing*, 2020

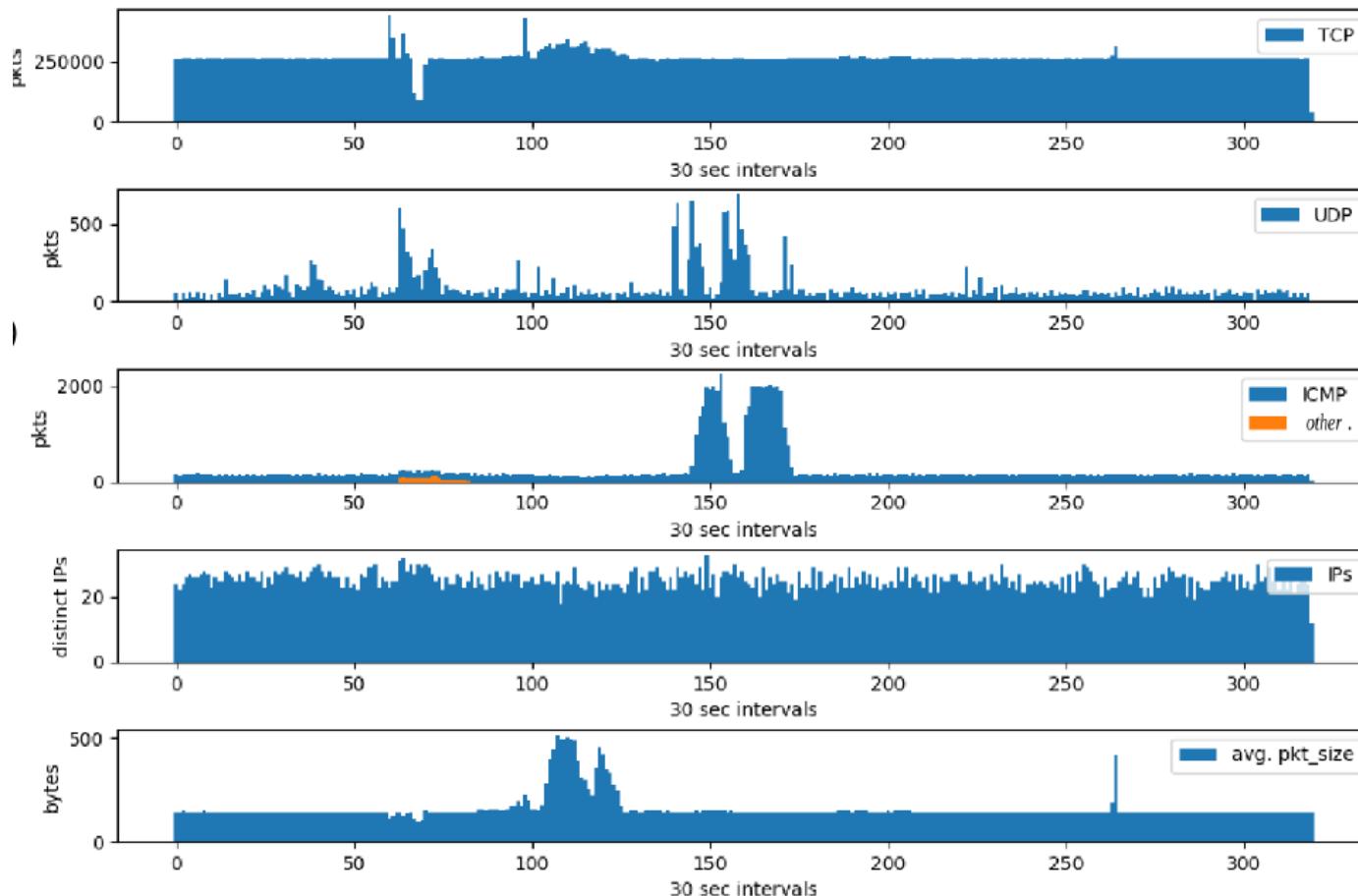
## HAST-IDS: Learning Hierarchical Spatial-Temporal Features Using Deep Neural Networks to Improve Intrusion Detection

- *IEEE ACCESS*, 2017

## Malware Traffic Classification Using Convolutional Neural Network for Representation Learning

- *IEEE ICOIN*, 2017

# A | Appendix: V Data Set Characteristic



# A | Appendix: Dimension Reduction Comparison

	mean	median	min	max	std	IQR
PCA	$-1.26 \cdot 10^{-7}$	<u>0.4823</u>	-4.0832	7.1323	0.3601	0.4690
k-PCA	$-1.85 \cdot 10^{-8}$	-0.0003	-4.0831	7.1324	0.3670	<u>0.4781</u>
ICA	$-1.51 \cdot 10^{-9}$	$-2.60 \cdot 10^{-5}$	<u>-0.0539</u>	0.0661	0.00617	0.008
Isomap	<i>n/a</i>	<i>n/a</i>	<i>n/a</i>	<i>n/a</i>	<i>n/a</i>	<i>n/a</i>
LLE	$1.66 \cdot 10^{-7}$	$-2.74 \cdot 10^{-5}$	-0.3514	0.3098	0.0061	0.0025
SE	<i>n/a</i>	<i>n/a</i>	<i>n/a</i>	<i>n/a</i>	<i>n/a</i>	<i>n/a</i>
AE	<u>0.0162</u>	0.0016	-0.833	<u>0.7694</u>	<u>0.341</u>	0.4591
<i>original</i>	0.4844	0.4823	0	1	0.2924	0.5098

