



European Research Council
Established by the European Commission

Self-assessment Oracles for Anticipatory Testing

TECHNICAL REPORT: TR-Precrime-2023-09

Tahereh Zohdinasab^a, Vincenzo Riccio^a, Paolo Tonella^a

^aUniversità della Svizzera italiana

An Empirical Study on Low- and High-Level Explanations of Deep Learning Misbehaviours

| | |
|-----------------------------------|---|
| Project no.: | 787703 |
| Funding scheme: | ERC-2017-ADG |
| Start date of the project: | January 1, 2019 |
| Duration: | 60 months |
| Technical report num.: | TR-Precrime-2023-09 |
| Date: | July, 2023 |
| Organization: | Università della Svizzera italiana |
| Authors: | Tahereh Zohdinasab ^a , Vincenzo Riccio ^a , Paolo Tonella ^a |
| Dissemination level: | ^a Università della Svizzera italiana Public |
| Revision: | 1.0 |

Disclaimer:

This Technical Report is a pre-print of the following publication:

Tahereh Zohdinasab, Vincenzo Riccio, Paolo Tonella: *An Empirical Study on Low- and High-Level Explanations of Deep Learning Misbehaviours*. Proceedings of the ACM/IEEE International Symposium on Empirical Software Engineering and Measurement (ESEM), New Orleans, Louisiana, USA, October 22–27, 2023

Please, refer to the published version when citing this work.





Università della Svizzera Italiana (USI)

Principal investigator: Prof. Paolo Tonella
E-mail: paolo.tonella@usi.ch
Address: Via Buffi, 13 – 6900 Lugano – Switzerland
Tel: +41 58 666 4848
Project website: <https://www.pre-crime.eu/>

Abstract

Background: Most quality assessment approaches for Deep Learning (DL) focus on finding misbehaviour-inducing inputs. However, it is difficult to clearly understand the causes of misbehaviours, due to the DL software opaqueness. Recent research proposed different techniques to explain DL misbehaviours, producing input explanations either at a “low level” (raw input elements) or at a “high level” (input features).

Aims: We aim to compare the similarity between different explanations and assess to what extent they are understandable.

Method: We have conducted an empirical study involving 3 state-of-the-art techniques for DL explanation in 13 configurations, applied to 2 different DL tasks. We have also collected answers from 48 questionnaires submitted to SE experts.

Results: Low- and high-level techniques provide dissimilar explanations for the same inputs. However, experts deemed none of the explanations as useful in 28% of the cases.

Conclusion: Despite the complementarity of existing explanations, further research is needed to produce better explanations.

Contents

| | | |
|----------|--|-----------|
| 1 | Introduction | 1 |
| 2 | Background | 2 |
| 2.1 | Feature Map (FM) | 2 |
| 2.2 | Integrated Gradients (IG) | 2 |
| 2.3 | Local Interpretable Model-agnostic Explanations (LIME) | 3 |
| 3 | Comparing High- and Low-level explanations | 3 |
| 4 | Experimental Evaluation | 6 |
| 4.1 | Research Questions | 7 |
| 4.2 | Subject Systems | 7 |
| 4.3 | Experimental Procedure | 8 |
| 4.3.1 | High- and Low-Level Explanations Comparison | 9 |
| 4.3.2 | Human Assessment of the Explanations | 9 |
| 5 | Results | 11 |
| 5.1 | RQ1 (Similarity) | 11 |
| 5.2 | RQ2 (Understandability) | 13 |
| 5.3 | Discussion | 14 |
| 5.4 | Threats to Validity | 15 |
| 6 | Related Work | 16 |
| 6.1 | Comparison of Explainable AI techniques | 16 |
| 6.2 | Explainable AI for DL Testing | 16 |
| 7 | Conclusions and Future Work | 17 |

1 Introduction

Deep Learning (DL) algorithms are pervasively adopted in modern software systems, as they proved their ability to solve complex tasks, including safety-critical ones. A large part of the literature focuses on evaluating the quality of DL systems by triggering misbehaviours, i.e., deviations from their expected behaviour [28, 43]. Quality improvement requires understanding the causes of such misbehaviours, which is still an open problem [12, 38]. In fact, the increasing complexity of DL models, trained on large datasets thanks to the availability of scalable and high-performance infrastructures, introduces the risk to make (mis-)predictions whose motivations are hard to understand for humans. Recent DL algorithms, such as Deep Neural Networks (DNNs), involve a massive amount of computations and, thus, are “black boxes” to SE practitioners, since it is hard for them to understand how DL software arrives at a final decision [9, 38].

The lack of explainability of DL models may impact the trust in their predictions and consequently hinder the adoption of such systems. A common issue of DL systems is that they may take wrong decisions based on spurious correlations or biases in the training data, e.g., classifying wolves images as huskies due to the presence of snow on the background rather than leveraging the characteristics of the animal [26].

Multiple techniques address the aforementioned issues and explain DL systems’ (mis-)behaviours. Existing explanatory techniques characterise the misbehaviour-inducing inputs either at a low level, by identifying the raw input elements (e.g., image pixels) that are most relevant for the prediction, or at a high level, in terms of abstract features that characterise the input as a whole. We will refer to these techniques as follows:

Low Level Explanations identify a portion of the input as relevant to a specific DL prediction. Such identification is performed directly in the input space of the DL model. Low-level techniques can provide explanations with different aggregations, reporting, e.g., atomic input elements (e.g., pixels/words for an image/text) or contiguous regions of input elements (e.g., sets of contiguous pixels/words). We considered two techniques with different aggregations, representing the state of the art in the eXplainable AI (XAI) community: LIME [26] and INTEGRATED GRADIENTS (IG) [37].

High Level Explanations identify relevant features that characterise the whole input. Such identification is performed in the feature space, a manually-defined abstraction of the input space. As high-level technique, we considered FEATURE MAPS (FM) [44] that characterise inputs through human-interpretable features defined by domain experts. As an example, for the hand-written digit recognition task, feature maps characterise a misbehaviour by indicating how much the hand stroke is bold or discontinuous. FM proved to be highly discriminative of misbehaviours [44], so they indeed provide an explanation for misbehaviours.

In this work, we conduct an empirical study to investigate the similarity between explanations provided by low- and high-level techniques for the same sets of inputs. In this way, we can understand whether and to what extent different techniques are overlapping or complementary in explaining DL misbehaviours. Explanatory techniques provide an input characterisation that practitioners (e.g., software testers) can use to understand why an input triggered a misbehaviour [22]. For this reason, we involved human experts in our study to assess the understandability of the produced explanations, i.e., whether human assessors consider the explanations as effectively pointing to a plausible cause of the misbehaviour.

By leveraging input clustering and a novel ad-hoc similarity metric, our results show that high- and low-level techniques provide different explanations, i.e., they partition the same inputs in different ways. For most of the inputs, human experts chose either a low- or a high-level technique as effective in explaining misbehaviours. These results show that high- and low-level explanations are highly complementary. Results also show that both explanations often do not match human judgement, especially when explaining misbehaviours for image classifiers, hence demanding for novel XAI techniques that can cover the mis-explained cases. In summary, our paper makes the

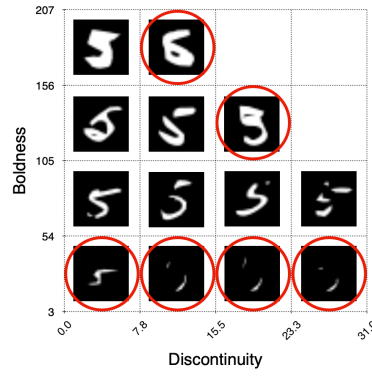


Figure 1: FM for a handwritten digit classifier. Axes quantify boldness and discontinuity of the digits. Cells report the most challenging input for the corresponding feature values, and the cells highlighted with red circles are mis-classified.

