

## A Preprocessing Data

We provide Python code for converting a flow-level dataset into a packet-level one. We also provide code to convert both the CID-IDS-2017 and UNSW-NB15 flows into a unified data structure for processing into a general purpose data transformer. Our code is also freely available.<sup>9</sup>

### A.1 Converting Flow-Level Data in Dictionaries

Flow-level datasets are typically comprised of one or more spreadsheets with columns associated with IP addresses, port numbers, protocols, timestamps, durations, and labels, among other things. This function produces a dictionary of these flows with which to pair PCAP packets. The following function parses these datasets, extracts all relevant flow information, and produces dictionaries where the keys are of the form:

$$\begin{aligned} & (\text{src\_ip}, \text{src\_port}, \text{dest\_ip}, \text{dest\_port}, \text{protocol}) \\ & (\text{dest\_ip}, \text{dest\_port}, \text{src\_ip}, \text{src\_port}, \text{protocol}) \end{aligned}$$

and the values corresponding to a list of flows (with timestamps, category label, and timestamp precision) that match this five-tuple. Since the flows are bidirectional in the spreadsheets, we need to record all flows with both the source and destination IP addresses and ports first. Otherwise we would miss any response packets, since the spreadsheets only acknowledge the initiator of the connection as the source, whereas in packet data, the source and destination fields oscillate based on the sender and receiver. There are a few miscellaneous artifacts of the datasets that we otherwise need to correct for before saving our dictionary. We sort the arrays for each dictionary value by the end time of the flow to produce a ordering of flows that finish first. The next two subsections provide additional information on the conversion of a spreadsheet into entries into the dictionary for both the UNSW-NB15 and CIC-IDS-2017 datasets, respectively.

```

1 def generate_flow_labels_hash(dataset):
2     """
3         Generate a hash for flows to labels with the given tuple:
4         (source_ip, destination_ip, source_port, destination_port).
5
6         @param dataset: ContextPCAPDataset to generate hashes for
7         """
8
9         # create a dictionary of flow hashes and a counter of the number of flows seen
10        flow_labels = {}
11        nflows = 0
12
13        # read all of the flow CSV files
14        for flow_filename in dataset.flow_datasets:
15            flow_start_time = time.time()
16            sys.stdout.write('Parsing flow file {}...'.format(flow_filename))
17
18            if dataset.flow_preprocessor == 'CIC-IDS':
19                data = pd.read_csv(flow_filename, delimiter = ',', encoding = 'cp1252')
20
21                nflows += cic_ids_flow_processor(dataset, data, flow_labels, nflows)
22            elif dataset.flow_preprocessor == 'UNSW-NB15':
23                data = pd.read_csv(flow_filename, dtype = {'attack_cat': str})
24
25                # attack category left blank for normal traffic
26                data['attack_cat'] = data['attack_cat'].fillna('Benign')
27
28                # update column dtypes for ports and clean data
29                dropped_ports = ['-', '0x20205321']
30                data = data[~data['sport'].isin(dropped_ports)]
31                data = data[~data['dsport'].isin(dropped_ports)]
32
33                data.reset_index(drop = True, inplace = True)
34
35                # convert hexadecimal ports into base 10 before column conversion
36                hex_ports = {
```

<sup>9</sup><https://github.com/SRI-CSL/trinity-packet>

```

36         '0x000b': 11,
37         '0x000c': 12,
38         '0xc0a8': 49320,
39         '0xcc09': 52233,
40     }
41     # replace the hex ports
42     data['sport'] = data['sport'].replace(hex_ports)
43     data['dsport'] = data['dsport'].replace(hex_ports)
44
45     data['sport'] = data['sport'].astype(np.int32)
46     data['dsport'] = data['dsport'].astype(np.int32)
47
48     nflows += unsw_nb15_flow_processor(dataset, data, flow_labels, nflows)
49 else:
50     assert ('Unknown flow preprocessor attribute.')
51
52 sys.stdout.write('done in {:.2f} seconds.\n'.format(time.time() - flow_start_time))
53 sys.stdout.flush()
54
55 # since the times are often at coarser resolution than the microsecond durations, sort
56 # by end time the last (start_time, end_time) that satisfies the constraint
57 # start_time <= timestamp <= end_time is the correct flow
58 for hash_tuple in flow_labels:
59     flow_labels[hash_tuple].sort(key = lambda x: x[1])
60
61 # save the hash filename
62 hash_filename = '{}/flow_labels_hash.pkl'.format(dataset.temp_directory)
63 with open(hash_filename, 'wb') as fd:
64     pickle.dump(flow_labels, fd, protocol = pickle.HIGHEST_PROTOCOL)

```

## A.2 UNSW-NB15 Specific Functions

When looking at flows in the UNSW-NB15 data, we only consider flows corresponding to the TCP and UDP protocols. We read the UNSW-NB15 spreadsheets into a pandas dataframe and focus on the following columns: srcip, sport, dstip, dsport, proto, Stime, Ltime, attack\_cat. We return a running count of the total number of flows processed thus far. Note, we need to provide the dictionary and number of flows as arguments/return values because the flow-level datasets are typically split into multiple spreadsheets based on attack and/or time of day. The timestamps given with the flows are already in Coordinated Universal Time, and additionally counting up from the Linux epoch.

```

1 def unsw_nb15_flow_processor(dataset, data, flow_labels, previous_flows):
2 """
3 Process a given flow from UNSW-NB15 and add hashes to the flow labels. Returns the
4 number of processed flows.
5
6 @param dataset: ContextPCAPDataset variable with category mapping attribute
7 @param data: pandas dataframe with relevant flow information
8 @param flow_labels: a dictionary of flow hashes to start_time, end_time, category, etc.
9 @param previous_flows: number of previous flows seen
10 """
11 # count the number of flows processed
12 nflows = 0
13
14 # remove leading and trailing whitespace from column names
15 data.columns = data.columns.str.strip()
16
17 # go through every flow in the CSV file
18 for (src_ip, src_port, dest_ip, dest_port, protocol, start_time, end_time, category) in
19     zip(
20         data['srcip'],
21         data['sport'],
22         data['dstip'],
23         data['dsport'],
24         data['proto'],
25         data['Stime'],
26         data['Ltime'],
27         data['attack_cat'],
28     )
29     # Create a tuple of (src_ip, src_port, dest_ip, dest_port, protocol, start_time, end_time, category)
30     # and add it to the flow_labels dictionary
31     flow_labels[(src_ip, src_port, dest_ip, dest_port, protocol, start_time, end_time, category)] =
32         flow_labels.get((src_ip, src_port, dest_ip, dest_port, protocol, start_time, end_time, category), [])
33         flow_labels[(src_ip, src_port, dest_ip, dest_port, protocol, start_time, end_time, category)].append((start_time, end_time, category))
34
35 # Return the total number of flows processed
36 return nflows

