

ROULETTE: A Neural Attention Multi-Output Model for Explainable Network Intrusion Detection

Giuseppina Andresini^{a,*}, Annalisa Appice^{a,b}, Francesco Paolo Caforio^c,
Donato Malerba^{a,b}, Gennaro Vessio^a

^a*Department of Computer Science, University of Bari Aldo Moro, Bari, Italy*

^b*Consorzio Interuniversitario Nazionale per l'Informatica - CINI, Bari, Italy*

^c*IISS Majorana, Martina Franca, Italy*

Abstract

¹ Network Intrusion Detection (NID) systems are one of the most powerful forms of defense for protecting public and private networks. Most of the prominent methods applied to NID problems consist of Deep Learning methods that have achieved outstanding accuracy performance. However, even though they are effective, these systems are still too complex to interpret and explain. In recent years this lack of interpretability and explainability has begun to be a major drawback of deep neural models, even in NID applications. With the aim of filling this gap, we propose ROULETTE: a method based on a new neural model with attention for an accurate, explainable multi-class classification of network traffic data. In particular, attention is coupled with a multi-output Deep Learning strategy that helps to discriminate better between network intrusion categories. We report the results of extensive experimentation on two benchmark datasets, namely NSL-KDD and UNSW-NB15, which show the beneficial effects of the proposed attention mechanism and multi-output learning

*Corresponding author

Email addresses: giuseppina.andresini@uniba.it (Giuseppina Andresini),
annalisa.appice@uniba.it (Annalisa Appice),
francescopaolo.caforio@posta.istruzione.it (Francesco Paolo Caforio),
donato.malerba@uniba.it (Donato Malerba), gennaro.vessio@uniba.it (Gennaro Vessio)

¹This version of the contribution has been accepted for publication, after peer review but is not the Version of Record and does not reflect post-acceptance improvements, or any corrections. The Version of Record is available online at: <https://doi.org/10.1016/j.eswa.2022.117144>. Accepted Version is subject to the publisher's Accepted Manuscript terms of use <https://www.elsevier.com/about/policies-and-standards/copyright>

strategy on both the accuracy and explainability of the decisions made by the method.

Keywords: network intrusion detection, multi-class classification, deep learning, attention, explainable artificial intelligence, multi-output learning

1 **1. Introduction**

2 Over the past decade, the predominance of Deep Learning in intrusion de-
3 tection has been repeatedly assessed in the cybersecurity literature (Berman
4 et al. 2019, Naseer et al. 2018). In particular, several methods based on Deep
5 Learning, such as Autoencoders (Naseer et al. 2018), Recurrent Neural Networks
6 (Folino et al. 2021), Long Short-Term Memory networks (Sovilj et al. 2020, Yin
7 et al. 2017), Generative Adversarial Networks (Yang et al. 2019, Andresini, Ap-
8 pice, De Rose & Malerba 2021), Convolutional Neural Networks (Al-Turaiki &
9 Altwaijry 2021, Andresini, Appice & Malerba 2021*b*), Siamese and Triplet net-
10 works (Bedi et al. 2020*b*, Andresini, Appice & Malerba 2021*a*), have recently
11 contributed to introducing advanced classification capabilities into Network In-
12 trusion Detection (NID) systems, sustaining the resilience of the security line of
13 private and public networks.

14 However, Deep Learning techniques train classification models that are typ-
15 ically “black-boxes”. Indeed, these models are implicitly represented in numer-
16 ical form as synaptic weights in the network and, in general, it is difficult, if
17 not impossible, to interpret these weights, due to the complexity of the net-
18 work structure. The opacity of these black-boxes was acceptable as long as the
19 dominant criterion for assessing the quality of NID systems was their accuracy,
20 measured in terms of standard evaluation metrics (e.g., f-score and predictive
21 accuracy). In some classification problems, accurate black-box models may be
22 acceptable. However, even though the priority of modern NID systems today
23 is still to provide accurate classifications of network traffic, easier-to-explain
24 models are becoming increasingly desirable in NID applications. In fact, “ex-
25 planations” can provide measurable factors on which characteristics influence

26 the prediction of a cyber-attack and to what extent (Mane & Rao 2021). These
27 factors can give the security analyst much better insight into why an alert was
28 reported. Furthermore, if the explanation corresponds to domain knowledge,
29 the analyst can easily and confidently approve it. Therefore, translating the
30 model outcome into feature contribution and analyzing the impact of certain
31 traffic characteristics can increase stakeholder confidence (Wali & Khan 2021).

32 The explainability of classifications refers partially to the human interpretabil-
33 ity of the processes underlying the decisions given by a classification model.
34 The more explainable a certain model is, the easier it will be for a human to
35 understand or explain the underlying reasoning. Therefore, explaining intru-
36 sion classifications can help turn predictions into actions to better achieve the
37 resilience of the network defense security line. We note that this need for ex-
38 plainable alerts also matches the emerging EU vision, which is extending the
39 “right to explanation” formulated by the GDPR to solutions based on Artificial
40 Intelligence, and especially Deep Learning (Sartor & Lagioia 2020).

41 *Attention* is an increasingly popular Deep Learning technique, oriented to-
42 wards the design of explainable deep neural models. This mechanism allows the
43 adaptive selection of the input where the network “sees” the most important
44 information. The attention mechanism was introduced to improve the accuracy
45 performance of the encoder-decoder model for machine translation (Bahdanau
46 et al. 2015). It has also been exploited in Computer Vision systems to improve
47 their performance for a variety of tasks, ranging from image captioning to visual
48 question answering (Guo et al. 2021, Komodakis & Zagoruyko 2017). However,
49 analyzing the results of the network attention branches could also provide in-
50 sight into how the black-box model works, by contributing to the enhancement
51 of the explainability of its decisions. Despite this relationship between the at-
52 tention mechanism and the challenge of explainability, recent network security
53 studies that have tested attention in NID applications have focused solely on the
54 gain in accuracy achieved (Liu et al. 2020, Tang et al. 2020, Zhao et al. 2022).
55 In addition, these studies apply the attention mechanism with feature-vector
56 representations of network traffic flow traces. On the other hand, based on hu-

57 man visual attention, attention mechanisms have recently proved very effective
58 in image classification applications (Joshi et al. 2021).

59 The present study is boosted by the interesting results recently presented
60 in (Andresini, Appice, De Rose & Malerba 2021, Andresini, Appice & Malerba
61 2021*b*), which assessed how imagery representations of network traffic data may
62 be adopted to accurately separate intrusions from normal traffic. Using this idea
63 as a springboard, we propose a new NID system, called ROULETTE (neuRal at-
64 tentiOn MULti-Output ModEl for explainable InTrusion DeTEction), which
65 applies a Convolutional Neural Network (CNN) with an attention mechanism
66 to images converted from flow characteristics of network traffic data. The main
67 contribution of this study is the evaluation of the effectiveness of the attention
68 mechanism in the CNN classification of network flow traces. Furthermore, con-
69 trary to (Andresini, Appice, De Rose & Malerba 2021, Andresini, Appice &
70 Malerba 2021*b*), where the classification was binary, i.e. normal vs. intrusion, in
71 this study the classification is multi-class, i.e. we separate intrusions from nor-
72 mal samples, but we also recognize the attack category. The effectiveness of the
73 proposed neural model with attention is quantified in terms of accuracy of the
74 classifications, as well as transparency of the decisions. In fact, thanks to the at-
75 tention mechanism, we are able to produce the attention map of a network flow
76 trace classification, and this map is expected to specify the flow characteristics
77 that are most relevant for classification. This distinction of input traffic data
78 allows us to identify the specific flow characteristics of intrusion categories and
79 can provide useful information for cyber-defenders with no prior knowledge. In
80 particular, we can consider such information as a hint to reduce the workload
81 in manual analysis.

82 An additional contribution is the improvement of the performance of the
83 classification model by using a multi-output architecture that is trained to si-
84 multaneously predict both the binary output and the multi-class output of a
85 network flow trace. The proposed multi-output architecture has two branches
86 that produce fully-connected heads at the end of the network: the branch with
87 the binary head allows us to learn features that are useful to separate normal

88 traffic from intrusions, which helps the branch with the multi-class head to
89 recognize different categories of intrusions.

90 In short, this paper provides the following contributions:

- 91 • The definition of an innovative neural methodology for NID applications,
92 which integrates the attention mechanism to achieve both accuracy and
93 transparency in multi-class classifications.
- 94 • The formulation of a multi-output architecture to predict the intrusion
95 category of any new network flow trace by taking advantage of features
96 learned from the binary classification of normal traffic data vs. intrusions.
- 97 • The presentation of the results of an extensive evaluation that investigates
98 the feasibility of the proposed learning components in the multi-class sce-
99 nario, as well as the ability of our methodology to obtain accuracy com-
100 parable to competitive, Deep Learning-based approaches taken from the
101 recent literature on NID systems.
- 102 • The exploration of the effect of several properties of explanations (i.e.,
103 compactness, robustness and separability), produced through the atten-
104 tion mechanism, on accuracy, and the analysis of particular information
105 disclosed by the produced explanations on the flow characteristics of spe-
106 cific attack categories.

107 This paper is organized as follows. Related works are presented in Section 2.
108 The proposed multi-output neural network with attention is described in Sec-
109 tion 3. The experimental setup is described in detail in Section 4. The results
110 of the evaluation of the proposed method are discussed in Sections 5 and 6,
111 regarding accuracy and explainability, respectively. Finally, Section 7 refocuses
112 on the purpose of our research, draws conclusions and illustrates possible future
113 developments.

114 2. Related Work

115 Recent trends in cybersecurity research have classified Deep Learning as a
116 prominent Artificial Intelligence paradigm for addressing NID problems. In this
117 paper we renew a Deep Learning-based approach that integrates a neural atten-
118 tion mechanism for multi-class classification of network flow traces. Therefore,
119 we mainly focus the literature overview on recent studies applying Deep Learn-
120 ing (see Section 2.1) and eXplainable Artificial Intelligence - XAI (see Section
121 2.2) to classify flow-based network traffic data.

122 2.1. Deep Learning

123 Several recent studies have investigated the performance of various deep
124 neural network architectures for multi-class classification of network flow traces.
125 Most of these studies conducted experimental studies using benchmark, network
126 flow-based datasets such as NSL-KDD (Tavallae et al. 2009) or UNSW-NB15
127 (Moustafa & Slay 2015). In particular, the empirical study illustrated in (Ka-
128 songo & Sun 2020) experimented various architectures, showing how deep neu-
129 ral networks can gain accuracy over various traditional classifiers (e.g., SVM,
130 KNN and Logistic Regression). Various deep neural network architectures were
131 tested in (Vinayakumar et al. 2019) also in combination with feature selection
132 analysis. Following this research direction, the accuracy performance of var-
133 ious multi-class, fully-connected, recurrent and convolutional neural network
134 architectures is compared in (Gao et al. 2020). Moreover, this study couples
135 Deep Learning-based classification with association rule discovery. Specifically,
136 the association rules with the “normal” class in the consequent are applied to
137 match network flow traces classified as malicious, and (possibly) to filter out the
138 misclassified normal network traffic. In (Gao et al. 2019) a multi-class NID ap-
139 proach is formulated by combining Incremental-Extreme Learning Machine and
140 Adaptive Principal Component Analysis. In (Al-Turaiki & Altwaijry 2021) a 2D
141 representation of network flow traces is adopted, in order to train a multi-class
142 CNN. As a variant of the proposed pipeline, a so-called Deep Feature Synthesis
143 is also introduced, in order to complete a feature engineering step.