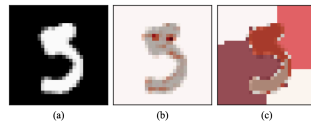


Figure 2: (a) A sample mis-classified handwritten digit; (b) Explanation by IG; (c) Explanation by LIME. (darker red highlights higher contribution to the DNN output)

following contributions:

- We introduce an empirical framework for comparing explanatory techniques with different granularity levels;
- We evaluate explanatory techniques on DNNs within two different domains, i.e., image and text classification;
- We perform a human study for comparing the understandability of explanations for DNN mis-behaviours;

2 Background

2.1 Feature Map (FM)

A FM represents the space of the features that are relevant for characterising the test inputs as a multi-dimensional grid, in which each axis corresponds to a considered feature. A feature can either be *structural* (i.e., characterising the input itself) or *behavioural* (i.e., characterising the system’s output when exercised by the given input). In particular, we refer to high-level, human-interpretable input features defined by experts [44], such as digit boldness and discontinuity for the handwritten digit classifier shown in Figure 1. Inputs are assigned a map cell, computed by measuring the metric that quantifies each feature. Therefore, each cell contains the inputs characterized by features that fall in the corresponding value interval. The map granularity (i.e., the number of map cells in each dimension) is a configurable parameter of a FM that can be changed depending on the desired level of discrimination.

2.2 Integrated Gradients (IG)

Associating the prediction of a DNN to the input elements that caused such prediction (i.e., *input attribution*) is a way of explaining the model’s behaviour. For instance, in an image classification

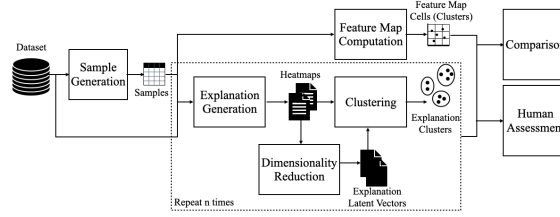


Figure 3: An overview of our evaluation pipeline

model, an attribution method could reveal which pixels of the image were responsible for a certain prediction.

The gradient of the output with respect to the input elements (e.g., pixels) quantifies the impact of each input element on the prediction of a DNN, therefore the gradient computation can be considered as the most basic attribution method.

Sundararajan et al. [37] took an axiomatic approach for attributing prediction of DNN to its input features, called IG. They introduced two axioms to be preserved, *Sensitivity* and *Implementation Invariance*, and proposed IG by considering the straight-line path in DNN from the baseline (for example, black image in image classification task) to the input and cumulating the gradients at all points along the path for a given number of steps. The *number of steps* is a hyperparameter of this technique, representing the steps needed by IG in the gradient approximation for each input image. The output of this technique is a heatmap highlighting the important elements in the input (e.g. pixels of an image or words in a text). IG can be implemented to generate explanation for a batch of inputs with predefined size. Figure 2 (b) shows a heatmap generated by IG for an image classification task.

2.3 Local Interpretable Model-agnostic Explanations (LIME)

LIME provides a low-level explanation for the predictions of any classification model by identifying the portions of the input that mostly affect the model’s output. It produces an explanation for each input instance by mimicking the behaviour of the model in the neighbourhood of the input. LIME perturbs the input to train a linear model around the input itself and generates a predefined number of samples, while the model under test is treated as a black box. The number of samples can be customized, as it can be set as a hyperparameter of the technique. Such perturbed input regions are called *super-pixels* and their impact on the linear model’s decision determines their importance in the final explanation. LIME is applicable to both images and texts and provides a visual representation of the generated explanations. It has been used in different studies to understand the reasons behind DNN decisions [8,17,21]. Figure 2 (c) shows an example of explanation generated by LIME for a hand written digit 5.

3 Comparing High- and Low-level explanations

To compare high-level and low-level explanations, we designed the evaluation pipeline shown in Figure 3. Initially, we gather the inputs to be explained. Such inputs can be either obtained directly from an existing dataset or can be automatically generated by test generators (*Sample Generation* step). For each input, we obtain high-level explanations in the *Feature Map Computation* step, i.e., we compute the FM cells it belongs to. In this way, each FM cell represents a cluster of inputs sharing similar high-level features. To obtain low-level explanations, we perform the *Explanation Generation* step, in which we apply low-level techniques that generate heatmaps. Moreover, we perform a *Dimensionality Reduction* step on low-level explanations to generate lower-dimensional latent vectors. A *latent vector* is the projection of an input onto the latent space. Using latent vectors allows us to avoid high variability and sparseness in the high-dimensional heatmap explanations.

Indeed, a sparse, highly variable representation of the inputs would prevent the construction of meaningful clusters of the inputs, as in the higher-dimensional space all distances/similarities tend to degenerate to the extreme values, giving rise to degenerate clusters. Since both low-level explanatory techniques and dimensionality reduction approaches are non-deterministic, we repeat the low-level explanation generation multiple times. To make the explanations provided by low-level techniques comparable to the high-level ones, we cluster low-level explanations by using a *Clustering* technique. In the *Clustering* step, we group together similar heatmap explanations, or the corresponding latent vectors (to address the sparseness problem). In the *Comparison* step, we assess how similar the clusters generated by high-level and low-level techniques are, based on our custom definition of a *Gini similarity* metric, which is close to zero when the elements grouped together by one technique are scattered across many clusters produced by the other technique and is equal to one when the elements grouped together by one technique are all in the same cluster produced by the other technique. Finally, we submit a questionnaire to experts in order to evaluate the understandability of explanations in the *Human Assessment* step. In the following, we provide a detailed explanation of each step of our evaluation pipeline.

Sample Generation: Since we are interested in investigating DL misbehaviours, we resort to failure-inducing inputs artificially crafted by test generators, besides the ones from the original dataset. In this way, we can collect enough misbehaviour-inducing inputs even for robust DL models, i.e., for which there are only a few misbehaviours in the original test set. By considering both types of inputs, we perform our study on a sufficient number of misbehaviours, covering a diversified variety of samples. We do not use the train set because we mimic the testing phase, when developers evaluate a DL model against new, unseen inputs

FM Computation: To provide high-level explanations, we used FM. FEATURE MAPS can be generated with different feature combinations and numbers of dimensions. For the sake of completeness, we considered all possible numbers of dimensions and feature combinations when computing the FM. In this way, we can discuss the similarity between high- and low-level explanations at increasing FM dimensionality (e.g., when moving from one isolated feature to a combination of multiple features). For each input obtained in the Sample Generation step, we compute its corresponding feature values, so we can assign each input to the corresponding FM cell, i.e., the map cell whose value intervals contain the measured input feature values. We use the non-empty FM cells, i.e., those containing at least one misbehaviour-inducing input, as high-level explanation clusters.