```

```

25     data['Ltime'],
26     data['attack_cat'],
27 ):
28     # TCP is protocol 6 and UDP is protocol 17 by the Internet Assigned Numbers
29     # Authority (IANA)
30     if not protocol == 'tcp' and not protocol == 'udp': continue
31     protocol = protocol.upper()
32
33     # remove white space around category and add 'Benign' for empty attacks (benign
34     # traffic)
35     category = category.strip()
36     if not len(category): category = 'Benign'
37
38     # update the category name if a category mapping exists
39     if hasattr(dataset, 'category_mapping'):
40         category = dataset.category_mapping[category]
41
42     # precision is zero second for the UNSW-NB15 datasets since the start and end
43     # time are both given (in seconds)
44     precision = 0
45
46     # create a list of this hash tuple if it doesn't already exist
47     hash_tuple = (src_ip, src_port, dest_ip, dest_port, protocol)
48     if not hash_tuple in flow_labels:
49         flow_labels[hash_tuple] = []
50     flow_labels[hash_tuple].append((start_time, end_time, precision, category,
51         previous_flows + nflows))
52     # since these flows are bidirectional, we also need to hash the reverse flow
53     # information
54     # when iterating through dpkt, it will give (src_ip, src_port, dest_ip, dest_port,
55     # protocol)
56     # but the CSV file only considers the src_ip as the initiator of the connection,
57     # will miss
58     # all response packets
59     hash_tuple = (dest_ip, dest_port, src_ip, src_port, protocol)
60     if not hash_tuple in flow_labels:
61         flow_labels[hash_tuple] = []
62     flow_labels[hash_tuple].append((start_time, end_time, precision, category,
63         previous_flows + nflows))
64
65     # increment the number of flows seen
66     nflows += 1
67
68 return nflows

```

### A.3 CIC-IDS-2017 Specific Functions

One of the issues with converting the CIC-IDS-2017 flow-level data into packet-level data is the timestamps associated with the flows. These timestamps are in a string format both with and without second precision. This code takes a timestamp in the CIC-IDS-2017 string format and converts it into a timestamp in Coordinated Universal Time used by various PCAP capturing systems. We also return the precision of the time given. PCAP data has significantly higher time precision than these flow-level datasets. Thus, we keep track of the precision of the flows to provide lower and upper bounds for their actual start and end times.

```

1 def cic_ids_time_parser(str_time):
2     """
3
4     Time parser specific to the CIC-IDS-2017 dataset. Different datasets require different
5     parsing functions based on the time format. For the CIC-IDS-2017 dataset, attack
6     times are given in the New Brunswick timezone GMT-3 during the summer months. AM/PM
7     are not given but since all traffic occurs between 9am and 5pm, we add 12 to
8     afternoon hours.
9
10    @param str_time: the string representation of the time.
11    """
12    # get the date and time of the attack

```

```
12 date, time = str_time.split()
13 day, month, year = date.split('/')
14
15 # determine if seconds are given or not
16 if time.count(':') == 1:
17     hour, minute = time.split(':')
18     second, precision_in_seconds = 0, 60
19 else:
20     hour, minute, second = time.split(':')
21     precision_in_seconds = 1
22
23 # if the hour is in the afternoon (less than 8pm since some attacks
24 # begin in the 8th hours), add 12 to convert to military time
25 hour = int(hour)
26 if hour < 8: hour += 12
27 # add three to the hour to get the time in GMT format to match the
28 # times in dpkt pcap format. current times are in New Brunswick time
29 # in summer months which has GMT-3. we do not need to worry about the
30 # day changing since the last attacks occurs at 5pm ADT (8pm GMT)
31 hour += 3
32 # convert into a timestamp to get the float value for return
33 timestamp = datetime(
34     day = int(day),
35     month = int(month),
36     year = int(year),
37     hour = int(hour),
38     minute = int(minute),
39     second = int(second),
40     tzinfo = timezone.utc,
41 )
42
43 # convert to a float value and return the precision in seconds
44 return timestamp.timestamp(), precision_in_seconds
```

When looking at flows in the CIC-IDS-2017 data, we only consider flows corresponding to TCP and UDP. We read the CIC-IDS-2017 spreadsheets into a pandas dataframe and focus on the following columns: Source IP, Source Port, Destination IP, Destination Port, Protocol, Timestamp, Flow Duration, and Label. We return a running count of the total number of flows processed thus far. Note, we need to provide the dictionary and number of flows as arguments/return values because the flow-level datasets are typically split into multiple spreadsheets based on attack and/or time of day.

```
1 def cic_ids_flow_processor(dataset, data, flow_labels, previous_flows):
2     """
3         Process a given flow from CIC-IDS and add hashes to the flow labels. Returns the number
4         of processed flows.
5
6     @param dataset: ContextPCAPDataset variable with category mapping attribute
7     @param data: pandas dataframe with relevant flow information
8     @param flow_labels: a dictionary of flow hashes to start_time, end_time, category, etc.
9     @param previous_flows: number of previous flows seen
10    """
11    # count the number of flows processed
12    nflows = 0
13
14    # remove leading and trailing whitespace from column names
15    data.columns = data.columns.str.strip()
16
17    # go through every flow in the CSV file
18    for (src_ip, src_port, dest_ip, dest_port, protocol, timestamp, duration, category) in
19        zip(
20            data['Source IP'],
21            data['Source Port'],
22            data['Destination IP'],
23            data['Destination Port'],
24            data['Protocol'],
25            data['Timestamp'],
```

```

25     data['Flow Duration'],
26     data['Label']
27   ):
28     # TCP is protocol 6 and UDP is protocol 17 by the Internet Assigned Numbers
29     # Authority (IANA)
30     if not protocol == 6 and not protocol == 17: continue
31     if protocol == 6: protocol = 'TCP'
32     else: protocol = 'UDP'
33
34     # get the start time as an integer conditioned on whether or not seconds are
35     # included
36     start_time, precision = cic_ids_time_parser(timestamp)
37     # CIC-IDS-2017 flow durations in microseconds (cannot add since the integral values
38     # are milliseconds)
39     end_time = (datetime.fromtimestamp(start_time) + timedelta(microseconds = int(
40         duration))).timestamp()
41
42     # update the category name if a category mapping exists
43     if hasattr(dataset, 'category_mapping'):
44       category = dataset.category_mapping[category]
45
46     # create a list of this hash tuple if it doesn't already exist
47     hash_tuple = (src_ip, src_port, dest_ip, dest_port, protocol)
48     if not hash_tuple in flow_labels:
49       flow_labels[hash_tuple] = []
50     flow_labels[hash_tuple].append((start_time, end_time, precision, category,
51         previous_flows + nflows))
52     # since these flows are bidirectional, we also need to hash the reverse flow
53     # information
54     # when iterating through dpkt, it will give (src_ip, src_port, dest_ip, dest_port,
55     # protocol)
56     # but the CSV file only considers the src_ip as the initiator of the connection,
57     # will miss
58     # all response packets
59     hash_tuple = (dest_ip, dest_port, src_ip, src_port, protocol)
60     if not hash_tuple in flow_labels:
61       flow_labels[hash_tuple] = []
62     flow_labels[hash_tuple].append((start_time, end_time, precision, category,
63         previous_flows + nflows))
64
65     # increment the number of flows seen
66     nflows += 1
67
68   return nflows