144 Furthermore, some recent studies have explored the performance of deep
145 metric learning architectures (e.g., Siamese networks or Triplet networks) in
146 multi-class NID applications. For example, the multi-class problem is studied
147 in (Bedi et al. 2020*b*), where a Siamese network is trained both on pairs of similar
148 samples (belonging to the same class) and pairs of dissimilar network flow traces
149 (belonging to opposite classes), to classify the intrusion class trace. A testing
150 network flow trace is classified according to a distance score computed for each
151 class: the class with the best score is predicted. An ensemble consisting of a
152 Siamese network, a classic deep network and an XGBoost binary classifier is de-
153 scribed in (Bedi et al. 2020*a*), where the goal is to separate normal network flow
154 traces from intrusions. A multi-class XGBoost classifier is then used to classify
155 traces detected as intrusions into different attack classes. In (Andresini, Appice
156 & Malerba 2021*a*) both the one-versus-all and the one-versus-one combination
157 strategy are applied to Triplet networks originally trained for binary classifica-
158 tion, in order to perform multi-class classification of network flow traces.

159 Finally, the multi-class classification has recently attracted attention also
160 in Adversarial Deep Learning. In this setting, the main objective is a classi-
161 fier that makes mistakes by making small changes to the training data. This
162 idea is investigated in (Caminero et al. 2019), where Variational Generative Au-
163 toencoders are experimented in an Adversarial Deep Reinforcement Learning
164 approach, formulated for multi-class classifications of network flow traces. A
165 specific conditional Variational Autoencoder architecture is also described in
166 (Lopez-Martin et al. 2017) for multi-class classification. This architecture in-
167 tegrates intrusion labels within the decoder layers. In particular, this study
168 classifies network flow traces according to the intuition that the autoencoder
169 best learns how to recover the original features when it processes the correct
170 label as input. Therefore, a testing network flow trace can be assigned to the
171 label that yields the recovered features that are closest to those observed.

172 *2.2. eXplainable Artificial Intelligence*

173 eXplainable Artificial Intelligence, or XAI, is a sub-field of Artificial Intel-
174 ligence that aims to enable humans to understand the decisions of artificial
175 systems by producing more explainable models, while maintaining a good level
176 of predictive accuracy (Lakkaraju et al. 2019). Significant interest in the XAI
177 research community has recently been observed in the development of “post-
178 hoc” explanations, in which an XAI technique can be applied to already trained
179 black-box Deep Learning models. Alternatively, an XAI technique can be incor-
180 porated into the Deep Learning algorithm. Today, several XAI techniques have
181 already been tested in many real-world applications such as business decision,
182 process optimization, medical diagnosis and investment recommendation, in or-
183 der to improve the reliability, transparency and fairness of Deep Learning-based
184 decisions (Xu, Uszkoreit, Du, Fan, Zhao & Zhu 2019, Antwarg et al. 2021).

185 Even in the cybersecurity field, security practitioners have begun to complain
186 about the black-box nature of Deep Learning-based decisions. The recent study
187 of Warnecke et al. (2020) started the investigation into how post-hoc XAI tech-
188 niques can be applied to produce explanations for the decisions of deep neural
189 networks, trained for both malware detection and vulnerability discovery ap-
190 plications. Specifically, this study compares several post-hoc XAI techniques
191 regarding the accuracy of explanations, as well as security-focused aspects, such
192 as completeness, efficiency and robustness. Notably, this study led to very re-
193 cent studies exploring the performance of various post-hoc XAI techniques also
194 in NID applications of Deep Learning.

195 In (Burkart et al. 2021) a surrogate model is trained to produce explanations
196 of binary decisions produced in NID applications. The surrogate model is a de-
197 cision tree, trained from network flow traces sampled around a counterfactual-
198 based local decision boundary. A surrogate model is also coupled to a deep
199 neural network in (Szczepański et al. 2020), in order to explain its black-box
200 decision. In (Sarhan et al. 2021) the transparency of a NID system is im-
201 proved by using SHAP (Lundberg & Lee 2017) to identify the input features
202 that contribute most to binary decisions produced through a neural network

203 black-box. SHAP is also used in (Wang et al. 2020) to perform analyses of the
204 most relevant features for detecting each category of intrusion. In (Andresini,
205 Pendlebury, Pierazzi, Loglisci, Appice & Cavallaro 2021) the analysis of the
206 feature relevance is conducted via DALEX (Biecek 2018). This explanation
207 analysis aims to explain how a binary model changes over a stream of network
208 flow traces to fit new attack categories. In (Marino et al. 2018) an adversar-
209 ial learning approach is experimented in order to understand why some of the
210 network flow traces are mis-classified by a deep neural network. The explana-
211 tions are formulated in terms of the minimal modifications required to change
212 the output of the black-box for any mis-classified network trace. Finally, in
213 (Caforio et al. 2021) Grad-CAM (Selvaraju et al. 2020) is applied to produce
214 gradient-based visual explanations of the CNN binary classification of imagery-
215 represented network flow traces. The main innovation of this study is that it
216 combines Grad-CAM explanations with the nearest-neighbour search, in order
217 to improve the accuracy of the CNN decisions in a post-hoc way, as well as to
218 increase the transparency of the CNN black-box decisions.

219 On the other hand, a few recent NID studies have also begun to investigate
220 the effect of intrinsic XAI techniques, such as the attention mechanism incorpo-
221 rated into the neural architectures and used instead of post-hoc explanations.
222 In (Liu et al. 2020) a Bidirectional Gated Recurrent Unit network with hierar-
223 chical attention is trained for the binary classification of network flow traces,
224 while a Temporal Convolutional Network is trained with the attention mecha-
225 nism in (Zhao et al. 2022) for the multi-class classification. Finally, a pipeline
226 with a Stacked Autoencoder and an attention mechanism is experimented in
227 (Tang et al. 2020). The use of attention in current NID studies has so far relied
228 only on feature vectors. Furthermore, although probing into the results of at-
229 tention branches of a deep neural network should provide insight into how the
230 black-box model works, current NID studies do not explore the improvement in
231 explainability achieved with attention, since they are focused solely on the gain
232 achieved in accuracy. The present work aims to address these issues.

233 3. Proposed Method

234 We assume that a dataset $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ of N training samples is avail-
235 able, where $\mathbf{x} \in \mathbb{R}^d$ is a d -dimensional vector of features (such as the number
236 of packets transmitted and the number of failed logins) that characterize flow
237 traces of historical network traffic, whereas $y \in \{1, \dots, K\}$ is the target vari-
238 able with K distinct classes: *normal* traffic data and various types of *intrusion*,
239 depending on those historically detected and labeled. The proposed intrusion
240 detection model, schematized in Fig. 1, is mainly based on three components:

- 241 • The reformulation of the network traffic classification task as an image
242 classification problem, which makes it possible to take advantage of con-
243 volution filters to learn a new discriminating representation.
- 244 • The use of an attention mechanism followed by an *average* layer, which
245 allows the extraction of an easy-to-explain attention map of the classifica-
246 tions.
- 247 • A multi-output learning strategy that allows the network to solve a binary,
248 auxiliary task, intended to improve the performance of the main multi-
249 class classification task.

250 In the pre-processing step image encoding is performed to transform the
251 feature vector \mathbf{x} of each network flow trace into a single-channel square image
252 $\mathbf{X} \in \mathbb{R}^{m \times m}$. This transformation is done by assigning each flow feature of \mathbf{x} to a
253 pixel frame of \mathbf{X} . Flow features are assigned to neighbouring pixels in the image,
254 depending on their correlation, and are not simply stacked in an arbitrary order.
255 A detailed description of how feature-pixel assignment is performed is reported
256 in (Andresini, Appice, De Rose & Malerba 2021). Thanks to this encoding,
257 we are able to reformulate the classical classification as an image classification
258 problem. Note that this reformulation has already proved useful in (Andresini,
259 Appice, De Rose & Malerba 2021, Caforio et al. 2021), although the previous
260 studies handled the binary classification task (to separate normal flow traces
261 from intrusions, regardless of the specific intrusion category). A convolution

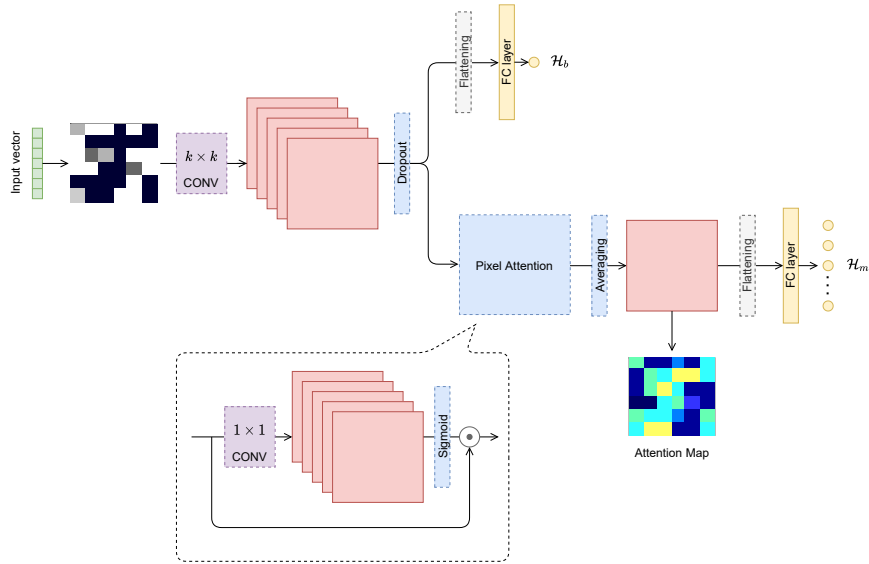


Figure 1: Schema of ROULETTE. Abbreviations: CONV = convolution; FC = fully-connected.

262 layer is then used to learn filters aimed at minimizing the classification error
 263 made in a backpropagated way. It is worth noting that, since the resolution of
 264 the input image is relatively low, and its content lacks the complexity of real
 265 photo-realistic images, a hierarchy of learned features is not needed, so a single
 266 convolutional layer is sufficient. For the same reason, this convolutional layer is
 267 followed by a dropout layer, which is used to mitigate early overfitting.

268 The *conv + dropout* block is then split into two branches. The first branch is
 269 responsible for the main multi-class classification task, aimed at discriminating
 270 network traffic data in multiple categories from a normal to a specific type of
 271 intrusion. Since we are interested in “explaining” the multi-class classifications
 272 made by the model, this branch includes a simple attention mechanism. This is
 273 analogous to the *pixel-attention* layer recently introduced in (Zhao et al. 2020),
 274 which is used to generate pixel-wise attention maps from the input volume.