	non-linearity	OOS processing	ability to retrain	space complexity	time complexity
PCA		<b>x</b>		$\mathcal{O}(N * D)$	$\mathcal{O}(N * D + D^3)$
k-PCA	<b>x</b>	<b>x</b>		$\mathcal{O}(N^2)$	$\mathcal{O}(N^3)$
ICA	<b>x</b>	<b>x</b>		$\mathcal{O}(N^2)$	$\mathcal{O}(D^2 * N)$
Isomap	<b>x</b>	<b>x</b>		$\mathcal{O}(d * N^2)$	$\mathcal{O}(N^2)$
LLE	<b>x</b>			$\mathcal{O}(N^2)$	$\mathcal{O}(d * N^2)$
SE	<b>x</b>			$\mathcal{O}(N^2)$	$\mathcal{O}(N^3)$
AE	<b>x</b>	<b>x</b>	<b>x</b>	$\mathcal{O}(2 * \theta)$	$\mathcal{O}(\text{epoch} * D * N)$

Table 5.1: Comparison of different dimension reduction methods

	<i>Train</i>	<i>Validation</i>	<i>Eval<sub>H</sub></i>	<i>Eval<sub>D</sub></i>	<i>Eval<sub>E</sub></i>	<i>Eval<sub>R</sub></i>
Duration [sec]	1139.7	162	2907	702	702	702
Size [GB]	1.7	0.24	0.72	1.2	1.2	1.2
# Frames	10,774,388	1,539,556	1,518,761	3,521,649	6,157,648	5,232,729
# Anomalies (%)	0	0	424,364 (27.9)	135,048 (3.8)	2,460 (0.04)	1,930 (0.037)
# MACs	39	33	110	37	57	54
# IP	42	34	34	38	40	40
# Protocols	25	19	21	20	22	22
# Flows [tcp/udp]	993/417	144/51	6,648/39	255/95	448/161	377/136
ATU (std)	146.3 (166)	145.6 (158)	477.8 (1228)	144.6 (12)	146.3(0.5)	146.3(0.5)
kpackets/sec	0.4	9.5	0.5	5.0	8.7	7.5

Table 4.1: V [redacted] traffic statistics per sub set

# A | Appendix: Data Characteristics

	<i>Train<sub>RL</sub></i>	<i>Validation<sub>RL</sub></i>	<i>Train<sub>AD</sub></i>	<i>Validation<sub>AD</sub></i>
Duration [sec]	898.6	4.1	408.5	4.1
Size [GB]	1.4	0.14	1.4	0.14
# Frames	10,000,000	100,000	10,000,000	100,000
# Anomalies [%]	0	0	0	0
# MACs	54	31	57	33
# IP	51	32	60	35
# Protocols	36	21	39	21
# Flows [tcp/udp]	87,961 / 172	971/19	87,368/236	1,572/21
ATU (std)	113.23 (105)	108.23 (104)	113.38 (105)	112.55 (105)
<i>k</i> packets/sec	11.1	71.4	24.0	24.0
	<i>Eval<sub>H</sub></i>	<i>Eval<sub>D</sub></i>	<i>Eval<sub>E</sub></i>	<i>Eval<sub>R</sub></i>
Duration [sec]	299.7	487.6	487.2	486.9
Size [GB]	1.1	1.3	1.3	1.3
# Frames	7,350,000	6,284,294	6,194,072	6,194,576
# Anomalies [%]	18,900 (0.25)	170,608 (2.7)	516 (0.008)	148 (0.002)
# MACs	56	56	67	67
# IP	52	48	48	48
# Protocols	37	34	36	33
# Flows [tcp/udp]	70,316/148	28,649/119	28,649/119	28,649/122
ATU (std)	114.44 (117)	113.65 (322)	113.12 (101)	113.41 (103)
<i>k</i> packets/sec	24.0	12.8	12.7	12.7

Table 4.2: SWaT traffic statistics per sub set

# A | Appendix: Autoencoder Hyperparameter search

RNN-Cell	Loss Function	Optimization	LR-Scheduler	Loss after Epoch 5
GRU	MSE	SGD	cycle	$2.82 \cdot 10^{-4}$
GRU	MSE	SGD	plateau	$1.72 \cdot 10^{-4}$
GRU	MSE	SGD	step	$1.72 \cdot 10^{-4}$
GRU	MSE	adamW	cycle	$2.34 \cdot 10^{-5}$
GRU	MSE	adamW	step	$6.11 \cdot 10^{-5}$
GRU	MSE	adamW	plateau	$1.26 \cdot 10^{-4}$
GRU	BCE	SGD	cycle	$1.14 \cdot 10^{-3}$
GRU	BCE	SGD	step	$2.93 \cdot 10^{-4}$
GRU	BCE	SGD	plateau	$6.34 \cdot 10^{-4}$
GRU	BCE	adamW	cycle	$3.15 \cdot 10^{-5}$
GRU	BCE	adamW	step	$6.92 \cdot 10^{-5}$
GRU	BCE	adamW	plateau	$1.71 \cdot 10^{-4}$
LSTM	MSE	SGD	cycle	$9.49 \cdot 10^{-5}$
LSTM	MSE	SGD	step	$5.81 \cdot 10^{-5}$
LSTM	MSE	SGD	plateau	$5.81 \cdot 10^{-5}$
LSTM	MSE	adamW	cycle	$2.21 \cdot 10^{-5}$
LSTM	MSE	adamW	step	$2.84 \cdot 10^{-4}$
LSTM	MSE	adamW	plateau	$7.16 \cdot 10^{-4}$
LSTM	BCE	SGD	cycle	$9.75 \cdot 10^{-5}$
LSTM	BCE	SGD	step	$6.52 \cdot 10^{-5}$
LSTM	BCE	SGD	plateau	$6.33 \cdot 10^{-5}$
LSTM	BCE	adamW	cycle	$2.39 \cdot 10^{-5}$
LSTM	BCE	adamW	step	$1.15 \cdot 10^{-3}$
LSTM	BCE	adamW	plateau	$9.35 \cdot 10^{-4}$