```

#### A.4 Matching Packet Data with Flow Labels

The following function takes a pcap file and the flow labels dictionaries generated in the previous sections and labels the packets themselves. There are a few main points to consider when constructing this function. The TCP flags are presented in a single integer value requiring us to do bit-wise operations to see which flags are set for each packet. For packets following the UDP protocols, we simply zero out those excluded values. We then identify the list of possible flows that match based on the five-tuple and its reciprocal:

$$\begin{aligned} & (\text{src\_ip}, \text{src\_port}, \text{dest\_ip}, \text{dest\_port}, \text{protocol}) \\ & (\text{dest\_ip}, \text{dest\_port}, \text{src\_ip}, \text{src\_port}, \text{protocol}) \end{aligned}$$

Flow-level data constructors use a small set number of IP addresses and port numbers to more easily identify labels. However, this causes problems when assigning packets to flows for labeling. Thus for each packet's five-tuple, we iterate over all matching flows and identify those valid given the start and end times recorded previously. We note, however, that the precision of the flow-level datasets is coarser than the PCAP files provided. Thus, we match a packet to any flow where the packet's timestamp falls after the start time but before the end time plus the (im-)precision for that particular flow. There are occasionally multiple flows that will match with a given packet. This often happens (in these academic datasets) in denial-of-service attacks,

for example, where flow durations can be quite short. If the flow categories do not match, we simply drop the packet from analysis to not pollute the data. At this point, we have the information needed to store into a dataframe for input into the SAFE-NID program. Note, a real-world implementation of SAFE-NID would not require this burdensome processing since there would be no labels needed.

```

1 def inet_to_str(inet):
2     """
3         Convert an inet object to a string representation.
4
5     @param inet: inet network address
6     """
7     # return either IPv4 or IPv6 address
8     try:
9         return socket.inet_ntop(socket.AF_INET, inet)
10    except ValueError:
11        return socket.inet_ntop(socket.AF_INET6, inet)
12
13 def read_pcap_data(pcap_file, flow_labels):
14     """
15         Given a pcap class, read the pcap data packet by packet.
16
17     @param pcap: type PCAPFile that contains filenames and types
18     @param flow_labels: labels for each flow identified by tuple (src ip, src port, dest ip,
19     dest port)
20     """
21     # store all packets into an array to save into a dataframe later
22     packets = []
23
24     # read this PCAP file and return the TCP/UDP packets
25     with open(pcap_file.filename, 'rb') as pcap_fd:
26         pcap_start_time = time.time()
27         sys.stdout.write('Reading PCAP file {}...'.format(pcap_file.filename))
28         # read either PCAP or PCAP next generation formats
29         if pcap_file.pcap_type == 'PCAP':
30             pcap = dpkt.pcap.Reader(pcap_fd)
31         elif pcap_file.pcap_type == 'PCAPNG':
32             pcap = dpkt.pcapng.Reader(pcap_fd)
33
34         npackets = 0
35         ndropped = 0
36         # read each packet in the pcap file
37         for timestamp, packet in pcap:
38             # Linux cooked capture
39             if packet.datalink() == dpkt.pcap.DLT_LINUX_SLL:
40                 eth = dpkt.sll.SLL(packet)
41             # ethernet format
42             else:
43                 eth = dpkt.ethernet.Ethernet(packet)
44
45             # skip any non IP IPv6 packets
46             if not isinstance(eth.data, dpkt.ip.IP) and not isinstance(eth.data, dpkt.ip6.IPv6):
47                 IP6):
48                 continue
49
50             # get the IP data
51             ip = eth.data
52             # get relevant fields from the IP packet header (convert to string
53             # representation)
54             src_ip = inet_to_str(ip.src)
55             dest_ip = inet_to_str(ip.dst)
56             # get specific attributes from the header based on the protocol type (IP v IPv6)
57             if isinstance(eth.data, dpkt.ip.IP):
58                 ttl = ip.ttl
59                 total_length = ip.len
60                 internet_layer_protocol = 'IPv4'
61             elif isinstance(eth.data, dpkt.ip6.IPv6):
62                 ttl = ip.hlim
63                 total_length = ip.plen