275 More precisely, given an input volume X^{l-1} of size $H \times W \times C$ (height \times
 276 width \times channels), the *pixel-attention* layer convolves a 1×1 convolution filter,

277 followed by a sigmoid activation. Formally:

$$\mathbf{X}^l = \sigma(\text{CONV}(\mathbf{X}^{l-1})) \odot \mathbf{X}^{l-1}, \quad (1)$$

278 where CONV is a 1×1 convolution, σ is the sigmoid function and \mathbf{X}^l is the
279 resulting output tensor at the new layer l . This mechanism basically serves to
280 generate attention coefficients for all pixels in each feature map, thus “weight-
281 ing” their contribution to the final classification. In fact, we first use C point-
282 wise filters so that all pixels are weighted; the resulting feature maps are then
283 squashed between 0 and 1 by the sigmoid function and multiplied element-wise
284 with the input tensor, effectively producing a new tensor \mathbf{X}^l of the same shape
285 $H \times W \times C$.

286 Furthermore, in order to be able to produce a single attention map that
287 explains the importance of the features in the original image, we propose using
288 an *average* layer which basically averages corresponding pixels in the different
289 feature maps. More precisely, given the tensor of shape $H \times W \times C$ produced by
290 the *pixel-attention*, the *average* layer produces a single-channel matrix of size
291 $H \times W$, for which each attention pixel α_{ij}^l is equal to:

$$\alpha_{ij}^l = \frac{1}{C} \sum_{c=1}^C \alpha_{ijc}^{l-1}, \quad (2)$$

292 where α_{ijc}^{l-1} is the attention pixel at position ij in the c -th feature map from the
293 preceding $l-1$ layer. The *average* layer directly provides a heatmap, highlighting
294 the input image regions the model attended on to perform the classification. In
295 other words, a “visual explanation” is obtained. It is worth noting that since
296 the mechanism is learnable and embedded into an end-to-end trainable model,
297 this is in contrast to “post-hoc” visual explanation methods, such as the popular
298 Grad-CAM (Selvaraju et al. 2017), which can only be applied after the training
299 process has ended. A learnable attention mechanism is desirable as it helps the
300 model focus on key parts during training. Furthermore, it does not require an
301 expensive two-step process to derive the heatmap. In addition, this layer acts
302 as a simple regularizer, as the feature vector to propagate would be scaled by a
303 factor of C , and it also reduces the computational cost.

304 The output of the *average* layer is finally flattened and fed into a fully-
 305 connected layer, plus an output layer with a softmax activation attached and
 306 as many units as there are classes to predict. This branch aims to minimize a
 307 cross-entropy loss function:

$$\mathcal{H}_m = - \sum_{k=1}^K y_k \log \hat{y}_k, \quad (3)$$

308 where y_k is the ground truth label, while \hat{y}_k is the network output for a single
 309 data sample.

310 Inspired by multi-output learning (Xu, Shi, Tsang, Ong, Gong & Shen 2019),
 311 we propose the use of a second branch that aims to perform a simpler binary
 312 classification between normal traffic and intrusions, regardless of the specific
 313 type of attack. The second predicted output is intended as an auxiliary classi-
 314 fication objective, aimed at supporting the main multi-class classification task.
 315 In fact, since the two branches share the same first convolutional layer and their
 316 classification heads are optimized simultaneously by backpropagation, in this
 317 way the network is forced to learn features that are useful to separate normal
 318 traffic from attacks. This helps the main branch to focus on how to discrimi-
 319 nate better between different types of attack that can be very similar and have
 320 overlapping characteristics. Multi-output learning is not uncommon in neural
 321 networks and has proven effective in some other domains, e.g. (Cao et al. 2018,
 322 Castellano et al. 2020). Specifically, the output of the *conv + dropout* block
 323 is directly flattened and given as input to a fully-connected layer followed by a
 324 single sigmoid-activated output neuron. This second output serves to minimize
 325 a binary cross-entropy:

$$\mathcal{H}_b = -(y \log \hat{y} + (1 - y) \log (1 - \hat{y})). \quad (4)$$

326 Note that this branch does not include an attention mechanism, as the goal
 327 is to learn how to explain multi-class classification. Overall, the network learns
 328 to minimize, by backpropagation, the joint loss:

$$\mathcal{L} = \lambda \mathcal{H}_m + (1 - \lambda) \mathcal{H}_b, \quad (5)$$

329 where $\lambda \in [0, 1]$ is a hyper-parameter that balances the contribution of the
330 individual losses.

331 4. Experimental Setup

332 In this section we describe the datasets used for evaluating the accuracy
333 and explainability of ROULETTE, i.e. NSL-KDD (Tavallae et al. 2009) and
334 UNSW-NB15 (Moustafa & Slay 2015), and implementation details.

335 4.1. Dataset Description

336 The NSL-KDD dataset (Tavallae et al. 2009)² comprises normal network
337 flow traces and four categories of attack: Denial of Service (DoS), User to Root
338 (U2R), Remote to Local (R2L) and Probing attack. The training set is made
339 up of 21 different attack sub-categories, while the test set is composed of 37
340 different attack sub-categories. This means there are 16 novel attacks in the
341 test set. Each trace in the NSL-KDD dataset has 41 features, and detailed
342 descriptions of these features are provided in (Tavallae et al. 2009). We note
343 that both U2R and R2L are rare attacks. While this dataset may not represent
344 perfectly existing real-world networks, recent, state-of-the-art studies still use it
345 as an effective benchmark dataset to help researchers compare different multi-
346 class classification NID methods.

347 The UNSW-NB15 (Moustafa & Slay 2015),³ on the other hand, includes re-
348 alistic modern normal activities and synthetic contemporary attack behaviours
349 extracted from network traffic monitored in 2015. Both the training set and
350 the testing set contain normal network flow traces and nine categories of at-
351 tacks: Fuzzers, Analysis, Backdoors, DoS, Exploits, Generic, Reconnaissance,
352 Shellcode and Worms. Each trace in the UNSW-NB15 dataset has 49 features:
353 detailed descriptions of these features are provided in (Moustafa & Slay 2015).

²<https://www.unb.ca/cic/datasets/nsl.html>

³<https://research.unsw.edu.au/projects/unsw-nb15-dataset>

354 The number of samples collected both in the training and testing set for
 355 each category per dataset is reported in Table 1. A complete description of the
 356 flow characteristics enclosed in both the NSL-KDD and UNSW-NB15 datasets
 357 is reported in (Choudhary & Kesswani 2020).

Table 1: Number of network flow traces per class type in both the training set and testing set of NSL-KDD and UNSW-NB15.

Class type	NSL-KDD		UNSW-NB15	
	Train	Test	Train	Test
Normal	67343	9711	56000	37000
DoS	45927	7458	12264	4089
Probe	11656	2421	-	-
R2L	995	2754	-	-
U2R	52	200	-	-
Fuzzers	-	-	18184	6062
Analysis	-	-	2000	677
Backdoors	-	-	1746	583
Exploits	-	-	33393	11132
Generic	-	-	40000	18871
Reconnaissance	-	-	10491	3496
Shellcode	-	-	1133	378
Worms	-	-	130	44
Total	125973	22544	175341	82332

Table 2: Hyper-parameter search space

Hyper-parameter	Values
Mini-batch size	$\{2^5, 2^6, 2^7, 2^8, 2^9\}$
Learning rate	$[0.0001, 0.001]$
Dropout rate	$[0, 1]$
# of filters	$\{2^5, 2^6, 2^7, 2^8, 2^9\}$
Kernel size	$[2, 4]$
# of neurons per hidden layer	$\{2^5, 2^6, 2^7, 2^8, 2^9\}$
λ	$[0.5, 1]$

358 4.2. Implementation Details

359 We developed ROULETTE in Python 3, using the high-level neural net-
 360 work API Keras 2.4 integrated in TensorFlow (Abadi et al. 2015).⁴ In the
 361 pre-processing step, the categorical input features were mapped into numeri-
 362 cal features using the one-hot-encoder strategy, and then the numerical features
 363 were scaled using the min-max normalization.

364 For each dataset, we optimized the hyper-parameters of the architecture us-
 365 ing the tree-structured Parzen estimator algorithm (Bergstra et al. 2011) as im-
 366 plemented in the Hyperopt library (Bergstra et al. 2013). The hyper-parameter
 367 optimization was performed using a random stratified split of 20% of the en-
 368 tire training as a validation set, following the Pareto principle. We selected
 369 the hyper-parameter configuration that achieved the lowest validation loss. The
 370 hyper-parameter search space is reported in Table 2. The configuration was
 371 completed by the commonly used ReLU (Glorot et al. 2011) as the activation
 372 function for each hidden layer.

373 We trained the network with mini-batches using back-propagation, and the
 374 gradient-based optimization was performed using the Adam update rule (Kingma

⁴<https://github.com/gsndr/ROULETTE>.

375 & Ba 2014). The weights were initialized following the Xavier scheme (Glorot
376 & Bengio 2010). In addition, a maximum number of epochs equal to 150 was
377 set, and an early stopping approach based on the lowest loss on the validation
378 set (the same set used for the hyper-parameter optimization) was used, in order
379 to retain the best classification models.

380 **5. Accuracy Performance Analysis**

381 In this section we show the results of an analysis aimed at evaluating the
382 accuracy performance of ROULETTE, in order to answer the following questions:

383 Q1 How does the accuracy of the proposed multi-output Deep Learning strat-
384 egy change by varying a few dimensions of the neural network architecture,
385 e.g. dropout layer, regularization penalty term and aggregation operation?

386 Q2 Is the proposed multi-output Deep Learning strategy able to achieve higher
387 accuracy than the single-output strategy?

388 Q3 How does the attention mechanism help the accuracy performance of the
389 predictive model?

390 Q4 Does the defined neural attention multi-output model outperform state-
391 of-the-art NID systems?

392 The multi-class accuracy metrics measured for this analysis are described in
393 Section 5.1. The results of the sensitivity study to answer Q1 are reported in
394 Section 5.2. The results of the ablation study to answer Q2 and Q3 are illus-
395 trated in Section 5.3. Finally, the results of the comparative analysis performed
396 to answer Q4 are shown in Section 5.4.