Explanation Generation: In this step, we consider two popular XAI techniques as our low-level explanatory techniques: IG for fine-grained explanations and LIME for coarse-grained explanations. We apply the two considered XAI techniques separately on the generated inputs and extract heatmaps as vectors representing the relevance of input elements/regions.

Dimensionality Reduction: Beside the heatmap explanations in the original input space, we consider explanations projected onto a latent space, i.e., latent vectors. In fact, by projecting heatmaps to latent space, we keep a lower number of dimensions, which are more representative of the meaningful directions of variability of the inputs, possibly avoiding the construction of degenerate clusters. For generating latent vectors, we choose the *t-SNE* algorithm [10,41], since it projects similar inputs to neighbouring points and dissimilar inputs to distant points with high probability.

We generate latent vectors in two modes: (1) *global latent*: the projections of explanations in the latent space are computed considering inputs from all output classes; (2) *local latent*: the projections of explanations in the latent space are computed considering only inputs from a specific class. While the local latent space can achieve good separation of inputs from a specific class, the global latent space aims at separating both different classes of inputs and inputs within each class. The former is better if the generated inputs stay confined within one class, while the latter might be beneficial when the generated inputs tend to cross the borders between classes.

Clustering: The explanations generated by high-level techniques are FM cells based on high-level feature values of the inputs, whereas low-level explanations are vectors of contributions of

low-level input elements/regions. Therefore, there is no way to directly measure the similarity between explanations at these two different abstraction levels. To make the comparison feasible, we propose to compare the *clusters of explanations* instead of the raw explanations (which are uncomparable). The underlying idea is the following: if two explanatory techniques group the inputs in a similar way, then they can be deemed similar; otherwise they are different.

For high-level explanations we use map cells as clusters. For low-level explanations, we need to group the explanations by using clustering approaches in such a way that objects in the same group, i.e., a cluster, are more similar to each other than to those in other groups. There exist multiple clustering techniques in the literature [4, 7, 14, 20]. For our study, we rely on the *Affinity Propagation* clustering technique [7] which recursively exchanges messages between data points; such messages encode the affinity of one point when choosing another point as its neighbor. The recursive exchange of messages continues until a set of highly-affine groups emerges. The main advantage of this clustering technique is that, unlike other techniques such as K-Means [20], (1) all the points are considered as possible centroids of the clusters, which avoids biasing the clusters to some randomly chosen points, used as the initial centroids; and (2) there is no need to provide the number of clusters to the algorithm in advance. In this step, we cluster the explanations in the three considered input spaces: *original*, consisting of raw explanation vectors; *global latent*, consisting of globally projected explanations; and *local latent*, consisting of locally projected explanations. Of course, each considered space may lead to different clusters, based on the distribution of the inputs in the corresponding space. The output of this step consists of the clusters generated for the fine-grained and the coarse-grained explanations, in each of the considered input spaces.

Comparison: In the previous steps, we generated high-level and low-level explanations for the considered inputs and processed them to generate clusters. The main evaluation step of our study is the comparison of these explanation clusters. Existing similarity or distance metrics to compare two sets of clusters, such as the MoJo metric [40], are based on the transformation of one cluster into the other. The computational complexity involved in the computation of these metrics is typically high (exponential in the worst case), but what is even more concerning in our usage scenario is that such metrics are highly sensitive to the number of clusters, and when there is a disparity in such number the distance tends to grow (similarity tends to decrease), because more transformation steps are needed to change one clustering into the other. However, a high disparity in the number of clusters is not necessarily an indicator of distance in our case. In fact, if all clusters in one set are *pure*, because their elements come from the same clusters in the other set, we deem the two clusterings very similar between each other, regardless of any disparity in the number of clusters. To capture such a notion of similarity between two sets of clusters, we define a novel metric based on *Gini Impurity (GI)* [31]. GI reflects the impurity level of a group of entities by indicating the probability that two samples from the given group have different labels, i.e., belong to different classes. So, to compute GI we need to define what are the *groups* and what are the *labels*. In our setting, when comparing two sets of clusters (low- vs high-level explanations), groups are the clusters identified in one set (source clustering), while labels are the cluster identifiers from the other set (target clustering). Hence, the impurity of a group of data D (i.e., a cluster from the source clustering) against the target clustering A can be measured as follows:

$$GI(D, A) = 1 - \sum_{i=1}^{|A|} p_{Ai}^2 \quad (1)$$

where $|A|$ is equal to the number of clusters in A , and p_{Ai} is the probability that cluster id i of clustering A occurs in dataset D . GI ranges between 0 and 1, being 0 when D contains no impurity (i.e., all its elements belong to the same cluster from A) and being minimum when D is uniformly impure (i.e., all its elements belong to a different cluster from A).

Let us consider cluster $CB1$ in Figure 4 (d). GI of $CB1$ against clustering A (Figure 4 (a)) is computed as follows: $GI(CB1, A) = 1 - \sum_{i=1}^3 p_{Ai}^2 = 1 - (1 + 0 + 0) = 0$. ($p_{A1} = 1$ is the probability of cluster $CA1$ to be found in cluster $CB1$; $p_{A2} = p_{A3} = 0$ is the probability to find $CA2, CA3$).

6

Table 1: Hyperparameters and configuration details

| | MNIST | IMDB |
|-----------------------------|----------------|--------------------|
| Number of cells per feature | 5 | 5 |
| Number of runs | 10 | 10 |
| t-SNE components | 2 | 2 |
| t-SNE perplexity | 1 | 1 |
| Similarity metric | Euclidean | Euclidean |
| Input size | 28×28 | 2000 |
| Vocabulary size | - | 10000 |
| Number of inputs | 250 | 250 |
| Target class label | "5" | "positive" |
| Number of steps for IG | 50 | 50 |
| Batch size for IG | 64 | 100 |
| Number of samples for LIME | 100 | 5000 |
| Explanation library | <i>Xplique</i> | <i>Alibi, Lime</i> |

4 Experimental Evaluation

4.1 Research Questions

RQ1 (Similarity): *How similar are high-level and low-level explanations of DL misbehaviours?*

High- and low-level approaches provide explanations about failure-inducing inputs from different perspectives. We aim to measure the similarity of these different explanations to assess to what extent they are comparable or complementary.

Metrics: To assess how similar the explanations generated by high-level and low-level techniques are, we compare the way inputs are grouped by similar explanations. To this aim, we measure (1) the number of clusters generated by each technique and (2) the similarity between these clusters according to our custom Gini Similarity (GS) metric. If two techniques produce explanations that group the inputs in a similar way (i.e., they produce nearly the same number of clusters that are pairwise very similar to each other), then they can be deemed similar; otherwise they differ, as they partition the input vectors in a different way.

RQ2 (Understandability): *How understandable are high-level and low-level explanations of DL misbehaviours?*

Since high-level techniques quantify high-level input features, while low-level approaches highlight low-level elements of the inputs, it is important to investigate whether these two types of explanations are equally understandable to humans.

Metrics: Effectively assessing the understandability of an explanation requires humans in the loop. Therefore, we designed a human study to assess if explanations are understandable and if they match the human expectations. We count the number of cases in which the explanation Matches with Human (MH), i.e., the number of times human assessors select a given explanation as possibly pointing to the cause of the misbehaviour.

4.2 Subject Systems

In our study, we consider two DL systems belonging to different domains, i.e., handwritten digit recognition and sentiment analysis, which have been widely used in the literature to generate explanations for DL systems [26, 37].