Table 6.1: Hyperparameter grid search results

## A | Appendix: Training Results (Histogram)

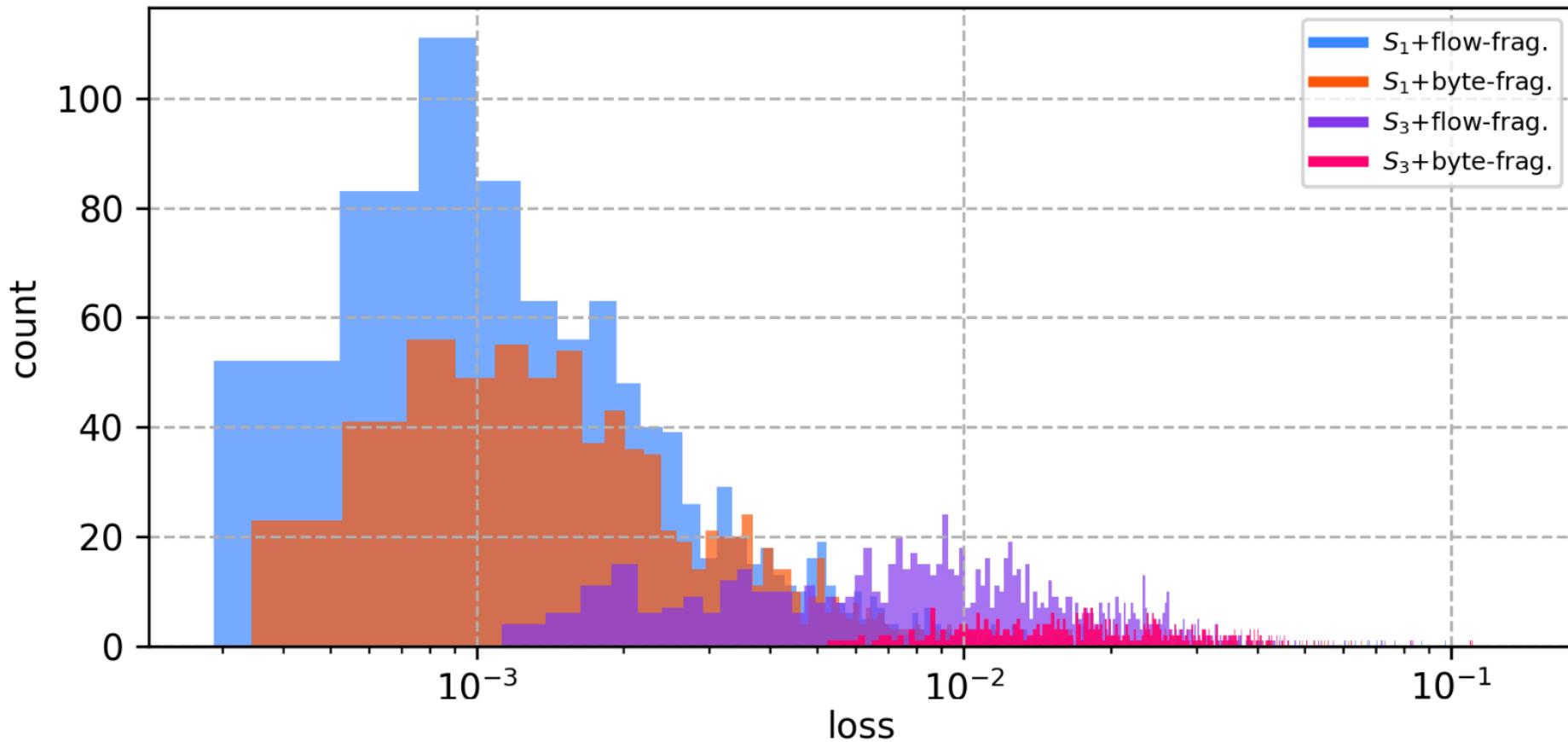


Figure 6.3: Training loss distribution