397 *5.1. Performance Metrics*

398 We measured the predictive performance of the compared methods by com-
399 puting standard multi-class classification metrics. All compared metrics were
400 computed to evaluate performance on the classifications produced on the testing
401 sets. Specifically, we considered the following metrics:

- 402 • **Precision**, which measures the precision of the classification per class type,
403 i.e. how many network flow traces are correctly classified for a particular
404 class type k , given all predictions of that class, i.e. $P_k = \frac{tp_k}{tp_k + fp_k}$.
- 405 • **Recall**, which measures the recall of the classification per class type k ,
406 i.e. how many network flow traces are correctly classified for a particular
407 class type k given all occurrences of that class type, i.e. $R_k = \frac{tp_k}{tp_k + fn_k}$.
- 408 • **F1**, which measures the harmonic mean of Precision and Recall per class
409 type k , i.e. $F1_k = 2 \frac{P_k \times R_k}{P_k + R_k}$. The higher the F1 score per class type k ,
410 the better the balance between the Precision and Recall achieved by the
411 method in predicting network flow traces of class k .
- 412 • **Macro-F1**, which measures the average F1 score per class type k , i.e. $\text{Macro-F1} =$
413 $\frac{1}{K} \sum_{k=1}^K F1_k$.
- 414 • **Weighted-F1**, which measures the weighted mean of the F1 score per class
415 type k , i.e. $\text{Weighted-F1} = \sum_{k=1}^K w_k F1_k$, where each weight w_k is equal to
416 the probability of class k in the testing set, i.e. $w_k = \frac{n_k}{N}$.
- 417 • **A**, which measures the overall accuracy, i.e. the proportion of correctly
418 classified network flow traces in all the classified traces, i.e. $A = \frac{1}{N} \sum_{k=1}^K tp_k$.

419 In the above formulation: tp_k is the number of network flow traces of class
420 k that are correctly predicted as belonging to class k ; fp_k is the number of
421 network flow traces that are incorrectly classified as belonging to k ; tn_k is the
422 number of network flow traces that are correctly classified as not belonging to
423 class k ; fn_k is the number of network flow traces that are incorrectly classified
424 as not belonging to k ; n_k is the ground truth size of class k (i.e., the number of
425 traces labeled with class k); $N = \sum_k N_k$ is the total size of the classifications; K
426 is the number of distinct class types. In calculating the macro-metric reported
427 above, we gave equal weights (i.e., $\frac{1}{K}$) to each class type. In this way, we avoided
428 our evaluation offsetting the possible impact of unbalanced data learning. On

Table 3: Macro-F1, Weighted-F1 and A obtained using L1, L2 or DROPOUT. The best results are in bold.

Dataset	Architecture	Macro-F1	Weighted-F1	A
NSL-KDD	L1	0.566	0.753	0.781
	L2	0.562	0.747	0.777
	DROPOUT	0.613	0.790	0.815
UNSW-NB15	L1	0.374	0.731	0.727
	L2	0.414	0.760	0.751
	DROPOUT	0.424	0.767	0.764

429 the contrary, when calculating the weighted metric, the highly populated classes
 430 had a higher weight compared to the smaller ones.

431 5.2. Sensitivity Study

432 This sensitivity analysis was performed on the specific neural network ar-
 433 chitecture trained by ROULETTE, in order to explore how the accuracy of the
 434 trained neural network can be influenced by:

- 435 • The use of the dropout layer in place of the L1 or L2 regularization to
 436 prevent overfitting.
- 437 • The addition of a regularization penalty term to the objective function.
- 438 • The aggregation operation (i.e., AVG or MAX) adopted after the pixel at-
 439 tention layer; in other words, if it is better to take the maximum attention
 440 pixel per channel than to average all of the pixels.

441 Table 3 collects the Macro-F1, Weighted-F1 and A achieved by replacing the
 442 dropout layer of ROULETTE with either the L1 regularization or the L2 reg-
 443 ularization. The results show that the dropout layer is able to gain accuracy

Table 4: Macro-F1, Weighted-F1 and A of ROULETTE with L1 Penalty, L2 Penalty and No Penalty in the objective function. The best results are in bold.

Dataset	Architecture	Macro-F1	Weighted-F1	A
NSL-KDD	L1 Penalty	0.551	0.731	0.760
	L2 Penalty	0.541	0.723	0.754
	No Penalty	0.613	0.790	0.815
UNSW-NB15	L1 Penalty	0.353	0.721	0.711
	L2 Penalty	0.390	0.735	0.717
	No Penalty	0.424	0.767	0.764

444 when compared to its counterparts that use the L1 and L2 regularization, re-
 445 spectively. These results were expected as dropout has become the standard
 446 choice as a regularization technique in today’s architectures (Guo et al. 2019,
 447 Phaisangittisagul 2016). Studies have shown the effectiveness of dropout train-
 448 ing also in combination with convolutional layers (Alex Kendall & Cipolla 2017,
 449 Phaisangittisagul 2016, Srivastava et al. 2014).

450 Table 4 collects the Macro-F1, Weighted-F1 and A achieved by adding the L1
 451 penalty, L2 penalty and No penalty to the objective functions of ROULETTE
 452 described in Equations 3 and 4. The results show that the introduction of
 453 a regularizer to apply a penalty on the layer outputs does not lead to any
 454 improvement in the accuracy in either NSL-KDD or in UNSW-NB15. This can
 455 be explained by considering that the aforementioned dropout is already effective
 456 in countering overfitting.

457 Finally, Fig. 2 reports the Macro-F1, Weighted-F1 and A achieved by ROULETTE
 458 using attention with either an average layer (AVG) or a max layer (MAX). The
 459 AVG configuration, which is the baseline described in Equation 2, averages the
 460 pixels across all channels of the attention layer to produce a single attention
 461 map. The MAX configuration determines the maximum pixel on all channels.

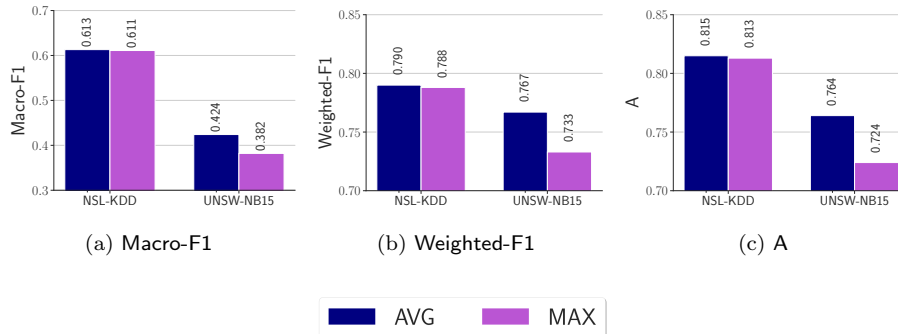


Figure 2: Macro-F1, Weighted-F1 and A of ROULETTE using attention with either an average layer (AVG) or a max layer (MAX).

462 The experimental results show that AVG outperforms MAX in both NSL-KDD
 463 and UNSW-NB15, although the difference between the performance is higher
 464 in UNSW-NB15 than in NSL-KDD. Our main intuition is that by averaging
 465 we can retain more pixel information across all the channels rather than simply
 466 selecting the maximum.

467 5.3. Ablation Study

468 The ablation study of ROULETTE was performed by evaluating the predic-
 469 tive performance of three architecture configurations identified as baselines of
 470 ROULETTE. These were in turn defined by discarding the branch with the bi-
 471 nary head or the attention mechanism from the training stage of ROULETTE.
 472 In particular, we considered the following baseline architectures:

- 473 • SO, which discards both the branch with the binary head and the attention
 474 mechanism. The deep neural network produces a single, multi-class output
 475 by minimizing the cross-entropy loss function.
- 476 • SO+A, which discards the branch with the binary head. The neural net-
 477 work produces a single, multi-class output with attention by minimizing
 478 the cross-entropy loss function.

Table 5: Macro-F1, Weighted-F1 and A both of ROULETTE and its baseline configurations SO, SO+A and MO. The best results are in bold.

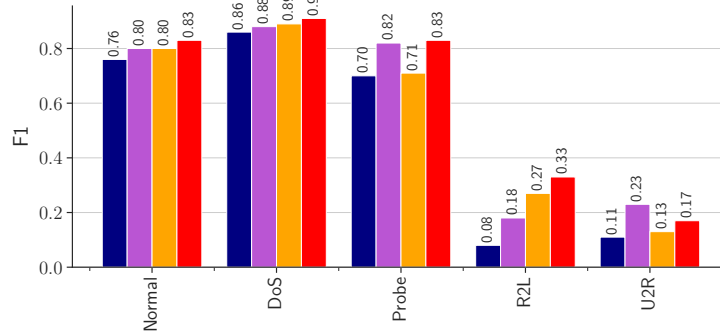
Dataset	Architecture	Macro-F1	Weighted-F1	A
NSL-KDD	SO	0.503	0.701	0.741
	SO+A	0.579	0.745	0.779
	MO	0.558	0.749	0.778
	ROULETTE	0.613	0.790	0.815
UNSW-NB15	SO	0.393	0.747	0.748
	SO+A	0.391	0.753	0.754
	MO	0.391	0.751	0.753
	ROULETTE	0.424	0.767	0.764

479 • MO, which discards the attention mechanism. The deep neural network
 480 produces both a multi-class and a binary output by minimizing the joint
 481 loss function.

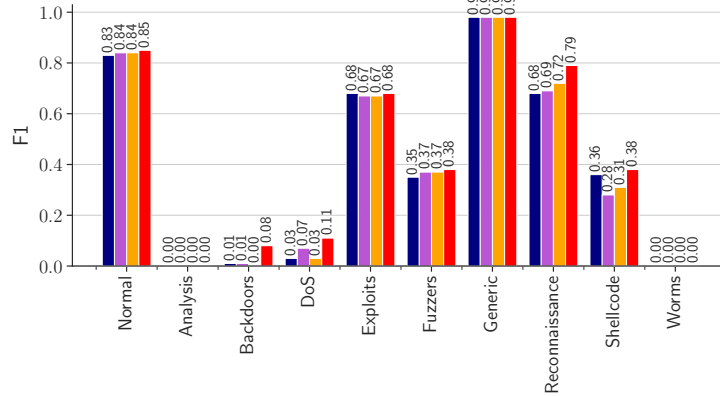
482 The results of Macro-F1, Weighted-F1 and A of SO, SO+A, MO and ROULETTE
 483 are reported in Table 5. They show that ROULETTE can take advantage of both
 484 the multi-output strategy and the attention mechanism achieving more accurate
 485 decisions than all its baselines.

486 Figure 3 reports the F1 scores computed per each class. These detailed
 487 results show that ROULETTE performs better than (or equal to) its baseline
 488 architectures in predicting all class categories of the multi-class problems con-
 489 sidered. Notably, ROULETTE is capable of achieving a good level of accuracy
 490 on various rare classes, such as R2L and U2R of NSL-KDD, and Shellcode of
 491 UNSW-NB15.

492 In general, this analysis shows the viability of our idea of exploiting the
 493 multi-output strategy, in order to obtain accuracy in the multi-class branch
 494 thanks to the knowledge learned in the auxiliary binary branch. In addition,



(a) NSL-KDD



(b) UNSW-NB15



Figure 3: F1 score per class of both ROULETTE and its baseline configurations SO, SO+A and MO.

495 the attention mechanism, which we have introduced in the neural network to
 496 see which input information is relevant for decisions, also allows ROULETTE to
 497 achieve higher levels of accuracy.

498 5.4. Competitor Analysis

499 The comparative analysis is performed to assess the significance of accuracy
 500 and novelty of ROULETTE compared to several competitors, selected from the

Table 6: A of ROULETTE vs. state-of-the-art, Deep Learning algorithms that perform multi-class classification.