```

```

63     internet_layer_protocol = 'IPv6'
64 else:
65     assert ('Unknown internet layer protocol.')
66
67 # skip any non UDP and TCP protocols (n.b., get_proto has a bug and does not
68 # recognize certain protocols, so we need to skip the others here)
69 if not isinstance(ip.data, dpkt.tcp.TCP) and not isinstance(ip.data, dpkt.udp.
70     UDP):
71     continue
72
73 # get the protocol name
74 transport_layer_protocol = ip.get_proto(ip.p).__name__
75 transport_layer_protocol = transport_layer_protocol.upper()
76
77 # count the number of packets before any drops
78 npackets += 1
79
80 # get the data from the transport layer and convert to bytes
81 transport = ip.data
82 data = transport.data
83
84 # get the source and destination ports
85 src_port = transport.sport
86 dest_port = transport.dport
87
88 # get relevant header information for TCP
89 if isinstance(ip.data, dpkt.tcp.TCP):
90     # need to divide by 2 raised to the bit location to get 0 and 1 values
91     cwr_flag = (transport.flags & dpkt.tcp.TH_CWR) // 128
92     ece_flag = (transport.flags & dpkt.tcp.TH_ECE) // 64
93     urg_flag = (transport.flags & dpkt.tcp.TH_URG) // 32
94     ack_flag = (transport.flags & dpkt.tcp.TH_ACK) // 16
95     psh_flag = (transport.flags & dpkt.tcp.TH_PUSH) // 8
96     rst_flag = (transport.flags & dpkt.tcp.TH_RST) // 4
97     syn_flag = (transport.flags & dpkt.tcp.TH_SYN) // 2
98     fin_flag = (transport.flags & dpkt.tcp.TH_FIN)
99 # all flags have zero value for UDP
100 else:
101     cwr_flag = 0
102     ece_flag = 0
103     urg_flag = 0
104     ack_flag = 0
105     psh_flag = 0
106     rst_flag = 0
107     syn_flag = 0
108     fin_flag = 0
109
110 # get the label for this packet
111 hash_tuple = (src_ip, src_port, dest_ip, dest_port, transport_layer_protocol)
112 # skip tuples that do not appear in the flow data
113 # for the CIC-IDS-2017 dataset, this will include all IPv6 traffic
114 if not hash_tuple in flow_labels:
115     ndropped += 1
116     continue
117
118 # go through possible hashed label values until identifying the proper flow
119 # based on the flow's start and end times. there can be multiple matches
120 # since the start time resolution is in seconds and the end time resolution is
121 # in microseconds so find the last flow that fits the conditions (must be
122 # the correct one)
123 packet_category = None
124 packet_flow_id = None
125 # the start time precision is only to the second (at best), we allow a small
126 # buffer for flows to be considered. if all flows in the buffer have the
127 # same category, the packet receives that category label
128 buffer_categories = set()
buffer_flow_ids = set()

```

```

129     for (start_time, end_time, precision, category, flow_id) in flow_labels[
130         hash_tuple]:
131             # consider all possible flows that can fall in this time frame
132             if start_time <= timestamp and timestamp <= end_time + precision:
133                 buffer_categories.add(category)
134                 buffer_flow_ids.add(flow_id)
135
136             # continue if there are overlapping flows with different category labels
137             if len(buffer_categories) == 1:
138                 packet_category = list(buffer_categories)[0]
139                 # take the first flow ID (not a perfect method)
140                 packet_flow_id = list(buffer_flow_ids)[0]
141             else:
142                 ndropped += 1
143                 continue
144
145             # only Benign category is labeled as non-attack
146             packet_label = 0 if packet_category == 'Benign' else 1
147
148             # get the hex representation of the payload
149             if isinstance(data, bytes):
150                 payload = data.hex()
151             else:
152                 payload = (data.pack()).hex()
153
154             # save the payload length (for pruning later)
155             payload_length = len(payload)
156
157             packets.append([
158                 src_ip,
159                 src_port,
160                 dest_ip,
161                 dest_port,
162                 timestamp,
163                 internet_layer_protocol,
164                 transport_layer_protocol,
165                 ttl,
166                 total_length,
167                 cwr_flag,
168                 ece_flag,
169                 urg_flag,
170                 ack_flag,
171                 psh_flag,
172                 rst_flag,
173                 syn_flag,
174                 fin_flag,
175                 payload,
176                 payload_length,
177                 packet_flow_id,
178                 packet_category,
179                 packet_label,
180             ])
181
182             sys.stdout.write('read {} packets (dropped: {}) in {:.2f} seconds.\n'.format(
183                 npackets, ndropped, time.time() - pcap_start_time))
184             sys.stdout.flush()
185
186     return packets

```

## B Complete Out-of-Distribution Results

Table 12: *FTP-Patator out-of-distribution results*. The transformer architecture with either the Gaussian kernel density or normalizing flows safeguard achieves very high AU ROC scores  $> 0.9936$ . The FNN model performs better than the CNN on this data, but still is significantly worse than when using a transformer architecture.