for the SWaT A6 data set

## A | Appendix: Training Results (Histogram)

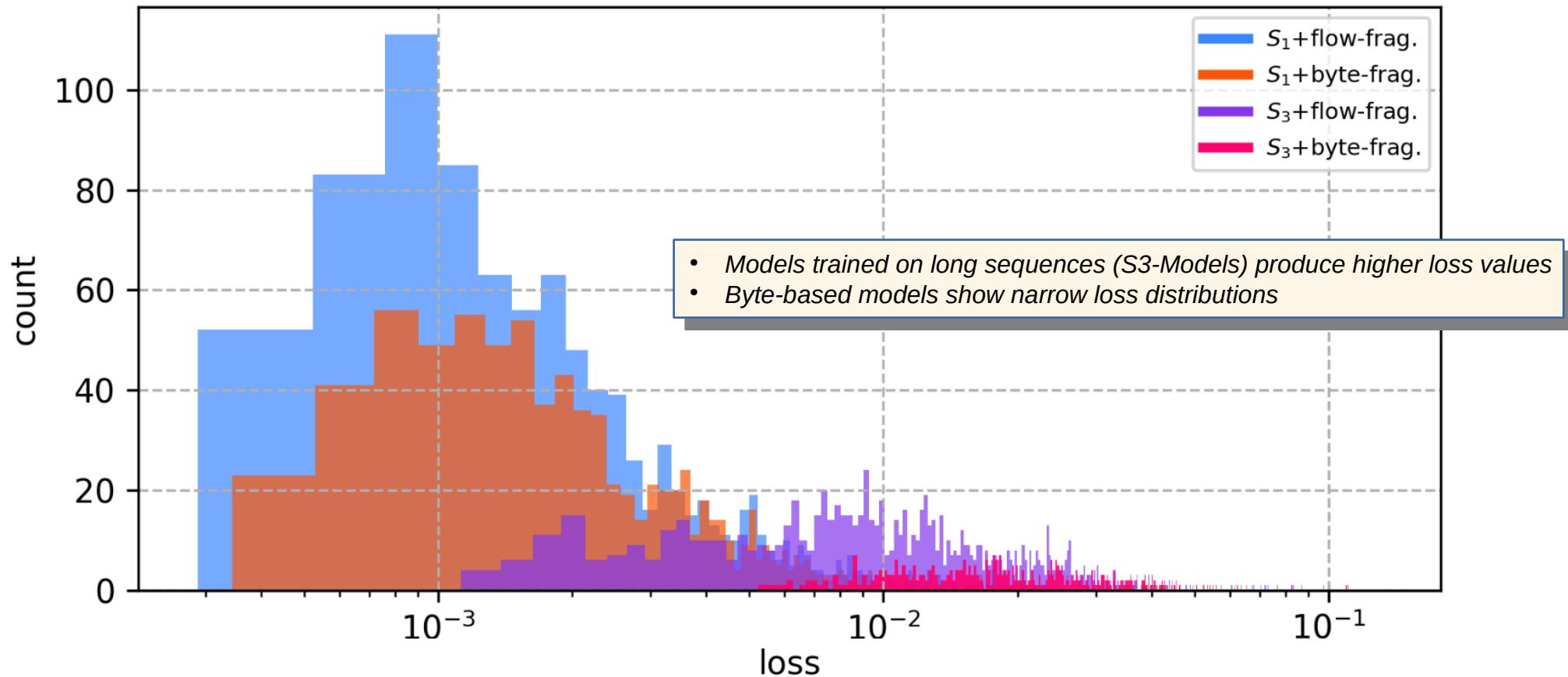


Figure 6.3: Training loss distribution for the SWaT A6 data set