Dataset	Approach	Description	A
NSL-KDD	Al-Turaiki & Altwaijry (2021)	CNN, Deep Feature Synthesis	0.814
	Andresini, Appice & Malerba (2021a)	Autoencoder, Triplet network, OVA	0.778
	Andresini, Appice & Malerba (2021a)	Autoencoder, Triplet network, OVO	0.766
	Bedi et al. (2020a)	Siamese network, Ensemble, XGBoost	0.79.9
	Bedi et al. (2020b)	Siamese network	0.769
	Cambrero et al. (2019)	Variational Generative Autoencoder	0.801
	Gao et al. (2019)	I-ELM, PCA	0.812
	Gao et al. (2020)	DNN	0.773
	Gao et al. (2020)	RNN	0.713
	Gao et al. (2020)	CNN	0.735
	Lopez-Martin et al. (2017)	Conditional Variation Autoencoder	0.801
	Tang et al. (2020)	Autoencoder, Attention, DNN	0.821
	Vinayakumar et al. (2019)	DNN	0.778
	Wang et al. (2020)	DNN, SHAP	0.803
	ROULETTE	CNN, Attention, Multi-output	0.815
UNSW-NB15	Al-Turaiki & Altwaijry (2021)	CNN, Deep Feature Synthesis	0.685
	Gao et al. (2019)	I-ELM, PCA	0.707
	Kasongo & Sun (2020)	DNN, Feature selection	0.756
	Vinayakumar et al. (2019)	DNN	0.660
	Zhao et al. (2022)	Temporal CNN, Attention	0.729
		ROULETTE	CNN, Attention, Multi-output

501 state of the art in NID literature. Note that the selected competitors differ
502 in the deep neural network architecture tested. The results of the competitors
503 are taken from the reference papers, as their code is not publicly available for
504 repeating the experiments. However, the comparison is safe as all methods have
505 been tested on the multi-class problem of the same training and testing sets
506 described in Section 4.1.

507 We point out that the competitors that integrate the attention mechanism
508 (Tang et al. 2020, Zhao et al. 2022) are closest to ROULETTE. Specifically,
509 the method described in (Tang et al. 2020), which has been tested on NSL-
510 KDD, integrates the attention layer into a deep neural network, trained with
511 the flow-based characteristics of the dataset encoded at the encoder level of an

512 autoencoder. On the other hand, the method described in (Zhao et al. 2022),
513 which was tested on UNSW-NB15, integrates the attention layer into a Temporal
514 CNN trained with the original flow-based characteristics of the dataset. Finally,
515 we note that the method described in (Wang et al. 2020) also experiments an
516 explanation mechanism. However, it uses post-hoc explanations based on SHAP,
517 which allow the author to achieve transparency of the Deep Learning decisions,
518 but it has no effect on the overall accuracy of these decisions.

519 For all methods in this comparative study we collect the A as this metric is
520 provided in all reference studies. The A results, reported in Table 6, show that
521 ROULETTE outperforms its competitors, including the attention-based com-
522 petitor evaluated on UNSW-NB15 (Zhao et al. 2022). The only exception is
523 the attention-based competitor defined in (Tang et al. 2020) which outperforms
524 ROULETTE on NSL-KDD. However, we note that the attention mechanism of
525 this competitor is trained on the encoded features of the dataset, instead of
526 the original input features. The encoded features commonly allow us to gain
527 accuracy in classification. This positive effect of autoencoders is also proved
528 in (Andresini, Appice & Malerba 2021b, Andresini et al. 2020) for binary for-
529 mulations of NID problems. However, training attention on encoded features
530 excludes the opportunity to explain the effect of input traffic features on deci-
531 sions.

532 **6. Explanation Property Analysis**

533 This analysis aimed to explore which properties of the intrinsic explana-
534 tions produced through the attention mechanism may be connected to their
535 ability to outperform post-hoc constructed counterparts. As post-hoc explana-
536 tions we considered the visual explanation maps produced through the popular
537 Grad-CAM technique (Selvaraju et al. 2017), after the neural training process
538 has been completed without attention. Moreover, we explored the relationship
539 between the observed properties of attention explanations and the accuracy
540 performance already investigated in Section 5. Specifically, this study explores

541 properties referred to as *compactness* (Section 6.2), *robustness* (Section 6.3) and
 542 *separability* (Section 6.4). The additional metrics considered to explore these
 543 explanation properties are introduced in Section 6.1.

544 6.1. Explanation Metrics

545 Two new metrics, namely *Inertia* and *Average link distance*, were measured for
 546 the analysis of compactness and robustness, respectively. The aforementioned
 547 Macro-F1, Weighted-F1 and A were measured for the analysis of separability.

548 *Inertia* measures the similarities of visual explanation maps clustered by class
 549 type. The lower the *Inertia*, the higher the ability to make transparent the
 550 common, intrinsic factors of the signature that is beyond the decisions produced
 551 for the network flow traces of the same class type. We measured *Inertia* as the
 552 averaged squared Euclidean distance computed for each visual explanation map
 553 to the visual explanation map centroid of its ground truth class, i.e. $Inertia =$
 554 $\frac{1}{K} \sum_{k=1}^K \frac{1}{|C_k|} \sum_{\mathbf{H}_j \in C_k} d(\mathbf{H}_j, \hat{\mathbf{H}}_k)$, where K is the number of classes, \mathbf{H}_j is a visual
 555 explanation map in class C_k and $\hat{\mathbf{H}}_k$ is the centroid of all visual explanation
 556 maps with class k .

557 *Average link distance* measures the average distance between the visual ex-
 558 planation maps produced for the same class type on both the training set and
 559 the testing set, respectively. For each class type k , we determined the mean of
 560 the Euclidean distances between each pair made up of the visual explanation
 561 map of a network flow trace of class k from the training set and the testing
 562 set, respectively, i.e. $d(k) = \frac{1}{|N_k| |N'_k|} \sum_{\mathbf{H}_j \in N_k, \mathbf{H}'_j \in N'_k} d(\mathbf{H}_j, \mathbf{H}'_j)$, where \mathbf{H}_j and \mathbf{H}'_j
 563 denote two network flow traces of class k belonging to the training set and test-
 564 ing set, while N_k and N'_k indicate the number of network flow traces of class k
 565 recorded in both the training set and the testing set, respectively. The lower
 566 the distance $d(k)$, the more robust the explanation of the learned classification
 567 model on the predictions produced for unseen network flow traces of class k .

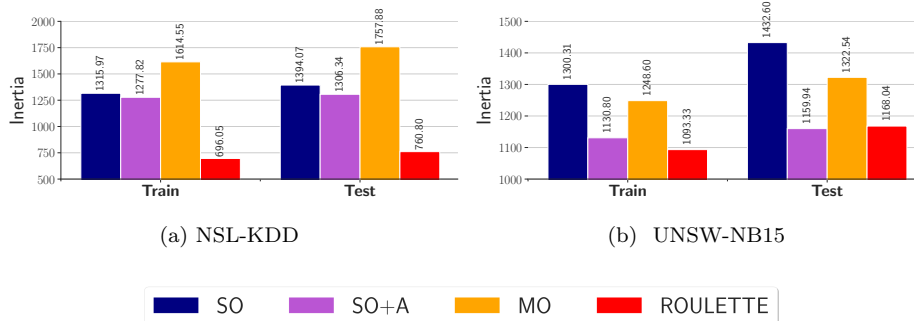


Figure 4: Inertia of the visual explanation maps produced with decisions yielded on the network flow traces of both the training set and the testing set, respectively.

568 6.2. Compactness

569 In principle, we expect a good multi-class classification model to be able to
 570 learn the distinctive signature of each class type. This signature should allow
 571 us to trigger the same decision process on network flow traces of the same class
 572 type. In this study we explored this ability related to the “compactness” of the
 573 decision explanations.

574 Moreover, we analyzed the Inertia of the visual explanation maps that were
 575 produced with attention in ROULETTE and SO+A, as well as with Grad-CAM
 576 in MO and SO. The results of Inertia, computed separately on the training and
 577 testing set of the performed experiments, are reported in Fig. 4. These re-
 578 sults show that the visual explanation maps produced through attention always
 579 achieve lower Inertia than their respective counterparts produced through Grad-
 580 CAM (i.e., the Inertia of ROULETTE is always lower than the Inertia of MO, just
 581 as the Inertia of SO+A is always lower than the Inertia of SO). This behaviour,
 582 which is observed equally in both the training set and the testing set, assesses
 583 that the explanations per class, produced with the intrinsic attention mecha-
 584 nism, are more compact than the explanations per class eventually produced
 585 post-hoc. In addition, we note that the attention mechanism, coupled with the
 586 multi-output Deep Learning strategy, commonly achieves the lowest Inertia in
 587 our study. The only exception is observed with the UNSW-NB15 testing set.

588 However, the difference between the best Inertia of SO+A and the runner-up
589 Inertia of ROULETTE is small (1159.94 vs. 1168.04). Therefore, this empirical
590 study suggests that a relationship may exist between the higher compactness
591 of the decision explanations produced by ROULETTE and the higher accuracy
592 (assessed in Section 5) of these decisions.

Table 7: Average link distance between pairs of visual explanation maps consisting of a network flow trace of the same class from the training set and the testing set, respectively. The best results are in bold.

Dataset	Class Type	SO	SO+A	MA	ROULETTE
NSL-KDD	Normal	1696.76	1775.18	2056.15	967.72
	DoS	2684.77	2512.99	2287.34	1263.12
	Probe	3157.55	2193.46	2075.19	1416.23
	U2R	2206.23	2211.46	3248.76	936.60
	R2L	2014.20	2433.27	2901.02	1038.04
	Average	2351.90	2225.28	2513.69	1124.34
UNSW-NB15	Normal	2645.61	2119.41	2660.66	2051.93
	Analysis	2361.67	2147.85	2306.63	2101.82
	Backdoors	1600.91	1442.60	1822.27	1343.21
	DoS	2500.13	2262.99	2421.57	2150.49
	Exploits	2259.04	2155.58	2496.73	1955.91
	Fuzzers	2492.79	1536.99	2140.86	1642.90
	Generic	2691.97	2492.47	2253.04	1997.52
	Reconnaissance	1608.62	1558.29	1807.41	1608.31
	Shellcode	1979.01	1704.12	2004.33	1571.44
	Worms	2320.29	1921.21	2571.25	2069.88
	Average	2246.00	1934.15	2248.48	1849.34

593 *6.3. Robustness*

594 As a further property, we explored the robustness of decision explanations
595 in relation to their ability to learn a classification model whose decisions are
596 still accurate on new attacks (e.g., variants of existing attacks). The robustness
597 of a classification model is commonly evaluated in terms of the accuracy of
598 classifications produced on unseen (test) data. In particular, our expectation
599 regarding the robustness of the decision explanations is that the explanations
600 of the decisions learned on the training network flow traces are roughly similar
601 to the explanations produced on unseen traces, which may be zero-day attacks.
602 We feel that this property of explanation robustness, if verified, helps predict
603 the correct class of unseen data.