The MNIST system is a classifier which recognises handwritten digits within image inputs from the MNIST dataset [16]. This DL system accepts greyscale images as inputs and predicts the

corresponding digit classes (from 0 to 9). We considered the convolutional DNN architecture provided by Keras [3] and trained it on the the MNIST training set using Tensorflow v2. More specifically, we used its default configuration, i.e. 12 epochs, batches of size 128, a learning rate equal to 1, and the *Adam* optimizer, which achieved 99.8% test accuracy.

The IMDB system classifies textual inputs based on their sentiments. It accepts the text of a movie review as input and predicts if it is positive or negative. We adopted a convolutional DNN with an embedding layer provided by Keras [25], and we used tokenised (with vocabulary size equals to 10000) and padded text with length limited to 2000 words. We trained the model on the IMDB training set in Tensorflow v2. We used 10 epochs, batches of size 32 with early stopping, and the *Adam* optimizer, achieving 88.0% test accuracy.

4.3 Experimental Procedure

To obtain the explanations, we adopted the pipeline described in section 3 (configuration reported in Table 1).

As misbehaviour-inducing inputs, for MNIST, we used a set of 250 inputs belonging to the same ground-truth class (i.e., digit “5”), either from the test set or generated by test input generators. We used the DeepJanus [29] and Sinvad [13] test generators since they (1) belong to different families of approaches (i.e., model-based and DL generative, respectively), and (2) have been demonstrated to produce the highest ratio of valid input that preserve their ground-truth label [30]. In particular, we could obtain only 11 misclassified inputs belonging to the class digit “5” from the MNIST test set, due to high accuracy of the considered model. Therefore, we also included 35 inputs generated by DeepJanus and 204 inputs generated by Sinvad, to have enough and diverse misbehaviour-inducing inputs for our study.

For IMDB, we also used a set of 250 inputs either from the test set or generated by test input generators, all belonging to the same class (“positive” sentiment). We obtained 1924 misclassified inputs belonging to the class “positive” from the test set. Then, to have more diverse inputs we also included 663 inputs generated by the considered input generator. In particular, we used the DeepHyperion [44] tool since it supports the generation of textual inputs. Finally, we randomly selected 250 inputs from the overall set of 2587 candidate inputs.

To obtain low-level explanations, we used three open source libraries for images and textual inputs, i.e., *Xqlique* [6] for image, and *Alibi* [15] and *Lime* [27] for text.

To obtain high-level explanations, we needed high-level features for each considered domain. For MNIST inputs, we used the three high-level features defined by Zohdinasab et al. [44] in their study involving human experts, i.e., (1) *Luminosity* (*Lum*): number of light pixels in the image, obtained by counting the pixels whose value is above 127; (2) *Orientation* (*Or*): vertical orientation of the digit, obtained by computing the angular coefficient of the linear regression of the non-black pixels; (3) *Moves* (*Mov*): sum of the Euclidean distances between pairs of consecutive segments of the digit, the distance being zero when consecutive segments are connected and greater than zero when there are discontinuities. For IMDB reviews, we defined three features: (1) *Positive word count* (*Pos*): number of words in the text with positive polarity, obtained by counting the words tagged as positive in the English Opinion Lexicon [11]; (2) *Negative word count* (*Neg*): number of words in the text with negative polarity, obtained by counting the words tagged as negative in the English Opinion Lexicon; (3) *Verb count* (*Verb*): number of verbs in the text, a proxy for the text complexity, computed by counting the words with a *verb* tag, according to the part-of-speech (POS) tagging produced by the NLTK library.

Our experimental procedure consists of the following two steps: (1) high- and low-level explanations comparison, using clustering; and (2) human assessment of the explanations, by means of two surveys: one for MNIST and one for IMDB.

4.3.1 High- and Low-Level Explanations Comparison

We generate FEATURE MAPS with 1, 2, and 3 dimensions re-scaled to size 5 per dimension. We selected 5 as number of cells since the number of clusters (i.e., filled cells) with this configuration is comparable to the number of clusters obtained through low-level techniques, which allowed us to interpret each cell as a cluster without any need for an additional clustering step. We ran FM generation only once since the clusters (i.e., the cells) are deterministic.

As regards low-level explanations, we ran IG and LIME approaches 10 times each, since the corresponding clusterings are obtained through t-SNE, which is non-deterministic. To obtain the best configuration for the *t-SNE* hyperparameters (i.e., number of components and perplexity), we performed preliminary runs with different configurations and selected the one producing the highest *silhouette score*. The best configurations of the hyperparameters of IG and LIME are reported in Table 1.

For clustering the low-level explanations, we used the *Euclidean distance* to compute the similarity matrices, where distances are computed in three different vector spaces: *original*, *global latent*, and *local latent* (see Figure 3). In the original space, we consider heatmaps as vectors to be clustered. For MNIST, a heatmap is a 28×28 matrix, where each vector component e_{ij} is the contribution value of the pixel occupying the i -th row and the j -th column in the image. The matrix is flattened into a vector before applying clustering. For IMDB, a heatmap is a vector of size 10000, the size of the vocabulary used by the tokenizer (which maps each word to the corresponding one-hot encoding) and each vector component e_i is the contribution value of the i -th word in the text. In both global latent and local latent spaces, we use vectors of size 2, the same number of dimensions that was found to be optimal for the t-SNE algorithm.

4.3.2 Human Assessment of the Explanations

To determine whether an explanation is human-understandable, we asked humans to determine which of the explanations provided them meaningful information about the reason why the model misbehaves for specific inputs. We published two surveys (one for each case study) by using Qualtrics, a survey platform commonly used also for software engineering research [18,42]. To ensure the assessment quality, we selected software engineering researchers and allowed each of them to devote up to 1 week to answer the questions and take at most 1 questionnaire for each case study.

For MNIST, we created a survey with 10 questions to be answered by human assessors. To this aim, we randomly selected 10 "5" digit images from the MNIST data set which are misclassified by the considered model and we computed high-level and low-level explanations for each of them. More specifically, for the high-level explanations we reported the feature values provided by the three-dimensional FM and for the low-level explanations we visualized IG and LIME heatmaps, overlaid on the original digit images. We used 3D FEATURE MAPS for the human study to include all the available features for the high-level explanations. For each MNIST image, we showed the human assessors the possible explanations (i.e. FM, IG, and LIME) and asked them to answer the following question: "*Below is a digit "5" that was incorrectly believed to be another digit by the computer. Please, select a possible explanation for such a mistake (you are allowed to select more than one explanation if they make sense, or none of them if they do not make sense)*". Figure 5 shows a misclassified digit "5" with the three different explanations. The assessors could select the explanations that they considered as a plausible reason for the misbehaviour. In particular, the assessors were allowed to choose zero or more explanations. We collected 29 answers from the human assessors in the MNIST survey.

For IMDB, we created a survey with 10 questions by randomly selecting 10 positive reviews from the IMDB dataset, which are misclassified by the model. For the high-level explanations, we reported the feature values provided by the three-dimensional FM and the explanations generated by LIME and IG and we visualised the words with highest contributions to negative and posi-

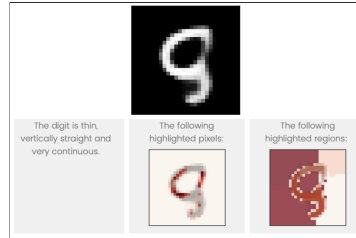


Figure 5: MNIST input from the human study for which the majority selected IG.