# A | Appendix: SWaT Data Set Comparison

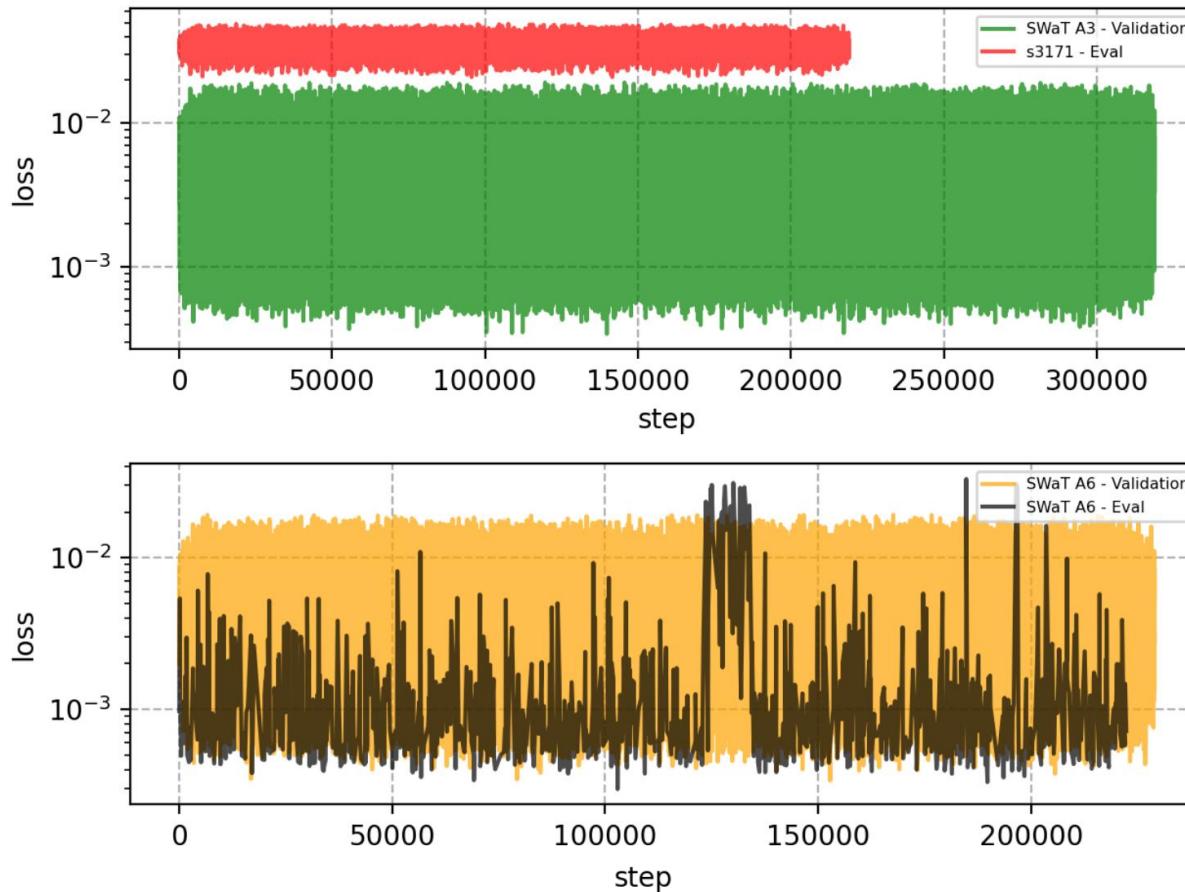
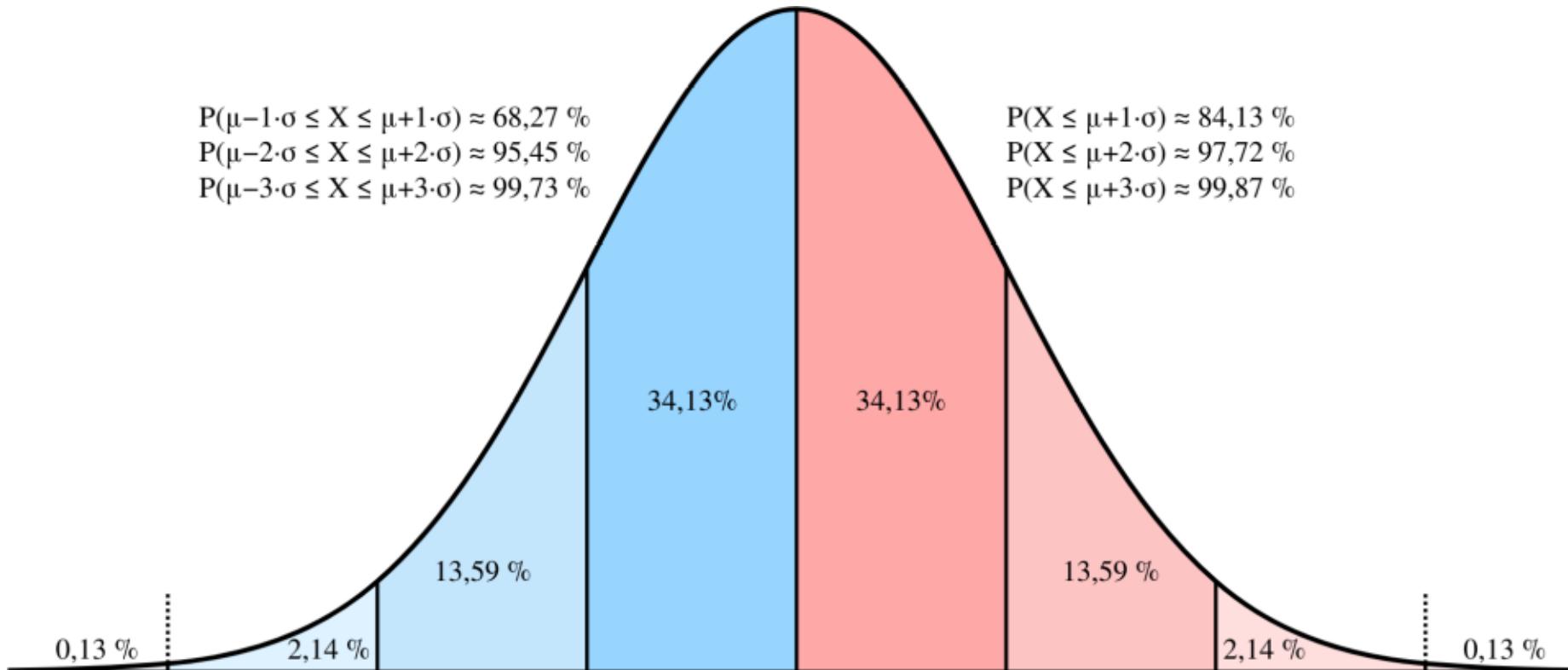
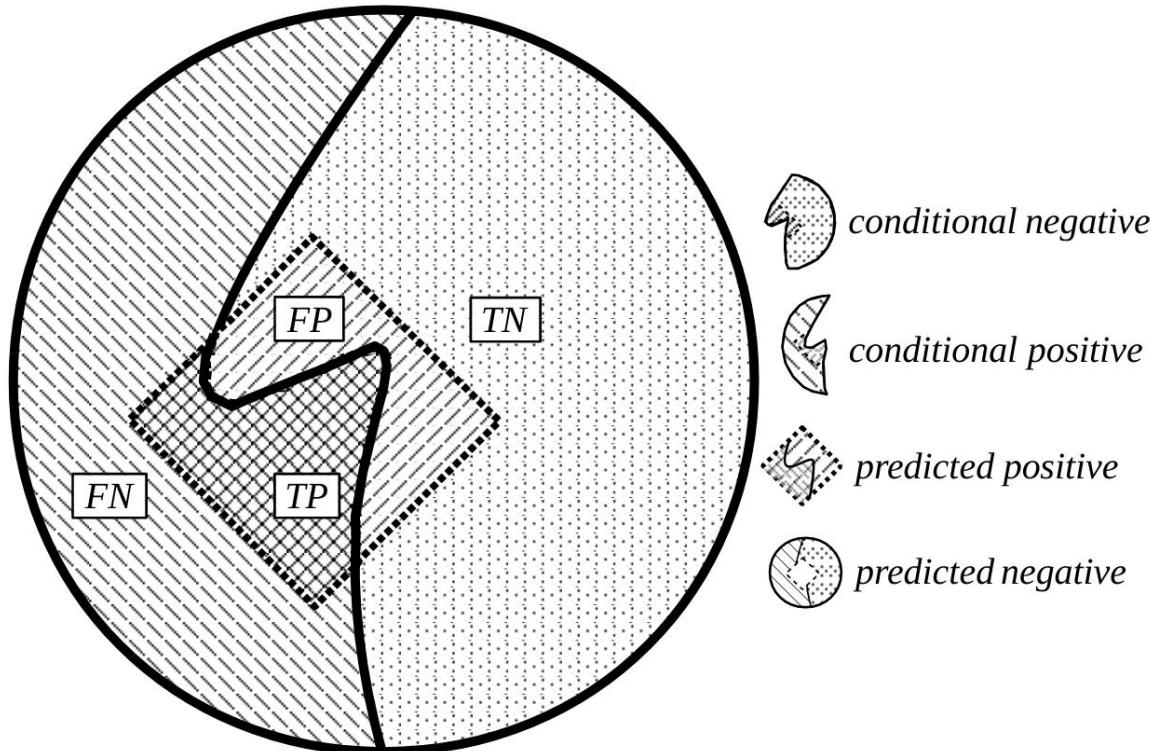