604 Table 7 reports the results of the Average link distance computed for each
605 single class, as well as the mean of the Average link distance computed on all
606 the classes. The mean results highlight that the overall robustness of the ex-
607 planations produced with the trainable attention mechanism (ROULETTE and
608 SO+A) is better than the robustness of the counterpart explanations produced
609 with the post-hoc Grad-CAM technique (MO and SO). The lowest Average link
610 distance is measured again with ROULETTE, which also achieves the highest ac-
611 curacy performance in Section 5. In these results we can see empirical evidence
612 that the increased robustness of decision explanations may be responsible for
613 the increased accuracy of the decisions.

614 The same conclusions can also be drawn by analyzing the results of Aver-
615 age link distance computed per class type. Indeed, ROULETTE produces the
616 most robust decision explanations in almost any class. The only exceptions
617 are observed for the Fuzzer and Worm classes of UNSW-NB15. However, the
618 differences between the best Average link distance of SO+A and the runner-up
619 Average link distance of ROULETTE are small in both of these classes (1536.99
620 vs. 1642.90 in Fuzzers and 1921.21 vs. 2069.88 in Worms). In addition, this
621 difference has negligible impact on accuracy. ROULETTE slightly outperforms
622 SO+A on Fuzzers ($F1 = 0.38$ in ROULETTE vs. $F1 = 0.37$ in SO+A in Fig. 3),
623 while both ROULETTE and SO+A (as well as SO and MO) fail to recognize all

624 Worm testing traces ($F1 = 0.0$ in Fig. 3).

625 *6.4. Separability*

Table 8: Macro-F1, Weighted-F1 and A of cluster-based classifications produced on both the training set and testing set. Clusters are computed with k -means on the training set populated with the original images of network flow traffic traces, as well as on the training set populated with the images of the visual explanation maps determined with the attention mechanism.

		Image	k	Macro-F1	Weighted-F1	A
NSL-KDD	Train	Attention	18	0.55	0.94	0.95
		Original	17	0.52	0.92	0.93
	Test	Attention	18	0.41	0.64	0.68
		Original	17	0.35	0.60	0.66
UNSW-NB15	Train	Attention	29	0.35	0.70	0.73
		Original	27	0.30	0.67	0.71
	Test	Attention	29	0.31	0.69	0.68
		Original	27	0.25	0.64	0.63

626 Finally, we explored how the explanation information synthesized through
627 the attention mechanism can actually help achieve accuracy in separating the
628 network flow traces of different class types. To this end, we compared the effect
629 of information enclosed in:

- 630 • The set of visual explanation maps that are intrinsically produced by the
631 neural model with attention for the training network flow traces.
- 632 • The original images of the training network flow traces.

633 Since the previous analyses had already assessed the performance of the
634 multi-output Deep Learning strategy, we considered here the visual explanation
635 maps produced through the attention mechanism of ROULETTE. We first per-
636 formed a clustering to isolate distinct prototypes hidden in both the original

637 images and the visual explanations that populate the training set. We used
638 the Elbow method with Inertia to determine the optimal value of k for k -means
639 (Ketchen & Shook 1996). Then, we assigned each cluster centroid to the ma-
640 jority class of the training network flow traces grouped in the cluster. Finally,
641 we measured the accuracy of cluster centroid-based decisions, in order to eval-
642 uate how the information contained in both the original images and the visual
643 explanation maps is actually useful for properly separating network flow traces
644 belonging to multiple classes. This evaluation was done on both the training
645 set and the testing set by classifying each network flow trace in the class type
646 associated with the nearest centroid.

647 The results of Macro-F1, Weighted-F1 and A computed on the decisions based
648 on the cluster centroids are collected in Table 8. These results show empirical
649 evidence that attention can actually gain accuracy as it is able to capture in-
650 formation that separate network flow traces better than the original data.

651 *6.5. Insights into the Classification Explanation*

652 We completed this study by exploring how the attention mechanism of
653 ROULETTE can help reveal important relationships between the flow charac-
654 teristics of the processed network traffic and the observed categories of the
655 observed intrusions. To this end, we analyzed the information contained in the
656 visual explanation maps produced for both the NSL-KDD and UNSW-NB15
657 datasets.

658 The heatmaps in Figs. 5 and 6 depict the ranked average feature relevance of
659 the flow characteristics that “drew” the attention of the neural models learned
660 from both NSL-KDD and UNSW-NB15, respectively. For each dataset, two
661 heatmaps report the top 15 flow characteristics ranked by class on both the
662 training set and the testing set, respectively. Notably, the set of the top-ranked
663 flow characteristics that the neural model mostly attended on in the training
664 set largely overlaps the set of the top-ranked flow characteristics attended on
665 in the testing set. In particular, 21 flow characteristics are enclosed in the
666 top 15 characteristics highlighted both in the training set and the testing set

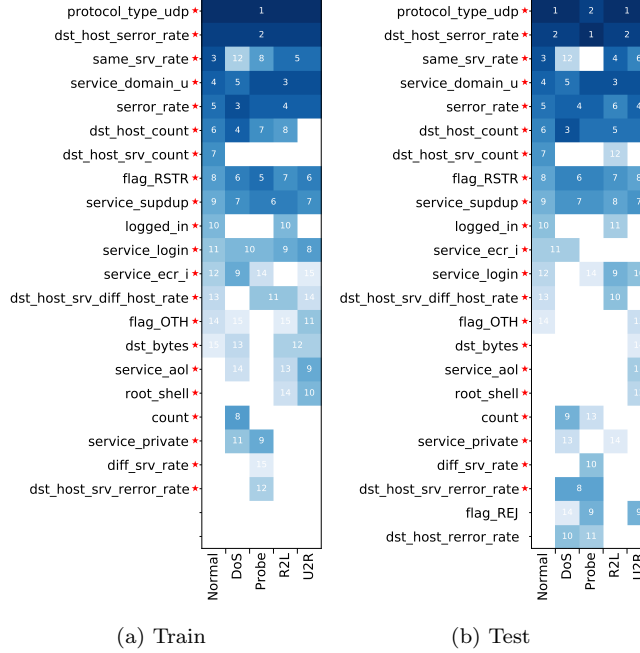


Figure 5: NSL-KDD: feature ranking map of the classification model learned with ROULETTE. We plot the ranking (1–15) of the flow characteristics (Y axis), which are ranked in the top 15 positions of the feature ranking determined with the attention mechanism over the various network flow traces grouped by class type (X axis) of both the training set (Fig. 5a) and the testing set (Fig. 5b). A star denotes the features that appear in the top 15 ranking of both the training set and testing set.

667 of NSL-KDD, while 20 flow characteristics are enclosed in the top 15 charac-
 668 teristics highlighted both in the training set and testing set of UNSW-NB15.
 669 Further considerations can be made by analyzing the meaning of the specific
 670 flow characteristics mostly attended on in both datasets.

671 6.5.1. NSL-KDD

672 Figure 5 shows that both *protocol_type_udp* and *dst_host_serror_rate* are rel-
 673 evant for classifying every class type in NSL-KDD, while the remaining flow
 674 characteristics are considered more or less relevant, depending on the specific
 675 type of intrusion to be classified.

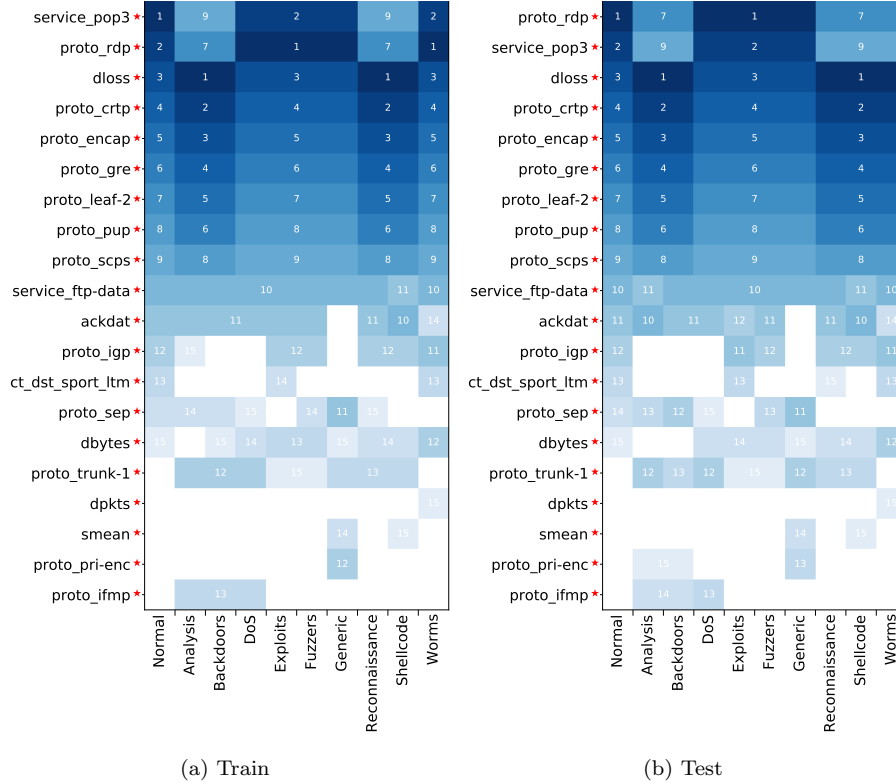


Figure 6: UNSW-NB15: feature ranking map of the classification model learned with ROULETTE. We plot the rank (1–15) of the flow characteristics (Y axis), which are ranked in the top 15 positions of the feature ranking determined with the attention mechanism over the various network flow traces grouped by class type (X axis) of both the training set (Fig. 6a) and the testing set (Fig. 6b). A star denotes the features that appear in the top 15 ranking of both the training set and testing set.

676 For example, *count* (i.e., “the number of connections to the same host as
677 the current connection in the past two seconds”) is seen as one of the most rele-
678 vant flow characteristics that the neural model takes into account to detect DoS
679 attacks, while it is less important when detecting other types of attacks. This
680 prominent role of *count* for DoS intrusion detection is consistent with the target
681 of a DoS attack which is to make a computer or network resource unavailable
682 (temporarily or indefinitely) to users by flooding the targeted machine with var-

683 ious connection requests. Similar considerations can be made on *dst_host_count*.
684 This flow characteristic, which counts “the number of connections having the
685 same destination host”, also conveys relevant information to detect DoS intru-
686 sions. In fact, a DoS attack can be in progress when a server is flooded by
687 sending numerous service requests to the target host.

688 Further considerations concern the attention paid to *service_ecr_i* to detect
689 DoS intrusions. In (Wang et al. 2020), *service_ecr_i* is recognized as a relevant
690 flow characteristic for the detection of Smurf attacks (a subcategory of DoS
691 intrusions) since, in this type of DoS, the targets are flooded with ECHO RE-
692 PLAY packets from each host on the broadcast address. Our study reveals that
693 the neural model attends on *service_ecr_i* to detect DoS intrusions. In addition,
694 *root_shell* (which equals 1 if the root shell is obtained; 0 otherwise) contains
695 relevant information for detecting U2R intrusions. U2R is a type of attack in
696 which the attacker tries to access network resources as a normal user, in order to
697 gain full access to the system. A U2R strategy might attempt to gain access to a
698 shell with administrator privilege (root shell). Finally, *dst_host_srv_error_rate*
699 and *diff_srv_rate* were considered relevant by the attention mechanism to de-
700 tect Probe attacks. The relationship between these two flow characteristics and
701 Probe has recently been discussed in (Wang et al. 2020).