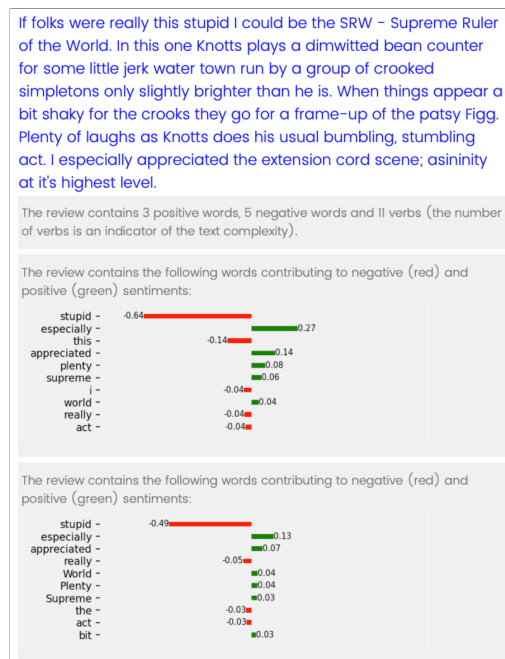


Figure 6: IMDB input from the human study for which the majority selected both FM and IG explanations.

Table 2: RQ1 - Number of clusters (NC) for high-level techniques for MNIST and IMDB.

| HL Technique | MNIST | | IMDB | |
|----------------|------------|----|--------------|----|
| | Features | NC | Features | NC |
| Feature map 3D | Mov-Lum-Or | 31 | Pos-Neg-Verb | 27 |
| Feature map 2D | Mov-Lum | 17 | Pos-Neg | 15 |
| | Lum-Or | 15 | Neg-Verb | 14 |
| | Or-Mov | 14 | Pos-Verb | 12 |
| Feature map 1D | Mov | 5 | Pos | 5 |
| | Lum | 5 | Neg | 5 |
| | Or | 5 | Verb | 5 |

Table 3: RQ1 - Number of clusters (NC) for low-level techniques for MNIST and IMDB.

| LL Technique | Input space | MNIST | IMDB |
|--------------|---------------|---------------|---------------|
| | | NC [min, max] | NC [min, max] |
| IG | Original | 44.0 [44, 44] | 35.1 [35, 36] |
| | Global Latent | 18.7 [17, 20] | 14.8 [12, 18] |
| | Local Latent | 20.5 [19, 21] | 18.2 [16, 22] |
| LIME | Original | 39 [37, 44] | 38.5 [37, 39] |
| | Global Latent | 17.7 [16, 19] | 17.3 [16, 19] |
| | Local Latent | 19.0 [17, 21] | 18.3 [16, 21] |

tive sentiments in a bar chart (i.e., we reported up to 10 most contributing words with non-zero contribution). For each text, we asked humans to answer the following question: “Below is a positive review that was incorrectly believed to be a negative review by the computer. Please, select a possible explanation for such a mistake (you are allowed to select more than one explanations if they make sense, or none of them if they do not make sense)”. Also in this case, we showed explanations provided by the FM, IG, and LIME (as shown in Figure 6). In the IMDB survey we collected 19 answers from the human assessors.

5 Results

5.1 RQ1 (Similarity)

Tables 2 and 3 report the number of clusters generated by each explanatory technique. For each low-level approach, we report the minimum and maximum number of clusters (NC) in order to show their variability due to non-determinism.

For MNIST, the number of clusters generated by FM 3D is 31, larger than the number of clusters generated by FM with any other dimensionality (i.e., 1D and 2D). Not surprisingly, the lower the dimensionality, the lower the number of clusters. In fact, all the FM dimensions were rescaled to size 5, leading to a lower number of cells for lower dimensionality maps (e.g., in 3D maps there are at most 125 cells to fill, while for 1D maps there are only 5 cells). The highest difference between high- and low-level techniques is between FM 1D and IG in the Original input space (39 clusters). Instead, the NC value for both FM 2D and 3D is comparable to the number of clusters generated by low-level techniques (resp. global/local latent space and original space).

As regards IMDB, the highest difference between high- and low-level techniques is the one between FM 1D and LIME in the Original input space, i.e., 33.5. Also for IMDB, the closest NC

Table 4: RQ1 - Comparing high-level and low-level explanations' Similarity (Sim); boldface indicate statistical significance when comparing original with latent space similarities.

| HL tech- nique | LL tech- nique | Input space | MNIST | | IMDB | |
|----------------------|----------------------|---------------|--------------------|-------------|----------------------|-------------|
| | | | Features | Sim | Features | Sim |
| Feature map 3D | IG | Original | Mov- Lum- Or | 0.70 | Pos- Neg- Verb | 0.74 |
| | | Global Latent | | 0.55 | | 0.68 |
| | | Local Latent | | 0.55 | | 0.66 |
| | LIME | Original | | 0.55 | | 0.81 |
| | | Global Latent | | 0.53 | | 0.66 |
| | | Local Latent | | 0.53 | | 0.68 |
| Feature map 2D | IG | Original | Mov- Lum | 0.71 | Pos- Neg | 0.74 |
| | | Global Latent | | 0.40 | | 0.52 |
| | | Local Latent | | 0.41 | | 0.53 |
| | LIME | Original | | 0.56 | | 0.77 |
| | | Global Latent | | 0.39 | | 0.55 |
| | | Local Latent | | 0.40 | | 0.56 |
| Feature map 2D | IG | Original | Lum- or | 0.76 | Neg- Verb | 0.78 |
| | | Global Latent | | 0.52 | | 0.60 |
| | | Local Latent | | 0.52 | | 0.60 |
| | LIME | Original | | 0.65 | | 0.82 |
| | | Global Latent | | 0.49 | | 0.61 |
| | | Local Latent | | 0.49 | | 0.61 |
| Feature map 2D | IG | Original | Mov- Or | 0.80 | Pos- Verb | 0.77 |
| | | Global Latent | | 0.54 | | 0.57 |
| | | Local Latent | | 0.54 | | 0.59 |
| | LIME | Original | | 0.64 | | 0.81 |
| | | Global Latent | | 0.52 | | 0.61 |
| | | Local Latent | | 0.51 | | 0.63 |
| Feature map 1D | IG | Original | Mov | 0.83 | Pos | 0.80 |
| | | Global Latent | | 0.49 | | 0.64 |
| | | Local Latent | | 0.52 | | 0.65 |
| | LIME | Original | | 0.68 | | 0.83 |
| | | Global Latent | | 0.50 | | 0.66 |
| | | Local Latent | | 0.54 | | 0.69 |
| Feature map 1D | IG | Original | Lum | 0.78 | Neg | 0.82 |
| | | Global Latent | | 0.47 | | 0.63 |
| | | Local Latent | | 0.53 | | 0.65 |
| | LIME | Original | | 0.69 | | 0.84 |
| | | Global Latent | | 0.45 | | 0.67 |
| | | Local Latent | | 0.50 | | 0.67 |
| Feature map 1D | IG | Original | Or | 0.94 | Verb | 0.87 |
| | | Global Latent | | 0.82 | | 0.71 |
| | | Local Latent | | 0.83 | | 0.73 |
| | LIME | Original | | 0.89 | | 0.89 |
| | | Global Latent | | 0.82 | | 0.72 |
| | | Local Latent | | 0.83 | | 0.74 |

values between the two classes of techniques is produced by FM 2D and low-level techniques in the global/local latent space. The number of clusters for FM 3D is comparable to low-level techniques in the original space.