Figure 6.1: Residual loss for different iTrust data sets



## A | Appendix: Binary Classification Venn Diagram



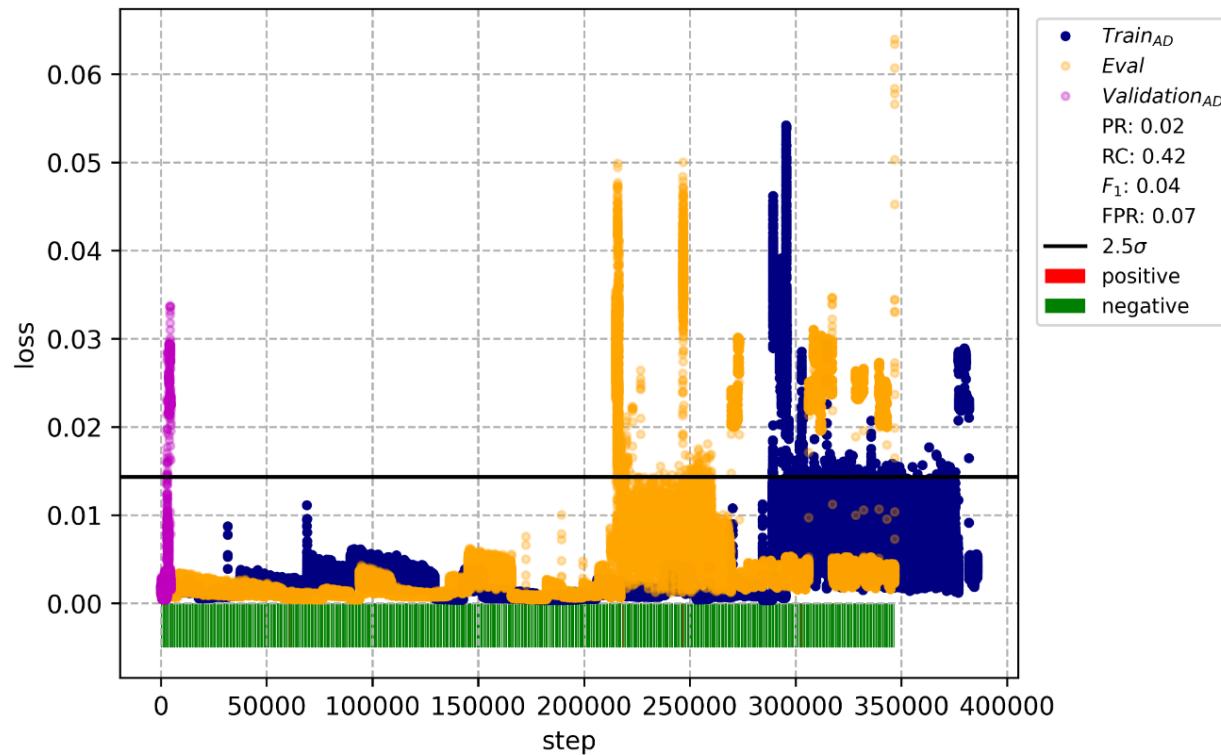


Figure B.3: SWaT A6  $Eval_H$  naive baseline loss visualisation for flow-fragments trained on an input sequence length of  $n = 1$