702 6.5.2. UNSW-NB15

703 Figure 6 shows that various *proto* and *service*-based flow characteristics “at-
704 tracted” the attention of the neural model in UNSW-NB15. These flow char-
705 acteristics correspond to the transaction protocol (e.g., RDP, CRTP) used in
706 the network flow trace and the type of connection service (e.g., FTP, HTTP),
707 respectively. While these flow characteristics appear relevant for the detection
708 of several class types, their relevance changes with the class type. For exam-
709 ple, *service_pop3*, which is the most relevant flow characteristic for detecting
710 network flow traces in the categories Normal, DoS, Exploits, Fuzzers, Generic
711 and Worms, is slightly relevant for detecting traces in the categories Analysis,
712 Backdoors, Reconnaissance and Shellcode.

713 In addition, *ackdat*, which refers to the TCP connection setup time (i.e., “the
714 time between the SYN ACK and the ACK response”) is relevant for detecting
715 shellcode intrusions, while it becomes less important when detecting other types
716 of attacks. Shellcode, in fact, is an exploiting attack in which the attacker
717 penetrates a piece of code from a shell to control a target machine using the
718 standard TCP/IP socket connections.

719 Finally, *dpkts*, i.e., “the count of the number of packets from source to desti-
720 nation”, is relevant for recognizing worm attacks. We recall that worm attacks
721 are self-replicating computer programs that spread automatically and can flood
722 the Internet in a very short time (Chen et al. 2003).

723 7. Conclusions

724 In this paper, we have presented ROULETTE: a system for multi-class clas-
725 sification of network traffic data. The proposed method learns a neural classifi-
726 cation model through a multi-output Deep Learning strategy that encompasses
727 both convolution and attention. Extensive experimentation was performed to
728 show the effectiveness of the proposed neural model with attention, quantified
729 in terms of accuracy of classifications, as well as transparency of decisions. In
730 particular, the results obtained indicate that ROULETTE is able to produce
731 decisions that are comparable to (or even more accurate than) decisions pro-
732 duced with competitive, Deep Learning-based approaches. Furthermore, the
733 results of the experimentation highlighted how the good accuracy performance
734 of ROULETTE can also be attributed to specific properties, such as the compact-
735 ness, robustness, and separability of the produced attention-based explanations.
736 Finally, the attention mechanism helps us to see particular characteristics of
737 network traffic that mainly help to recognize specific intrusion categories. This
738 may support the dissemination of useful information to cyber-defenders, thus
739 reducing the workload in manual analysis.

740 One limitation of the proposed method is the absence of a specific mech-
741 anism for dealing with rare classes. A research direction is to explore data

742 augmentation techniques, in order to reach a balancing condition in the learn-
743 ing stage. For example, GANs have recently helped to increase accuracy in the
744 binary classification of images of network traffic (Andresini, Appice, De Rose &
745 Malerba 2021).

746 Another limitation is that the proposed method performs the learning stage
747 in a batch fashion, without integrating any concept drift detection mechanism
748 to properly fit the learned model to an evolving streaming environment. This
749 is an issue for all adversary-facing security systems. Recent studies have begun
750 to investigate this problem both in NID applications (Andresini, Pendlebury,
751 Pierazzi, Loglisci, Appice & Cavallaro 2021) and malware detection problems
752 (Pendlebury et al. 2019). These studies propose incremental, semi-supervised
753 security systems to process cyber-data streams by reducing labeling overhead
754 and continuously updating the underlying model as the data characteristics are
755 affected by concept drift.

756 Finally, recent studies in Computer Vision have achieved amazing results
757 with transformer-based architectures, such as the increasingly popular ViT
758 (Dosovitskiy et al. 2020). The idea behind transformers is to define the lay-
759 ers of the neural network entirely on the attention mechanism. At each layer
760 a new hidden representation is generated for each position in the input data
761 by using multiple attention heads that calculate attention weights for all pairs
762 of positions in the input. Although the image encoding adopted in this study
763 is already robust to arbitrary permutations of features, it is still constrained
764 by the inductive “locality” bias of standard convolutions. The multi-headed
765 self-attention strategy of vision transformers could further improve the general-
766 ization of the model.

767 **CRedit Authorship Contribution Statement**

768 **Giuseppina Andresini:** Conceptualization, Methodology, Software, Data
769 curation, Investigation, Validation, Visualization, Writing - original draft, Writ-
770 ing - review & editing. **Annalisa Appice:** Conceptualization, Methodology,

771 Investigation, Validation, Supervision, Writing - original draft, Writing - review
772 & editing. **Francesco Paolo Caforio**: Conceptualization, Software, Investiga-
773 tion. **Donato Malerba**: Conceptualization, Project administration, Writing -
774 review & editing. **Gennaro Vessio**: Conceptualization, Methodology, Valida-
775 tion, Writing - original draft, Writing - review & editing, Supervision.

776 **Acknowledgments**

777 We acknowledge the support of the Italian Ministry of University and Re-
778 search through the project “TALIsMan - Tecnologie di Assistenza personALiz-
779 zata per il Miglioramento della qualità della vita” (Grant ID: ARS01.01116),
780 funding scheme PON RI 2014-2020, as well as the project “Modelli e tecniche
781 di data science per l’analisi di dati strutturati” funded by the University of
782 Bari “Aldo Moro”. The authors wish to thank Jerzy Stefanowski for his helpful
783 comments which inspired the evaluation of the separability property and Lynn
784 Rudd for her help in reading the manuscript.

785 **References**

- 786 Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado,
787 G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., ... & Zheng, X. (2015),
788 ‘TensorFlow: Large-scale machine learning on heterogeneous systems’.
789 **URL:** <https://www.tensorflow.org/>
- 790 Al-Turaiki, I. & Altwaijry, N. (2021), ‘A convolutional neural network for im-
791 proved anomaly-based network intrusion detection’, *Big Data* **9**(3), 233–252.
- 792 Alex Kendall, V. B. & Cipolla, R. (2017), Bayesian segnet: Model uncertainty
793 in deep convolutional encoder-decoder architectures for scene understanding,
794 *in* ‘Proceedings of the British Machine Vision Conference (BMVC)’, BMVA
795 Press, pp. 57.1–57.12.

- 796 Andresini, G., Appice, A., De Rose, L. & Malerba, D. (2021), ‘Gan augmen-
797 tation to deal with imbalance in imaging-based intrusion detection’, *Future*
798 *Generation Computer Systems* **123**, 108–127.
- 799 Andresini, G., Appice, A. & Malerba, D. (2021a), ‘Autoencoder-based deep
800 metric learning for network intrusion detection’, *Inf. Sci.* **569**, 706–727.
- 801 Andresini, G., Appice, A. & Malerba, D. (2021b), ‘Nearest cluster-based in-
802 trusion detection through convolutional neural networks’, *Knowledge-Based*
803 *Systems* **216**, 106798.
- 804 Andresini, G., Appice, A., Mauro, N. D., Loglisci, C. & Malerba, D. (2020),
805 ‘Multi-channel deep feature learning for intrusion detection’, *IEEE Access*
806 **8**, 53346–53359.
- 807 Andresini, G., Pendlebury, F., Pierazzi, F., Loglisci, C., Appice, A. & Cavallaro,
808 L. (2021), Insomnia: Towards concept-drift robustness in network intrusion
809 detection, in ‘Proceedings of the 14th ACM Workshop on Artificial Intelli-
810 gence and Security’, AISec’21, Association for Computing Machinery, New
811 York, USA, p. 111–122.
- 812 Antwarg, L., Miller, R. M., Shapira, B. & Rokach, L. (2021), ‘Explaining anoma-
813 lies detected by autoencoders using shapley additive explanations’, *Expert*
814 *Systems with Applications* **186**, 115736.
- 815 Bahdanau, D., Cho, K. & Bengio, Y. (2015), Neural machine translation by
816 jointly learning to align and translate, in Y. Bengio & Y. LeCun, eds, ‘3rd In-
817 ternational Conference on Learning Representations, ICLR 2015, Conference
818 Track Proceedings’.
- 819 Bedi, P., Gupta, N. & Jindal, V. (2020a), ‘I-SiamIDS: an improved siam-ids
820 for handling class imbalance in network-based intrusion detection systems’,
821 *Applied Intelligence* **abs/2009.10940**.

- 822 Bedi, P., Gupta, N. & Jindal, V. (2020b), ‘Siam-IDS: Handling class imbalance
823 problem in intrusion detection systems using siamese neural network’,
824 *Procedia Computer Science* **171**, 780 – 789.
- 825 Bergstra, J., Bardenet, R., Bengio, Y. & Kégl, B. (2011), Algorithms for hyper-
826 parameter optimization, *in* ‘Advances in Neural Information Processing Sys-
827 tems (NeurIPS)’.
- 828 Bergstra, J., Yamins, D. & Cox, D. D. (2013), Making a science of model search:
829 Hyperparameter optimization in hundreds of dimensions for vision architec-
830 tures, *in* ‘Proc. of the International Conference on Machine Learning (ICML)’.
- 831 Berman, D. S., Buczak, A. L., Chavis, J. S. & Corbett, C. L. (2019), ‘A survey
832 of deep learning methods for cyber security’, *Information* **10**(4), 1–35.
- 833 Biecek, P. (2018), ‘Dalex: Explainers for complex predictive models in r’, *Jour-
834 nal of Machine Learning Research* **19**(84), 1–5.
- 835 Burkart, N., Franz, M. & Huber, M. F. (2021), Explanation framework for
836 intrusion detection, *in* J. Beyerer, A. Maier & O. Niggemann, eds, ‘Machine
837 Learning for Cyber Physical Systems’, Springer Berlin Heidelberg, Berlin,
838 Heidelberg, pp. 83–91.
- 839 Caforio, F., Andresini, G., Vessio, G., Appice, A. & Malerba, D. (2021), Leverag-
840 ing grad-cam to improve the accuracy of network intrusion detection systems,
841 *in* ‘Discovery Science’, Springer International Publishing, pp. 385–400.
- 842 Caminero, G., Lopez-Martin, M. & Carro, B. (2019), ‘Adversarial environment
843 reinforcement learning algorithm for intrusion detection’, *Computer Networks*
844 **159**, 96 – 109.
- 845 Cao, L., Li, L., Zheng, J., Fan, X., Yin, F., Shen, H. & Zhang, J. (2018), ‘Multi-
846 task neural networks for joint hippocampus segmentation and clinical score
847 regression’, *Multimedia Tools and Applications* **77**(22), 29669–29686.