Architecture	Layer	AU ROC ( $\uparrow$ )	TPR (TNR 95%) ( $\uparrow$ )	TPR (TNR 85%) ( $\uparrow$ )
<b>FTP-Patator</b>				
<b>Maximum Softmax Probability</b>				
CNN	Output	0.7931 ( $\pm 0.1534$ )	15.45% ( $\pm 15.48\%$ )	64.19% ( $\pm 37.97\%$ )
FNN	Output	0.6149 ( $\pm 0.0893$ )	0.12% ( $\pm 0.37\%$ )	9.89% ( $\pm 11.39\%$ )
Transformer	Output	0.4034 ( $\pm 0.2232$ )	0.21% ( $\pm 0.54\%$ )	14.61% ( $\pm 13.37\%$ )
<b>Energy-Based</b>				
CNN	Output	0.7628 ( $\pm 0.1657$ )	12.51% ( $\pm 12.78\%$ )	52.82% ( $\pm 32.85\%$ )
FNN	Output	0.6060 ( $\pm 0.1500$ )	0.19% ( $\pm 0.56\%$ )	16.55% ( $\pm 19.86\%$ )
Transformer	Output	0.4133 ( $\pm 0.1955$ )	0.21% ( $\pm 0.51\%$ )	15.98% ( $\pm 17.58\%$ )
<b>Gaussian Kernel Density</b>				
CNN	dense1	0.7563 ( $\pm 0.2059$ )	0.53% ( $\pm 0.51\%$ )	57.66% ( $\pm 26.93\%$ )
CNN	dense2	0.8207 ( $\pm 0.0684$ )	4.09% ( $\pm 2.69\%$ )	51.64% ( $\pm 30.72\%$ )
CNN	dense3	0.8297 ( $\pm 0.1157$ )	26.47% ( $\pm 24.29\%$ )	58.92% ( $\pm 29.76\%$ )
FNN	linear3	0.8473 ( $\pm 0.0259$ )	3.55% ( $\pm 8.46\%$ )	59.20% ( $\pm 19.59\%$ )
FNN	linear4	0.7511 ( $\pm 0.0406$ )	2.56% ( $\pm 4.37\%$ )	19.70% ( $\pm 15.77\%$ )
FNN	linear5	0.7585 ( $\pm 0.0591$ )	0.37% ( $\pm 0.47\%$ )	24.53% ( $\pm 17.33\%$ )
Transformer	linear1	<b>0.9950</b> ( $\pm 0.0018$ )	<b>100.00%</b> ( $\pm 0.00\%$ )	<b>100.00%</b> ( $\pm 0.00\%$ )
Transformer	linear2	0.9717 ( $\pm 0.0083$ )	89.63% ( $\pm 16.73\%$ )	<b>100.00%</b> ( $\pm 0.00\%$ )
Transformer	linear3	0.9255 ( $\pm 0.0406$ )	44.22% ( $\pm 33.00\%$ )	89.51% ( $\pm 16.49\%$ )
<b>Normalizing Flows</b>				
CNN	dense1	0.7442 ( $\pm 0.2052$ )	0.50% ( $\pm 0.72\%$ )	51.70% ( $\pm 35.34\%$ )
CNN	dense2	0.8058 ( $\pm 0.1224$ )	4.53% ( $\pm 5.42\%$ )	48.96% ( $\pm 31.88\%$ )
CNN	dense3	0.7998 ( $\pm 0.1178$ )	15.46% ( $\pm 14.41\%$ )	56.90% ( $\pm 29.89\%$ )
FNN	linear3	0.8640 ( $\pm 0.0152$ )	6.88% ( $\pm 12.78\%$ )	61.51% ( $\pm 14.50\%$ )
FNN	linear4	0.7917 ( $\pm 0.0389$ )	3.37% ( $\pm 8.52\%$ )	29.55% ( $\pm 14.60\%$ )
FNN	linear5	0.7732 ( $\pm 0.0580$ )	7.25% ( $\pm 7.74\%$ )	38.37% ( $\pm 13.60\%$ )
Transformer	linear1	0.9936 ( $\pm 0.0014$ )	<b>100.00%</b> ( $\pm 0.00\%$ )	<b>100.00%</b> ( $\pm 0.00\%$ )
Transformer	linear2	0.9726 ( $\pm 0.0064$ )	92.71% ( $\pm 13.52\%$ )	<b>100.00%</b> ( $\pm 0.00\%$ )
Transformer	linear3	0.9480 ( $\pm 0.0199$ )	58.83% ( $\pm 26.97\%$ )	97.50% ( $\pm 4.21\%$ )

Table 13: *Infiltration out-of-distribution results*. The transformer with normalizing flows safeguard performs worse here than with the other two zero-day exploits. However, it still outperforms baseline methods on AU ROC. The CNN architecture performs well here compared to FTP-Patator and SSH-Patator results.

Architecture	Layer	AU ROC ( $\uparrow$ )	TPR (TNR 95%) ( $\uparrow$ )	TPR (TNR 85%) ( $\uparrow$ )
<b>Infiltration</b>				
<b>Maximum Softmax Probability</b>				
CNN	Output	0.6611 ( $\pm 0.1897$ )	26.16% ( $\pm 20.54\%$ )	53.95% ( $\pm 23.83\%$ )
FNN	Output	0.7676 ( $\pm 0.0567$ )	12.39% ( $\pm 11.58\%$ )	47.52% ( $\pm 16.07\%$ )
Transformer	Output	0.3167 ( $\pm 0.2055$ )	3.37% ( $\pm 6.66\%$ )	9.56% ( $\pm 19.88\%$ )
<b>Energy-Based</b>				
CNN	Output	0.6744 ( $\pm 0.1859$ )	25.33% ( $\pm 21.43\%$ )	53.56% ( $\pm 26.27\%$ )
FNN	Output	0.7592 ( $\pm 0.0750$ )	11.41% ( $\pm 10.28\%$ )	47.91% ( $\pm 15.66\%$ )
Transformer	Output	0.3037 ( $\pm 0.2291$ )	4.33% ( $\pm 11.43\%$ )	10.13% ( $\pm 23.19\%$ )
<b>Gaussian Kernel Density</b>				
CNN	dense1	0.9581 ( $\pm 0.0062$ )	<b>83.74% (<math>\pm 10.23\%</math>)</b>	98.70% ( $\pm 0.53\%$ )
CNN	dense2	0.9511 ( $\pm 0.0218$ )	75.73% ( $\pm 19.16\%$ )	96.55% ( $\pm 2.93\%$ )
CNN	dense3	0.9622 ( $\pm 0.0217$ )	79.11% ( $\pm 15.19\%$ )	96.29% ( $\pm 2.99\%$ )
FNN	linear3	0.8740 ( $\pm 0.0337$ )	27.72% ( $\pm 13.36\%$ )	67.78% ( $\pm 13.74\%$ )
FNN	linear4	0.8748 ( $\pm 0.0405$ )	41.70% ( $\pm 15.30\%$ )	61.45% ( $\pm 16.68\%$ )
FNN	linear5	0.8345 ( $\pm 0.0439$ )	29.58% ( $\pm 19.39\%$ )	47.85% ( $\pm 18.80\%$ )
Transformer	linear1	0.9455 ( $\pm 0.0248$ )	62.58% ( $\pm 20.95\%$ )	93.79% ( $\pm 10.34\%$ )
Transformer	linear2	0.9552 ( $\pm 0.0318$ )	71.26% ( $\pm 28.99\%$ )	94.93% ( $\pm 13.73\%$ )
Transformer	linear3	0.8996 ( $\pm 0.0580$ )	35.48% ( $\pm 28.28\%$ )	75.45% ( $\pm 25.56\%$ )
<b>Normalizing Flows</b>				
CNN	dense1	0.9464 ( $\pm 0.0061$ )	67.72% ( $\pm 15.48\%$ )	98.90% ( $\pm 0.05\%$ )
CNN	dense2	0.9519 ( $\pm 0.0057$ )	72.07% ( $\pm 18.42\%$ )	98.91% ( $\pm 0.06\%$ )
CNN	dense3	0.9638 ( $\pm 0.0153$ )	77.08% ( $\pm 17.79\%$ )	98.22% ( $\pm 0.78\%$ )
FNN	linear3	0.8655 ( $\pm 0.0115$ )	5.89% ( $\pm 4.88\%$ )	66.73% ( $\pm 11.63\%$ )
FNN	linear4	0.8820 ( $\pm 0.0300$ )	30.41% ( $\pm 15.53\%$ )	67.72% ( $\pm 14.22\%$ )
FNN	linear5	0.8530 ( $\pm 0.0385$ )	29.11% ( $\pm 18.50\%$ )	59.60% ( $\pm 14.84\%$ )
Transformer	linear1	0.9555 ( $\pm 0.0146$ )	63.20% ( $\pm 24.30\%$ )	<b>99.64% (<math>\pm 0.62\%</math>)</b>
Transformer	linear2	<b>0.9668 (<math>\pm 0.0182</math>)</b>	77.89% ( $\pm 23.41\%$ )	99.18% ( $\pm 2.25\%$ )
Transformer	linear3	0.9361 ( $\pm 0.0455$ )	56.13% ( $\pm 31.83\%$ )	89.54% ( $\pm 15.68\%$ )