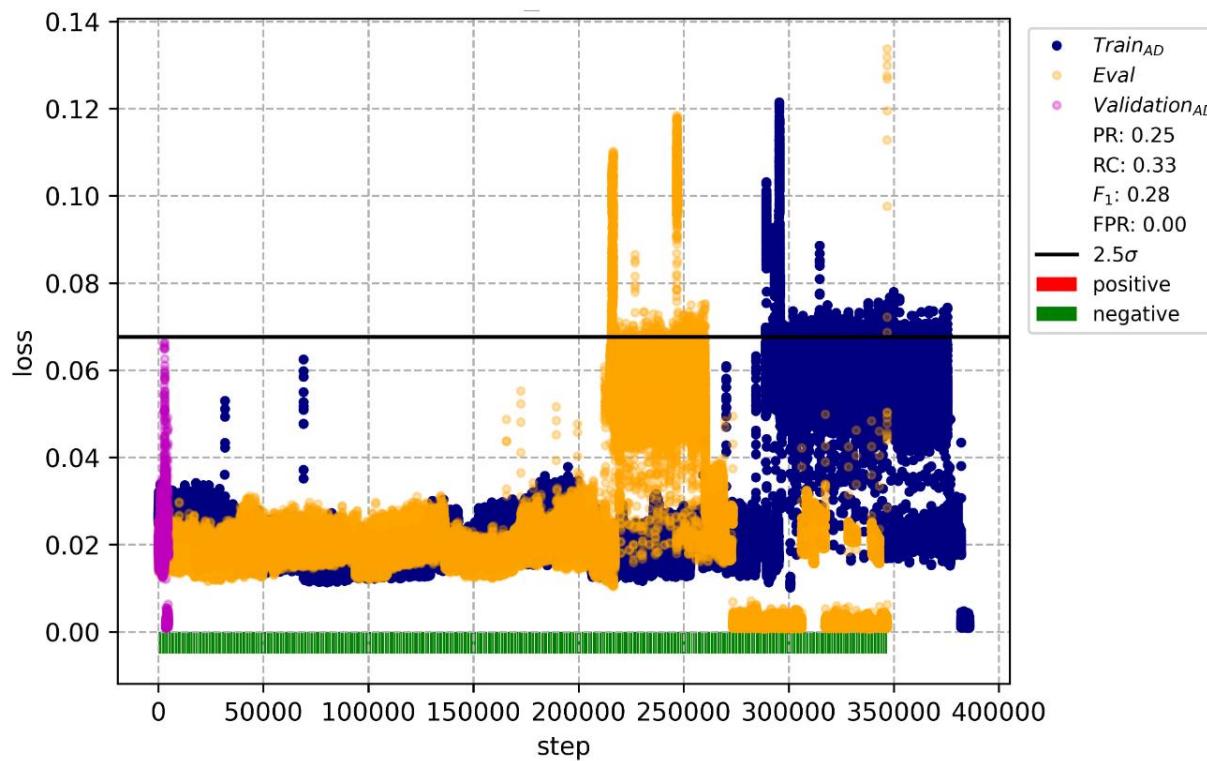


Figure B.4: SWaT A6  $Eval_H$  naive baseline loss visualisation for flow-fragments trained on an input sequence length of  $n = 3$