- 848 Castellano, G., Castiello, C., Mencar, C. & Vessio, G. (2020), ‘Crowd detection
849 in aerial images using spatial graphs and fully-convolutional neural networks’,
850 *IEEE Access* **8**, 64534–64544.
- 851 Chen, Z., Gao, L. & Kwiat, K. (2003), Modeling the spread of active worms,
852 in ‘IEEE INFOCOM 2003. Twenty-second Annual Joint Conference of the
853 IEEE Computer and Communications Societies (IEEE Cat. No.03CH37428)’,
854 Vol. 3, pp. 1890–1900 vol.3.
- 855 Choudhary, S. & Kesswani, N. (2020), ‘Analysis of kdd-cup’99, nsl-kdd and
856 unsw-nb15 datasets using deep learning in iot’, *Procedia Computer Science*
857 **167**, 1561–1573. International Conference on Computational Intelligence and
858 Data Science.
- 859 Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Un-
860 terthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S. et al. (2020),
861 ‘An image is worth 16x16 words: Transformers for image recognition at scale’,
862 *arXiv preprint arXiv:2010.11929* .
- 863 Folino, F., Folino, G., Guarascio, M., Pisani, F. & Pontieri, L. (2021), ‘On
864 learning effective ensembles of deep neural networks for intrusion detection’,
865 *Information Fusion* **72**, 48–69.
- 866 Gao, J., Chai, S., Zhang, B. & Xia, Y. (2019), ‘Research on network intru-
867 sion detection based on incremental extreme learning machine and adaptive
868 principal component analysis’, *Energies* **12**(7).
- 869 Gao, M., Ma, L., Liu, H., Zhang, Z., Ning, Z. & Xu, J. (2020), ‘Malicious net-
870 work traffic detection based on deep neural networks and association analysis’,
871 *Sensors* **20**(5), 1452.
- 872 Glorot, X. & Bengio, Y. (2010), ‘Understanding the difficulty of training deep
873 feedforward neural networks’, *Journal of Machine Learning Research - Pro-
874 ceedings Track* **9**, 249–256.

- 875 Glorot, X., Bordes, A. & Bengio, Y. (2011), Deep sparse rectifier neural net-
876 works., *in* ‘AISTATS’, JMLR.org, pp. 315–323.
- 877 Guo, M.-H., Xu, T.-X., Liu, J.-J., Liu, Z.-N., Jiang, P.-T., Mu, T.-J., Zhang, S.-
878 H., Martin, R. R., Cheng, M.-M. & Hu, S.-M. (2021), ‘Attention mechanisms
879 in computer vision: A survey’, *arXiv preprint arXiv:2111.07624* .
- 880 Guo, Y., Cao, H., Bai, J. & Bai, Y. (2019), ‘High efficient deep feature extraction
881 and classification of spectral-spatial hyperspectral image using cross domain
882 convolutional neural networks’, *IEEE Journal of Selected Topics in Applied
883 Earth Observations and Remote Sensing* **12**(1), 345–356.
- 884 Joshi, G., Walambe, R. & Kotecha, K. (2021), ‘A review on explainability in
885 multimodal deep neural nets’, *IEEE Access* **PP**, 1–1.
- 886 Kasongo, S. & Sun, Y. (2020), ‘Performance analysis of intrusion detection
887 systems using a feature selection method on the unsw-nb15 dataset’, *J Big
888 Data* **7**(105), 1–20.
- 889 Ketchen, D. J. & Shook, C. L. (1996), ‘The application of cluster analysis in
890 strategic management research: An analysis and critique’, *Strategic Manage-
891 ment Journal* **17**(6), 441–458.
- 892 Kingma, D. P. & Ba, J. (2014), Adam: A method for stochastic optimization,
893 *in* ‘ICLR’.
- 894 Komodakis, N. & Zagoruyko, S. (2017), Paying more attention to attention:
895 improving the performance of convolutional neural networks via attention
896 transfer, *in* ‘ICLR’.
- 897 Lakkaraju, H., Kamar, E., Caruana, R. & Leskovec, J. (2019), Faithful and cus-
898 tomizable explanations of black box models, *in* ‘Proc. of the 2019 AAAI/ACM
899 Conference on AI, Ethics, and Society’, pp. 131–138.
- 900 Liu, C., Liu, Y., Yan, Y. & Wang, J. (2020), ‘An intrusion detection model with
901 hierarchical attention mechanism’, *IEEE Access* **8**, 67542–67554.

- 902 Lopez-Martin, M., Carro, B., Sanchez-Esguevillas, A. & Lloret, J. (2017), ‘Con-
903 ditional variational autoencoder for prediction and feature recovery applied
904 to intrusion detection in IoT’, *Sensors* **17**(9), 1–17.
- 905 Lundberg, S. M. & Lee, S.-I. (2017), A unified approach to interpreting model
906 predictions, *in* ‘Proceedings of the 31st International Conference on Neural
907 Information Processing Systems’, NIPS’17, Curran Associates Inc., pp. 4768–
908 4777.
- 909 Mane, S. & Rao, D. (2021), ‘Explaining network intrusion detection system
910 using explainable AI framework’, *arXiv preprint arXiv:2103.07110* .
- 911 Marino, D. L., Wickramasinghe, C. S. & Manic, M. (2018), An adversarial
912 approach for explainable ai in intrusion detection systems, *in* ‘IECON 2018 -
913 44th Annual Conference of the IEEE Industrial Electronics Society’, pp. 3237–
914 3243.
- 915 Moustafa, N. & Slay, J. (2015), Unsw-nb15: a comprehensive data set for net-
916 work intrusion detection systems (unsw-nb15 network data set), *in* ‘2015 Mil-
917 itary Communications and Information Systems Conference (MilCIS)’, pp. 1–
918 6.
- 919 Naseer, S., Saleem, Y., Khalid, S., Bashir, M. K., Han, J., Iqbal, M. M. &
920 Han, K. (2018), ‘Enhanced network anomaly detection based on deep neural
921 networks’, *IEEE Access* **6**, 48231–48246.
- 922 Pendlebury, F., Pierazzi, F., Jordaney, R., Kinder, J. & Cavallaro, L. (2019),
923 TESSERACT: Eliminating Experimental Bias in Malware Classification
924 across Space and Time, *in* ‘28th USENIX Security Symposium’, USENIX
925 Association. USENIX Sec.
- 926 Phaisangittisagul, E. (2016), An analysis of the regularization between l2 and
927 dropout in single hidden layer neural network, *in* ‘2016 7th International Con-
928 ference on Intelligent Systems, Modelling and Simulation (ISMS)’, pp. 174–
929 179.

- 930 Sarhan, M., Layeghy, S. & Portmann, M. (2021), ‘An explainable machine
931 learning-based network intrusion detection system for enabling generalisabil-
932 ity in securing iot networks’, *CoRR* **abs/2104.07183**.
- 933 Sartor, G. & Lagioia, F. (2020), *The impact of the General Data Protection*
934 *Regulation (GDPR) on artificial intelligence*.
- 935 Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R., Parikh, D. & Batra,
936 D. (2017), Grad-CAM: Visual explanations from deep networks via gradient-
937 based localization, *in* ‘Proceedings of the IEEE international conference on
938 computer vision’, pp. 618–626.
- 939 Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R., Parikh, D. & Batra,
940 D. (2020), ‘Grad-cam: Visual explanations from deep networks via gradient-
941 based localization’, *Int. J. Comput. Vis.* **128**(2), 336–359.
- 942 Sovilj, D., Budnarain, P., Sanner, S., Salmon, G. & Rao, M. (2020), ‘A compar-
943 ative evaluation of unsupervised deep architectures for intrusion detection in
944 sequential data streams’, *Expert Systems with Applications* **159**, 113577.
- 945 Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I. & Salakhutdinov, R.
946 (2014), ‘Dropout: A simple way to prevent neural networks from overfitting’,
947 *Journal of Machine Learning Research* **15**(56), 1929–1958.
- 948 Szczepański, M., Choraś, M., Pawlicki, M. & Kozik, R. (2020), Achieving ex-
949 plainability of intrusion detection system by hybrid oracle-explainer approach,
950 *in* ‘2020 International Joint Conference on Neural Networks (IJCNN)’, pp. 1–
951 8.
- 952 Tang, C., Luktarhan, N. & Zhao, Y. (2020), ‘SAAE-DNN: Deep learning method
953 on intrusion detection’, *Symmetry* **12**(10), 1–20.
- 954 Tavallae, M., Bagheri, E., Lu, W. & Ghorbani, A. A. (2009), A detailed analysis
955 of the KDD CUP 99 data set, *in* ‘CISDA’, IEEE, pp. 1–6.

- 956 Vinayakumar, R., Alazab, M., Soman, K. P., Poornachandran, P., Al-Nemrat,
957 A. & Venkatraman, S. (2019), ‘Deep learning approach for intelligent intrusion
958 detection system’, *IEEE Access* **7**, 41525–41550.
- 959 Wali, S. & Khan, I. (2021), ‘Explainable AI and random forest based reliable
960 intrusion detection system’.
- 961 Wang, M., Zheng, K., Yang, Y. & Wang, X. (2020), ‘An explainable machine
962 learning framework for intrusion detection systems’, *IEEE Access* **8**, 73127–
963 73141.
- 964 Warnecke, A., Arp, D., Wressnegger, C. & Rieck, K. (2020), Evaluating ex-
965 planation methods for deep learning in security, *in* ‘2020 IEEE European
966 Symposium on Security and Privacy (EuroS&P)’, IEEE, pp. 158–174.
- 967 Xu, D., Shi, Y., Tsang, I. W., Ong, Y.-S., Gong, C. & Shen, X. (2019), ‘Sur-
968 vey on multi-output learning’, *IEEE Transactions on Neural Networks and*
969 *Learning Systems* **31**(7), 2409–2429.
- 970 Xu, F., Uszkoreit, H., Du, Y., Fan, W., Zhao, D. & Zhu, J. (2019), Explainable
971 ai: A brief survey on history, research areas, approaches and challenges, *in*
972 J. Tang, M.-Y. Kan, D. Zhao, S. Li & H. Zan, eds, ‘Natural Language Pro-
973 cessing and Chinese Computing’, Springer International Publishing, Cham,
974 pp. 563–574.
- 975 Yang, J., Li, T., Liang, G., He, W. & Zhao, Y. (2019), ‘A simple recurrent
976 unit model based intrusion detection system with DCGAN’, *IEEE Access*
977 **7**, 83286–83296.
- 978 Yin, C., Zhu, Y., Fei, J. & He, X. (2017), ‘A deep learning approach for intrusion
979 detection using recurrent neural networks’, *IEEE Access* **5**, 21954–21961.
- 980 Zhao, H., Kong, X., He, J., Qiao, Y. & Dong, C. (2020), Efficient image super-
981 resolution using pixel attention, *in* ‘European Conference on Computer Vi-
982 sion’, Springer, pp. 56–72.

983 Zhao, P., Fan, Z., Cao, Z. & Li, X. (2022), ‘Intrusion detection model using tem-
984 poral convolutional network blend into attention mechanism’, *International*
985 *Journal of Information Security and Privacy* **16**(1), 1–20.