For both case studies, IG produces a more stable number of clusters in the original space due to its deterministic nature (i.e. it is based on a mathematical equation which computes linear interpolation from a baseline), whereas there are some variations for global and local latent spaces due to the non-deterministic nature of the *t-SNE* algorithm. Instead, LIME always generates slightly different explanations across multiple runs, because of the randomness introduced by input sampling and surrogate model training during the explanation computation, and thus produces a variable number of clusters.

Table 4 reports the similarity between clusterings. To assess the statistical significance of the comparisons between different configurations, we performed the Mann-Whitney U-test and measured the effect size by means of the Vargha-Delaney’s A12 statistic [1]. For both MNIST and IMDB (fifth and seventh columns), the similarity values between high- and low-level techniques are almost always significantly higher (p -value < 0.05 , large effect size) when considering the original input space, rather than global and local latent spaces within the same configuration (i.e., up to 0.34 more in the comparison between FM 1D and IG for MNIST). This result may be due to the fact that the heatmaps in the original space retain more feature-related information than their projections into the latent space. For MNIST, the highest similarity is achieved between FM 1D (Orientation) and IG in the original space (0.94). For IMDB, the highest clustering similarity is between FM that considers the Verb feature only and LIME explanations in the original space (0.89). The most different partitions between low- and high-level techniques (i.e., lowest similarity) are obtained when considering FM 2D and low-level explanations in the latent space (either local or global). In fact, for MNIST we have the lowest similarity when considering Mov-Lum maps and LIME in the global latent space (0.39), whereas for IMDB the lowest similarity is between Pos-Neg maps and IG in the global latent space.

On average, the similarity between FM and IG, across input spaces and subjects, is 0.65, while it is 0.64 with LIME. These are relatively low values. As a reference, such similarity values are obtained when a cluster with 100 elements contains 23 impure elements, i.e., 23 elements that are assigned a different cluster by the other technique (see Equation 2).

RQ1: High-level and low-level techniques partition inputs in different ways. Despite FM 2D and low-level techniques in the global latent space produce nearly the same number of clusters, those show low similarity. Likewise, FM 1D produces clusterings similar to the ones obtained by low-level techniques in the original space, but with very different number of clusters.

5.2 RQ2 (Understandability)

Table 5 shows the results extracted from questionnaires for MNIST and IMDB. Each column reports the number of assessors who chose the explanation provided by each considered technique (choices were not mutually exclusive). The last column highlights the cases for which no explanation was selected by the assessor. Since each assessor can select 0 to 3 explanations in each question, the sum of the values in each row does not correspond to the number of assessors.

Table 5 (left) reports the answers we collected for MNIST. IG was selected more than other explanations 4 times, e.g., for Q2 IG has been chosen 18 times more than the others. The explanations by FM are selected more than the other explanations 4 times. LIME was selected more than the others in only 2 questions. We can conclude that there is no technique clearly more understandable than the others for explaining digit misclassifications, according to human assessors. Overall, the explanations generated by IG match more frequently with the human expectations, as they have been selected 30% of the times by the assessors. FM explanations and LIME heatmaps have been

Table 5: RQ2 - Number of Matches with Human Explanations (MH); ‘None’ indicates the number of cases when no match was found.

| Q# | MH | | | | | | | |
|-----|-------|----|------|------|-------|-----|------|------|
| | MNIST | | | | IMDB | | | |
| | FM 3D | IG | LIME | None | FM 3D | IG | LIME | None |
| Q1 | 12 | 2 | 10 | 8 | 2 | 13 | 3 | 3 |
| Q2 | 5 | 23 | 5 | 1 | 5 | 3 | 3 | 9 |
| Q3 | 4 | 7 | 9 | 12 | 11 | 17 | 8 | 0 |
| Q4 | 6 | 7 | 5 | 14 | 6 | 15 | 3 | 2 |
| Q5 | 3 | 15 | 2 | 11 | 10 | 14 | 7 | 1 |
| Q6 | 7 | 6 | 4 | 14 | 12 | 15 | 8 | 0 |
| Q7 | 7 | 11 | 7 | 10 | 0 | 14 | 5 | 3 |
| Q8 | 9 | 5 | 8 | 13 | 11 | 14 | 12 | 1 |
| Q9 | 5 | 1 | 9 | 17 | 6 | 8 | 4 | 4 |
| Q10 | 13 | 10 | 5 | 8 | 11 | 12 | 7 | 2 |
| Sum | 71 | 87 | 64 | 108 | 74 | 125 | 60 | 25 |

selected 24% and 22% of the times. However, in more than one third of the answers the assessors chose none of the explanations. This result suggests that there is large room for improvement in explaining image misclassifications. Table 5 (right) reports the answers for IMDB. IG was selected more than the other techniques 9 times, while FM was selected more than the others in the remaining case (i.e., Q2). Instead, LIME explanations were never selected more often than the other techniques. IG was selected the highest number of times by human assessors, i.e., 66% , followed by FM with 39% and LIME with 31%. Explanations were more understandable for textual inputs, than for handwritten digit images. In fact, only 13% of the times the assessors did not choose any explanation for a question on IMDB. For MNIST, only for 2 questions we have a majority, i.e., more than half of the assessors selected the same answer. Instead, for IMDB in 8 questions we have answers selected by the majority of assessors. This indicates that the explanations provided for IMDB find more consensus among humans than the explanations provided for MNIST.

RQ2: Both high- and low-level techniques produce human-understandable explanations of misbehaviours of text classifiers. In particular, IG explanations were selected 125 times out of 190 answers for IMDB. On the other hand, the explanations for misbehaviours of image classifiers poorly matched with the human judgement. In fact, in more than one third of the answers the assessors chose none of the proposed explanations for MNIST. On the remaining two thirds, high- and low-level explanations are chosen approximately the same number of times.

5.3 Discussion

The answers provided by the human assessors offer several insights that we discuss qualitatively in the following.

- For MNIST, the two most understandable explanations are IG and FM. Not only are these explanations at different levels (i.e., low- and high-level, respectively), but they show also some degree of complementarity. In fact, for Q1, Q8 and Q9, FM were selected more times than IG with a MH difference higher than 4. On the other hand, for Q2, Q5 and Q7, IG was selected more times than FM. Remarkably, for question Q2 (reported in Figure 5), IG was selected 18

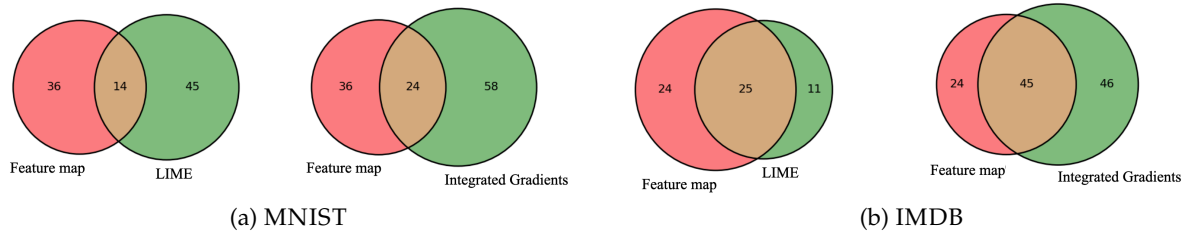


Figure 7: Comparison of the number of matches with human between high-level and low-level techniques.

times more than FM. In this case, the assessors did not find the FM explanation useful, while IG highlights the pixels that make the upper part of the five look like a nine, i.e., the leftmost pixels, making the upper part round, and the rightmost pixels, close the circle.