## A | Appendix: **Decision Analysis**

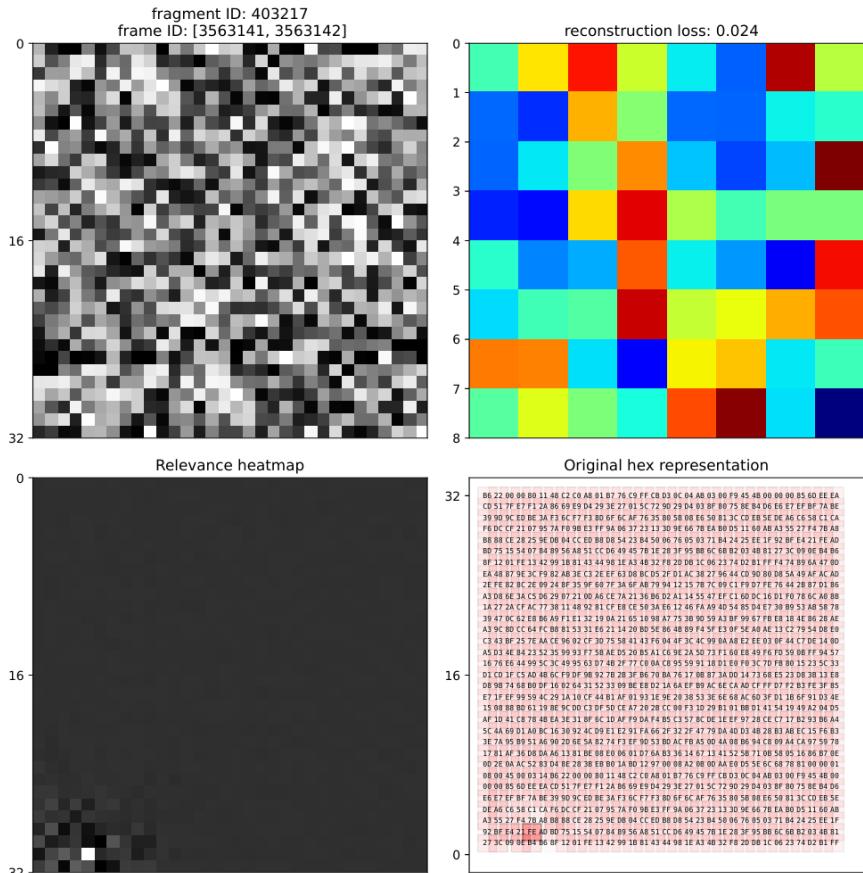


Figure B.9: Byte-fragments relevance heatmaps of the unusual UDP packet

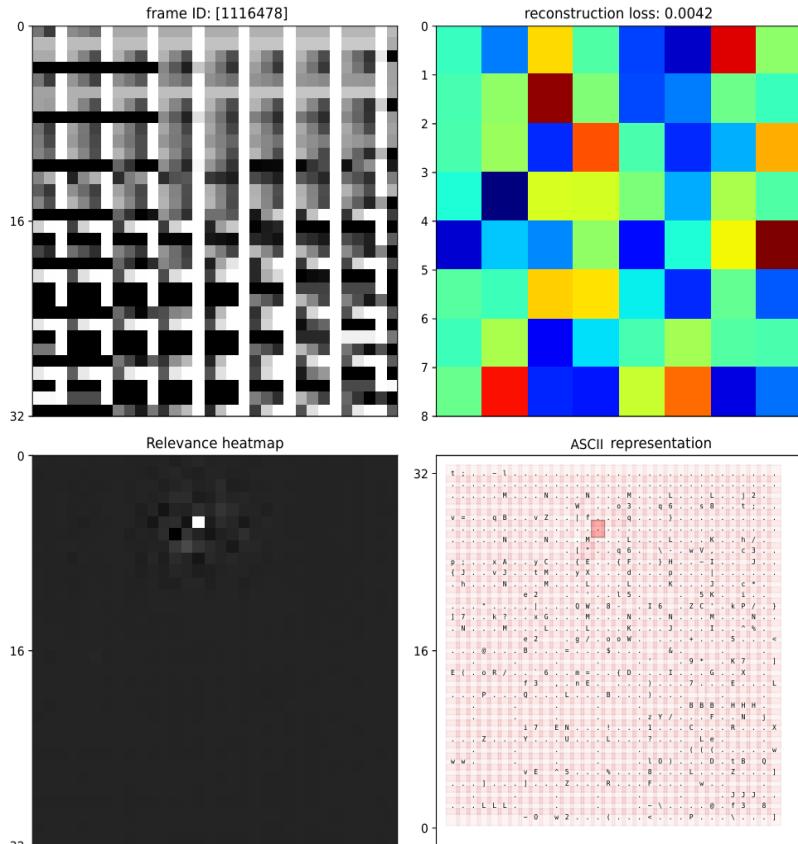


Figure 6.10: Byte-fragments relevance heatmaps of the SWaT file download event

# A | Appendix: Decision Analysis

	Time	Source	Destination	Protocol	Leng	Info
844...	3.420098	192.168.1.183	118.201.255.203	UDP	802	54028 → 1195 Len=760
84498	3.420099	192.168.1.183	118.201.255.203	UDP	806	54028 → 1195 Len=760
84499	3.420186	192.168.1.183	239.255.255.250	UDP	698	54029 → 3702 Len=656
<hr/>						
Frame 84497: 802 bytes on wire (6416 bits), 802 bytes captured (6416 bits) on interface 0						
Ethernet II, Src: Giga-Byt_6c:68:78 (e0:d5:5e:6c:68:78), Dst: AdiEngin_0b:0d:aa (00:08:a2:0b:0d:aa)						
Internet Protocol Version 4, Src: 192.168.1.183, Dst: 118.201.255.203						
User Datagram Protocol, Src Port: 54028, Dst Port: 1195						
Data (760 bytes)						
000	00 08 a2 0b 0d aa e0 d5	5e 6c 68 78 08 00 45 00				..... ^lhx..E.
010	03 14 b6 22 00 00 80 11	48 c2 c0 a8 01 b7 76 c9				...". .... H.....v.
020	ff cb d3 0c 04 ab 03 00	f9 45 4b 00 00 00 85 6d				..... .EK....m
030	ee ea cd 51 7f e7 f1 2a	86 69 e9 d4 29 3e 27 01				...Q....* .i..)>'.
040	5c 72 9d 29 d4 03 8f 80	75 8e b4 d6 e6 e7 ef bf				\r.).... u.....
050	7a be 39 9d 9c ed be 3a	f3 6c f7 f3 8d 6f 6c af				z.9....: l...ol.
060	76 35 80 5b 08 e6 50 81	3c cd eb 5e de a6 c6 58				v5.[..P. <..^...X
070	c1 ca f6 dc cf 21 07 95	7a f0 9b e3 ff 9a 06 37				.....!.. z.....7
080	23 13 3d 9e 66 7b ea b0	d5 11 60 ab a3 55 27 f4				#.=.f{.. . `..U'.
090	7b a8 b8 88 ce 28 25 9e	db 04 cc ed b8 d8 54 23				{.....(%.....T#
0a0	b4 50 06 76 05 03 71 b4	24 25 ee 1f 92 bf e4 21				.P.v..q. \$%.....!

Figure 6.9: Wireshark excerpt from a UDP packet with an unusual payload