Table 14: *SSH-Patator out-of-distribution results*. The transformer classifier with normalizing flows safeguard outperforms all existing methods with a very high AU ROC of 0.9901.

Architecture	Layer	AU ROC ( $\uparrow$ )	TPR (TNR 95%) ( $\uparrow$ )	TPR (TNR 85%) ( $\uparrow$ )
<b>SSH-Patator</b>				
<b>Maximum Softmax Probability</b>				
CNN	Output	0.7821 ( $\pm 0.0511$ )	17.19% ( $\pm 10.35\%$ )	66.61% ( $\pm 15.90\%$ )
FNN	Output	0.6568 ( $\pm 0.0727$ )	1.53% ( $\pm 0.87\%$ )	31.75% ( $\pm 11.64\%$ )
Transformer	Output	0.1924 ( $\pm 0.1103$ )	0.20% ( $\pm 0.35\%$ )	2.24% ( $\pm 2.45\%$ )
<b>Energy-Based</b>				
CNN	Output	0.7807 ( $\pm 0.0701$ )	18.90% ( $\pm 9.81\%$ )	64.77% ( $\pm 14.63\%$ )
FNN	Output	0.5932 ( $\pm 0.0896$ )	1.38% ( $\pm 0.92\%$ )	23.38% ( $\pm 12.69\%$ )
Transformer	Output	0.1833 ( $\pm 0.1142$ )	0.19% ( $\pm 0.34\%$ )	2.12% ( $\pm 2.57\%$ )
<b>Gaussian Kernel Density</b>				
CNN	dense1	0.6158 ( $\pm 0.1096$ )	7.14% ( $\pm 1.62\%$ )	27.81% ( $\pm 13.33\%$ )
CNN	dense2	0.6993 ( $\pm 0.1041$ )	11.45% ( $\pm 7.93\%$ )	39.42% ( $\pm 11.65\%$ )
CNN	dense3	0.7254 ( $\pm 0.1269$ )	23.72% ( $\pm 10.20\%$ )	46.08% ( $\pm 12.92\%$ )
FNN	linear3	0.8967 ( $\pm 0.0100$ )	26.14% ( $\pm 4.38\%$ )	77.99% ( $\pm 5.01\%$ )
FNN	linear4	0.8381 ( $\pm 0.0169$ )	19.34% ( $\pm 3.42\%$ )	52.81% ( $\pm 6.62\%$ )
FNN	linear5	0.8298 ( $\pm 0.0265$ )	18.36% ( $\pm 3.01\%$ )	56.14% ( $\pm 7.42\%$ )
Transformer	linear1	0.9821 ( $\pm 0.0037$ )	94.90% ( $\pm 2.61\%$ )	99.05% ( $\pm 1.09\%$ )
Transformer	linear2	0.9629 ( $\pm 0.0048$ )	76.63% ( $\pm 9.36\%$ )	<b>100.00% (<math>\pm 0.00\%</math>)</b>
Transformer	linear3	0.9097 ( $\pm 0.0313$ )	27.03% ( $\pm 14.36\%$ )	86.99% ( $\pm 12.27\%$ )
<b>Normalizing Flows</b>				
CNN	dense1	0.5739 ( $\pm 0.1273$ )	6.90% ( $\pm 0.51\%$ )	18.94% ( $\pm 6.43\%$ )
CNN	dense2	0.7086 ( $\pm 0.1117$ )	11.23% ( $\pm 4.00\%$ )	39.33% ( $\pm 13.11\%$ )
CNN	dense3	0.7010 ( $\pm 0.1336$ )	21.97% ( $\pm 11.52\%$ )	44.46% ( $\pm 15.20\%$ )
FNN	linear3	0.9119 ( $\pm 0.0062$ )	30.81% ( $\pm 3.04\%$ )	84.98% ( $\pm 3.99\%$ )
FNN	linear4	0.8810 ( $\pm 0.0144$ )	30.83% ( $\pm 3.27\%$ )	70.58% ( $\pm 6.78\%$ )
FNN	linear5	0.8596 ( $\pm 0.0312$ )	30.49% ( $\pm 3.33\%$ )	67.67% ( $\pm 8.18\%$ )
Transformer	linear1	<b>0.9901 (<math>\pm 0.0007</math>)</b>	<b>100.00% (<math>\pm 0.00\%</math>)</b>	<b>100.00% (<math>\pm 0.00\%</math>)</b>
Transformer	linear2	0.9668 ( $\pm 0.0062$ )	84.62% ( $\pm 12.51\%$ )	<b>100.00% (<math>\pm 0.00\%</math>)</b>
Transformer	linear3	0.9323 ( $\pm 0.0180$ )	45.45% ( $\pm 13.92\%$ )	92.96% ( $\pm 6.04\%$ )

Table 15: *CIC-IDS-2017 out-of-distribution results*. The highest performing combination for detecting the out-of-distribution samples is the the transformer descriminative classifier with the normalizing flows safe-guard.