- For IMDB, FM and IG explanations are the most chosen by human assessors. In particular, for half of the questions, IG and FM were chosen by at least 10 assessors and the difference in MH was lower than or equal to 4. For instance, Figure 6 reports Q8, where the assessors selected FM 11 times and IG 14 times. In this case, the assessors acknowledged that the review actually contained more negative than positive words as reported by the FM explanation. At the same time, the assessors probably agreed with the IG heatmap that the word “stupid” contained in the text can negatively affect the prediction.
- For some questions, all the explanations were judged as poorly understandable by humans, who did not choose any of the provided explanations as possible causes of the misbehaviour. For MNIST, Q9 can be considered the most challenging question, since only 15 times out of 87 times an explanation was selected, while the majority of the assessors selected none of the explanations. Similarly, in the IMDB survey, Q2 is the question with the lowest number of selected explanations, since the assessors chose an explanation only 11 times out of 57, while the majority chose none.
- As shown in Figure 7 (a), for MNIST, assessors selected only FM 36 times and only LIME 45 times while they selected both explanations 14 times. Similarly, they selected only FM 36 times, only IG 58 times and both FM and IG 24 times. This suggests that high-level and low-level explanations for MNIST digit misclassifications are highly complementary. Figure 7 (b) shows that for IMDB, most of the times both high- and low-level explanations have been selected together. However, the number of times assessors selected one of the two and not the other remains quite high, which indicates substantial complementarity.

In summary, in several instances either high- or low-level explanations look understandable for humans and provide hints on what are the properties of the inputs that make the DNN misbehave. In most cases, either high- or low-level explanations, but not both, were chosen as useful explanations, indicating that such explanations are highly complementary. On the other hand, there is still a large fraction of cases in which none of the two are deemed useful by humans, which points to the need for better DNN explanations.

5.4 Threats to Validity

Construct Validity: FM depends on the features selected to generate the map dimensions. For feature selection, we relied on previous work [44], while for the combination of features, we exhaustively considered all available dimensions and all their combinations. Moreover, the selection of different hyper-parameters for our study can affect the final results. To reduce this risk, we developed our framework in order to be configurable and we performed some preliminary runs with different hyper-parameters to select the best configuration.

External Validity: The selection of subjects could be a threat to external validity. Therefore, we considered two different domains, i.e., image and text. To represent different types of explanations, we considered one state-of-the-art, openly available technique for each explanation type.

Conclusion Validity: The stochastic nature of some of the considered low-level techniques, the dimensionality reduction approach and the clustering algorithm might affect the final outcome. To mitigate this threat, we ran these techniques multiple times and reported average value across runs.

Reproducibility: To make our results reproducible, We make our implementation and the experimental data openly accessible online: <https://github.com/testingautomated-usi/unboxer>.

6 Related Work

6.1 Comparison of Explainable AI techniques

During the last decade, a variety of techniques have been proposed to provide explanations for DL [19, 23, 26, 33–35, 37]. Although these techniques share a common goal, they differ in various aspects. Therefore, there is a growing need for comparisons of such explanatory techniques.

Tjoa and Guan [39] proposed a survey on explanatory techniques and categorised them based on their interpretability. Their analysis is based solely on the mathematical properties of each technique and on experimental results available from the literature, while we explicitly ran the explanatory techniques and compared them both quantitatively and qualitatively. Samek et al. [32] provided a survey on explanatory techniques applicable to DNNs. Unlike our work, they did not evaluate the similarity among explanations produced by different techniques and they analysed human interpretability of explanations without conducting any human assessment involving humans. In fact, they measured the size of the explanation file generated by each technique as a measure of human interpretability. Ledel and Berthold [17] examined the quality of explanations provided by LIME and SHAP for a different task, i.e., bug prediction, through an empirical study involving human assessors. Unlike them, we compared different techniques also quantitatively.

To the best of our knowledge, existing empirical comparisons only considered low-level explanations, without evaluating high-level explanations.

6.2 Explainable AI for DL Testing

In the literature, both low- and high-level explanations have been leveraged to interpret the output of DL testing. As for low-level explanations, HUDD [5] identifies root causes of DNN failures by using heatmap explanations generated by the LRP technique and then clustering inputs with similar heatmaps. Stocco et al. [36] proposed ThirdEye, a monitor for autonomous driving systems which relies on heatmaps produced by the SmoothGrad [34] technique to predict unsafe conditions in advance. High-level explanations provided by FM have been used for different testing tasks. DEEPHYPERION [44–46] is a test generator that leverages FM for extensively exploring the feature space and finding failure-inducing inputs with different characteristics. Nguyen et al. [24] used FM for test selection. FM has also been used to assess test suite adequacy, measured as the number of cells covered by the test inputs [2].

Existing works either considered low- or high-level explanations. Instead, we consider both explanation levels and compare them to assess their usefulness and understandability in explaining DL misbehaviours.

7 Conclusions and Future Work

In this paper, we provide an in-depth comparison of explanatory techniques for DL misbehaviours. Our empirical results show that high- and low-level explanations are both understandable for humans, but provide different and complementary insights. Our human study suggests that current explanations are not always satisfactory, as they may not provide human-interpretable causes of misbehaviours.

As future work, we will perform a larger evaluation, involving additional explanatory techniques, clustering solutions, DL tasks, and experts from the industry.

References

- [1] Andrea Arcuri and Lionel Briand. A hitchhiker’s guide to statistical tests for assessing randomized algorithms in software engineering. *Software Testing, Verification and Reliability*, 24(3):219–250, 2014.
- [2] Matteo Biagiola, Stefan Klikovits, Jarkko Peltomaki, and Vincenzo Riccio. SBFT tool competition 2023 - cyber-physical systems track. In *16th IEEE/ACM International Workshop on Search-Based And Fuzz Testing, SBFT 2023, Melbourne, Australia, May 14, 2023*, 2023.
- [3] François Chollet. Simple mnist convnet. https://github.com/keras-team/keras-io/blob/master/examples/vision/mnist_convnet.py, 2020.
- [4] Martin Ester, Hans-Peter Kriegel, Jörg Sander, and Xiaowei Xu. A density-based algorithm for discovering clusters in large spatial databases with noise. In *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining, KDD’96*, page 226–231. AAAI Press, 1996.
- [5] Hazem Fahmy, Fabrizio Pastore, Mojtaba Bagherzadeh, and Lionel Briand. Supporting deep neural network safety analysis and retraining through heatmap-based unsupervised learning. *IEEE Transactions on Reliability*, 70(4):1641–1657, 2021.
- [6] Thomas Fel, Lucas Hervier, David Vigouroux, Antonin Poche, Justin Plakoo, Remi Cadene, Mathieu Chalvidal, Julien Colin, Thibaut Boissin, Louis Bethune, Agustin Picard, Claire Nicodeme, Laurent Gardes, Gregory Flandin, and Thomas Serre. Xplique: A deep learning explainability toolbox. *Workshop on Explainable Artificial Intelligence for Computer Vision (CVPR)*, 2022.
- [7] Brendan J Frey and Delbert Dueck. Clustering by passing messages between data points. *science*, 315(5814):972–976, 2007.
- [8] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6904–6913, 2017.
- [9] Riccardo Guidotti, Anna Monreale, Salvatore Ruggieri, Franco Turini, Fosca Giannotti, and Dino Pedreschi. A survey of methods for explaining black box models. *ACM computing surveys (CSUR)*, 51(5):1–42, 2018.
- [10] Geoffrey E Hinton and Sam Roweis. Stochastic neighbor embedding. *Advances in neural information processing systems*, 15, 2002.
- [11] Minqing Hu and Bing Liu. Opinion lexicon. <https://www.cs.uic.edu/~liub/FBS/sentiment-analysis.html>, 2004.
- [12] Nargiz Humbatova, Gunel Jahangirova, Gabriele Bavota, Vincenzo Riccio, Andrea Stocco, and Paolo Tonella. Taxonomy of real faults in deep learning systems. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering, ICSE ’20*, page 1110–1121. ACM, 2020.
- [13] Sungmin Kang, Robert Feldt, and Shin Yoo. Sinvad: Search-based image space navigation for dnn image classifier test input generation. In *Proceedings of the IEEE/ACM 42nd International Conference on Software Engineering Workshops*, pages 521–528, 2020.
- [14] Ronald S King. *Cluster analysis and data mining: An introduction*. Mercury Learning and Information, 2015.