Architecture	Layer	AU ROC ( $\uparrow$ )	TPR (TNR 95%) ( $\uparrow$ )	TPR (TNR 85%) ( $\uparrow$ )
<b>CIC-IDS-2017</b>				
<b>Maximum Softmax Probability</b>				
CNN	Output	0.7818 ( $\pm 0.0254$ )	36.43% ( $\pm 4.48\%$ )	63.02% ( $\pm 4.47\%$ )
FNN	Output	0.7187 ( $\pm 0.0289$ )	23.25% ( $\pm 3.21\%$ )	50.69% ( $\pm 3.18\%$ )
Transformer	Output	0.6731 ( $\pm 0.0737$ )	23.41% ( $\pm 8.93\%$ )	38.37% ( $\pm 14.72\%$ )
<b>Energy-Based</b>				
CNN	Output	0.7589 ( $\pm 0.0436$ )	36.79% ( $\pm 6.40\%$ )	60.53% ( $\pm 3.58\%$ )
FNN	Output	0.7114 ( $\pm 0.0498$ )	22.90% ( $\pm 3.63\%$ )	49.71% ( $\pm 3.14\%$ )
Transformer	Output	0.6419 ( $\pm 0.0966$ )	21.88% ( $\pm 9.79\%$ )	35.48% ( $\pm 14.45\%$ )
<b>Gaussian Kernel Density</b>				
CNN	dense1	0.8814 ( $\pm 0.0144$ )	49.93% ( $\pm 4.44\%$ )	75.21% ( $\pm 3.38\%$ )
CNN	dense2	0.8981 ( $\pm 0.0170$ )	55.39% ( $\pm 5.94\%$ )	77.16% ( $\pm 5.28\%$ )
CNN	dense3	0.8580 ( $\pm 0.0394$ )	55.26% ( $\pm 4.30\%$ )	65.82% ( $\pm 5.45\%$ )
FNN	linear3	0.6986 ( $\pm 0.0464$ )	25.80% ( $\pm 5.13\%$ )	40.43% ( $\pm 6.12\%$ )
FNN	linear4	0.7343 ( $\pm 0.0597$ )	34.96% ( $\pm 4.17\%$ )	45.25% ( $\pm 5.73\%$ )
FNN	linear5	0.6801 ( $\pm 0.0520$ )	25.72% ( $\pm 2.39\%$ )	35.15% ( $\pm 3.33\%$ )
Transformer	linear1	0.8323 ( $\pm 0.0649$ )	50.61% ( $\pm 8.24\%$ )	64.29% ( $\pm 10.40\%$ )
Transformer	linear2	0.8753 ( $\pm 0.0552$ )	58.05% ( $\pm 11.80\%$ )	73.07% ( $\pm 11.11\%$ )
Transformer	linear3	0.7893 ( $\pm 0.0541$ )	42.14% ( $\pm 6.95\%$ )	53.56% ( $\pm 9.77\%$ )
<b>Normalizing Flows</b>				
CNN	dense1	0.8744 ( $\pm 0.0121$ )	30.50% ( $\pm 3.95\%$ )	71.26% ( $\pm 5.71\%$ )
CNN	dense2	0.8960 ( $\pm 0.0151$ )	52.88% ( $\pm 3.63\%$ )	76.19% ( $\pm 5.61\%$ )
CNN	dense3	0.8384 ( $\pm 0.0594$ )	53.77% ( $\pm 5.61\%$ )	67.10% ( $\pm 8.80\%$ )
FNN	linear3	0.8130 ( $\pm 0.0226$ )	27.27% ( $\pm 6.58\%$ )	50.10% ( $\pm 2.03\%$ )
FNN	linear4	0.7625 ( $\pm 0.0210$ )	35.61% ( $\pm 1.12\%$ )	48.22% ( $\pm 2.08\%$ )
FNN	linear5	0.7139 ( $\pm 0.0582$ )	31.64% ( $\pm 2.98\%$ )	42.69% ( $\pm 5.90\%$ )
Transformer	linear1	0.9011 ( $\pm 0.0115$ )	66.87% ( $\pm 3.37\%$ )	76.70% ( $\pm 2.31\%$ )
Transformer	linear2	<b>0.9259</b> ( $\pm 0.0046$ )	<b>73.01%</b> ( $\pm 2.10\%$ )	<b>81.55%</b> ( $\pm 1.11\%$ )
Transformer	linear3	0.8945 ( $\pm 0.0085$ )	63.15% ( $\pm 4.94\%$ )	77.45% ( $\pm 2.06\%$ )

Table 16: *UNSW-NB15 out-of-distribution results*. Although transformer with normalizing flows performs well, the best performer comes from the transformer with the Gaussian kernel density safeguard.