- [15] Janis Klaise, Arnaud Van Looveren, Giovanni Vacanti, and Alexandru Coca. Alibi explain: Algorithms for explaining machine learning models. *Journal of Machine Learning Research*, 22(181):1–7, 2021.
- [16] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [17] Benjamin Ledel and Steffen Herbold. Studying the explanations for the automated prediction of bug and non-bug issues using lime and shap. *arXiv preprint arXiv:2209.07623*, 2022.
- [18] Mario Linares-Vásquez, Boyang Li, Christopher Vendome, and Denys Poshyvanyk. Documenting database usages and schema constraints in database-centric applications. In *Proceedings of the 25th International Symposium on Software Testing and Analysis*, pages 270–281, 2016.
- [19] Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. *Advances in neural information processing systems*, 2017.
- [20] J MacQueen. Classification and analysis of multivariate observations. In *5th Berkeley Symp. Math. Statist. Probability*, pages 281–297, 1967.
- [21] Pavan Rajkumar Magesh, Richard Delwin Myloth, and Rijo Jackson Tom. An explainable machine learning model for early detection of parkinson’s disease using lime on datscan imagery. *Computers in Biology and Medicine*, 126:104041, 2020.
- [22] Tim Miller. Explanation in artificial intelligence: Insights from the social sciences. *Artificial intelligence*, 267:1–38, 2019.
- [23] Grégoire Montavon, Sebastian Lapuschkin, Alexander Binder, Wojciech Samek, and Klaus-Robert Müller. Explaining nonlinear classification decisions with deep taylor decomposition. *Pattern recognition*, 65:211–222, 2017.
- [24] Vuong Nguyen, Stefan Huber, and Alessio Gambi. Salvo: Automated generation of diversified tests for self-driving cars from existing maps. In *2021 IEEE International Conference on Artificial Intelligence Testing (AITest)*, pages 128–135. IEEE, 2021.
- [25] Mark Omernick and François Chollet. Text classification from scratch. https://github.com/keras-team/keras-io/blob/master/examples/nlp/text_classification_from_scratch.py, 2020.
- [26] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. “why should i trust you?” explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1135–1144, 2016.
- [27] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. Lime. <https://github.com/marcotcr/lime>, 2017.
- [28] Vincenzo Riccio, Gunel Jahangirova, Andrea Stocco, Nargiz Humbatova, Michael Weiss, and Paolo Tonella. Testing machine learning based systems: a systematic mapping. *Empir. Softw. Eng.*, 25(6):5193–5254, 2020.
- [29] Vincenzo Riccio and Paolo Tonella. Model-based exploration of the frontier of behaviours for deep learning system testing. In *Proceedings of the ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering, ESEC/FSE ’20*. ACM, 2020.
- [30] Vincenzo Riccio and Paolo Tonella. When and why test generators for deep learning produce invalid inputs: an empirical study. In *Proceedings of the IEEE/ACM International Conference on Software Engineering, ICSE ’23*. IEEE/ACM, 2023.

- [31] Lior Rokach and Oded Maimon. Top-down induction of decision trees classifiers-a survey. *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, 35(4), 2005.
- [32] Wojciech Samek, Grégoire Montavon, Sebastian Lapuschkin, Christopher J Anders, and Klaus-Robert Müller. Explaining deep neural networks and beyond: A review of methods and applications. *Proceedings of the IEEE*, 2021.
- [33] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, pages 618–626, 2017.
- [34] Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. Smoothgrad: removing noise by adding noise. *arXiv preprint arXiv:1706.03825*, 2017.
- [35] Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin Riedmiller. Striving for simplicity: The all convolutional net. *arXiv preprint arXiv:1412.6806*, 2014.
- [36] Andrea Stocco, Paulo J Nunes, Marcelo d’Amorim, and Paolo Tonella. Thirdeye: Attention maps for safe autonomous driving systems. In *Proceedings of 37th IEEE/ACM International Conference on Automated Software Engineering, ASE*, volume 22, 2022.
- [37] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *International conference on machine learning*, pages 3319–3328. PMLR, 2017.
- [38] Chakkrit Tantithamthavorn and Jirayus Jiarpakdee. *Explainable AI for Software Engineering*. Monash University, 2021. Retrieved 2021-05-17.
- [39] Erico Tjoa and Cuntai Guan. A survey on explainable artificial intelligence (xai): Toward medical xai. *IEEE transactions on neural networks and learning systems*, 32(11):4793–4813, 2020.
- [40] Vassilios Tzerpos and Richard C. Holt. Mojo: A distance metric for software clusterings. In *Sixth Working Conference on Reverse Engineering, WCRE '99, Atlanta, Georgia, USA, October 6-8, 1999*, page 187, 1999.
- [41] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008.
- [42] Melina Vidoni. Evaluating unit testing practices in r packages. In *2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE)*, pages 1523–1534. IEEE, 2021.
- [43] Jie M Zhang, Mark Harman, Lei Ma, and Yang Liu. Machine learning testing: Survey, landscapes and horizons. *IEEE Transactions on Software Engineering*, 2020.
- [44] Tahereh Zohdinasab, Vincenzo Riccio, Alessio Gambi, and Paolo Tonella. Deephyperion: exploring the feature space of deep learning-based systems through illumination search. In *Proceedings of the 30th ACM SIGSOFT International Symposium on Software Testing and Analysis*, 2021.
- [45] Tahereh Zohdinasab, Vincenzo Riccio, Alessio Gambi, and Paolo Tonella. Efficient and effective feature space exploration for testing deep learning systems. *ACM Transactions on Software Engineering and Methodology*, 2023.
- [46] Tahereh Zohdinasab, Vincenzo Riccio, and Paolo Tonella. Deepatash: Focused test generation for deep learning systems. In *Proceedings of the 32nd ACM SIGSOFT International Symposium on Software Testing and Analysis*, 2023.