Architecture	Layer	AU ROC ( $\uparrow$ )	TPR (TNR 95%) ( $\uparrow$ )	TPR (TNR 85%) ( $\uparrow$ )
<b>UNSW-NB15</b>				
<b>Maximum Softmax Probability</b>				
CNN	Output	0.5590 ( $\pm 0.0063$ )	34.87% ( $\pm 2.45\%$ )	44.47% ( $\pm 3.19\%$ )
FNN	Output	0.7349 ( $\pm 0.0417$ )	25.98% ( $\pm 2.04\%$ )	43.43% ( $\pm 5.41\%$ )
Transformer	Output	0.6255 ( $\pm 0.1033$ )	24.21% ( $\pm 11.24\%$ )	40.14% ( $\pm 14.78\%$ )
<b>Energy-Based</b>				
CNN	Output	0.5381 ( $\pm 0.0480$ )	32.24% ( $\pm 5.82\%$ )	41.78% ( $\pm 8.91\%$ )
FNN	Output	0.7177 ( $\pm 0.0977$ )	26.62% ( $\pm 2.73\%$ )	44.61% ( $\pm 8.65\%$ )
Transformer	Output	0.6153 ( $\pm 0.1003$ )	24.26% ( $\pm 11.29\%$ )	39.55% ( $\pm 14.19\%$ )
<b>Gaussian Kernel Density</b>				
CNN	dense1	0.9137 ( $\pm 0.0208$ )	53.07% ( $\pm 2.51\%$ )	79.71% ( $\pm 7.07\%$ )
CNN	dense2	0.9263 ( $\pm 0.0139$ )	61.61% ( $\pm 4.44\%$ )	83.14% ( $\pm 4.44\%$ )
CNN	dense3	0.9217 ( $\pm 0.0118$ )	62.52% ( $\pm 3.50\%$ )	81.29% ( $\pm 4.04\%$ )
FNN	linear3	0.8600 ( $\pm 0.0178$ )	27.30% ( $\pm 4.04\%$ )	65.76% ( $\pm 6.38\%$ )
FNN	linear4	0.9079 ( $\pm 0.0109$ )	56.06% ( $\pm 5.32\%$ )	80.31% ( $\pm 2.68\%$ )
FNN	linear5	0.8759 ( $\pm 0.0235$ )	43.55% ( $\pm 9.04\%$ )	71.70% ( $\pm 6.78\%$ )
Transformer	linear1	0.8762 ( $\pm 0.0227$ )	43.21% ( $\pm 6.52\%$ )	64.55% ( $\pm 8.58\%$ )
Transformer	linear2	<b>0.9636 (<math>\pm 0.0126</math>)</b>	<b>79.23% (<math>\pm 8.25\%</math>)</b>	<b>95.73% (<math>\pm 2.68\%</math>)</b>
Transformer	linear3	0.9550 ( $\pm 0.0094$ )	71.49% ( $\pm 8.26\%$ )	95.27% ( $\pm 1.73\%$ )
<b>Normalizing Flows</b>				
CNN	dense1	0.9071 ( $\pm 0.0047$ )	48.95% ( $\pm 1.34\%$ )	74.30% ( $\pm 2.58\%$ )
CNN	dense2	0.9263 ( $\pm 0.0063$ )	55.98% ( $\pm 1.29\%$ )	83.03% ( $\pm 2.72\%$ )
CNN	dense3	0.8982 ( $\pm 0.0130$ )	58.78% ( $\pm 1.38\%$ )	72.69% ( $\pm 3.38\%$ )
FNN	linear3	0.8567 ( $\pm 0.0137$ )	22.42% ( $\pm 3.62\%$ )	62.69% ( $\pm 6.80\%$ )
FNN	linear4	0.8963 ( $\pm 0.0112$ )	45.10% ( $\pm 8.46\%$ )	77.66% ( $\pm 3.28\%$ )
FNN	linear5	0.8549 ( $\pm 0.0219$ )	36.02% ( $\pm 7.18\%$ )	64.46% ( $\pm 5.96\%$ )
Transformer	linear1	0.9263 ( $\pm 0.0071$ )	77.83% ( $\pm 0.39\%$ )	81.84% ( $\pm 1.61\%$ )
Transformer	linear2	0.9536 ( $\pm 0.0072$ )	77.96% ( $\pm 3.03\%$ )	91.07% ( $\pm 2.59\%$ )
Transformer	linear3	0.9583 ( $\pm 0.0090$ )	77.71% ( $\pm 4.49\%$ )	94.44% ( $\pm 2.24\%$ )

## C ODIN Results

One critical component of our system is its ability to work on zero day exploits. Thus, OOD methods that require outlier exposure would fail in this set up. Furthermore, methods that require some level of parameter tuning generally fall out of scope for the same reason. Thus, we use the parameter-less version of energy-based detection. Nevertheless, we test ODIN (Liang et al., 2017a) as our model safeguard as well and provide those results here. Since we do not allow for parameter tuning, we use the default parameters ( $T = 1000, \epsilon = 0.05$ ) from the `pytorch-ood` library (Kirchheim et al., 2022).

We find a large number of identical output values from the ODIN detector. Thus, the values for TPR (TNR 95%) or TPR (TNR 85%) are skewed low. That is, the threshold value encompasses a much larger percentage of true negatives (since in some cases over 20% of ID output values are identical). We include the values even still for comparison in Table 17. However, we note that this would make such a safeguard unusable in practice because the threshold would be too hard to tune.

Table 17: *ODIN results*. We use ODIN as a model safeguard with default parameters from the `pytorch-ood` library.

Architecture	Layer	AU ROC ( $\uparrow$ )	TPR (TNR 95%) ( $\uparrow$ )	TPR (TNR 85%) ( $\uparrow$ )
<b>FTP-Patator</b>				
CNN	Output	0.4881 ( $\pm 0.2614$ )	0.00% ( $\pm 0.00\%$ )	0.05% ( $\pm 0.13\%$ )
FNN	Output	0.9193 ( $\pm 0.0265$ )	0.00% ( $\pm 0.00\%$ )	73.05% ( $\pm 38.27\%$ )
Transformer	Output	0.5633 ( $\pm 0.1748$ )	0.00% ( $\pm 0.00\%$ )	0.00% ( $\pm 0.00\%$ )
<b>Infiltration</b>				
CNN	Output	0.7131 ( $\pm 0.1616$ )	11.14% ( $\pm 14.62\%$ )	29.59% ( $\pm 30.75\%$ )
FNN	Output	0.5441 ( $\pm 0.1439$ )	0.00% ( $\pm 0.00\%$ )	6.22% ( $\pm 6.42\%$ )
Transformer	Output	0.6231 ( $\pm 0.1990$ )	0.00% ( $\pm 0.00\%$ )	10.68% ( $\pm 21.46\%$ )
<b>SSH-Patator</b>				
CNN	Output	0.3665 ( $\pm 0.2301$ )	0.85% ( $\pm 1.98\%$ )	4.56% ( $\pm 3.79\%$ )
FNN	Output	0.7255 ( $\pm 0.0555$ )	0.00% ( $\pm 0.00\%$ )	26.86% ( $\pm 15.71\%$ )
Transformer	Output	0.7277 ( $\pm 0.1070$ )	2.65% ( $\pm 7.94\%$ )	8.13% ( $\pm 18.65\%$ )
<b>CIC-IDS-2017</b>				
CNN	Output	0.4914 ( $\pm 0.0251$ )	0.72% ( $\pm 2.15\%$ )	8.44% ( $\pm 11.05\%$ )
FNN	Output	0.5687 ( $\pm 0.0486$ )	0.00% ( $\pm 0.00\%$ )	7.28% ( $\pm 7.35\%$ )
Transformer	Output	0.3072 ( $\pm 0.0586$ )	0.83% ( $\pm 1.66\%$ )	7.11% ( $\pm 5.73\%$ )
<b>UNSW-NB15</b>				
CNN	Output	0.6321 ( $\pm 0.0251$ )	22.50% ( $\pm 18.47\%$ )	24.91% ( $\pm 20.40\%$ )
FNN	Output	0.4167 ( $\pm 0.0565$ )	0.00% ( $\pm 0.00\%$ )	12.59% ( $\pm 6.48\%$ )
Transformer	Output	0.4422 ( $\pm 0.0694$ )	1.84% ( $\pm 3.87\%$ )	5.67% ( $\pm 7.41\%$ )