
TOPOMAP: A FEATURE-BASED SEMANTIC DISCRIMINATOR OF THE TOPOGRAPHICAL REGIONS IN THE TEST INPUT SPACE

A PREPRINT

✉ **Gianmarco De Vita**
Università della Svizzera italiana
Lugano, Switzerland
gianmarco.de.vita@usi.ch

✉ **Nargiz Humbatova**
Università della Svizzera italiana
Lugano, Switzerland
nargiz.humbatova@usi.ch

✉ **Paolo Tonella**
Università della Svizzera italiana
Lugano, Switzerland
paolo.tonella@usi.ch

September 4, 2025

ABSTRACT

Testing Deep Learning (DL)-based systems is an open challenge. Although it is relatively easy to find inputs that cause a DL model to misbehave, the grouping of inputs by features that make the DL model under test fail is largely unexplored. Existing approaches for DL testing introduce perturbations that may focus on specific failure-inducing features, while neglecting others that belong to different regions of the feature space.

In this paper, we create an explicit topographical map of the input feature space. Our approach, named TOPOMAP, is both black-box and model-agnostic as it relies solely on features that characterise the input space. To discriminate the inputs according to the specific features they share, we first apply dimensionality reduction to obtain input embeddings, which are then subjected to clustering. Each DL model might require specific embedding computations and clustering algorithms to achieve a meaningful separation of inputs into discriminative groups. We propose a novel way to evaluate alternative configurations of embedding and clustering techniques. We used a deep neural network (DNN) as an approximation of a human evaluator who could tell whether a pair of clusters can be discriminated based on the features of the included elements. We use such a DNN to automatically select the optimal topographical map of the inputs among all those that are produced by different embedding/clustering configurations.

The evaluation results show that the maps generated by TOPOMAP consist of distinguishable and meaningful regions. In addition, we evaluate the effectiveness of TOPOMAP using mutation analysis. In particular, we assess whether the clusters in our topographical map allow for an effective selection of mutation-killing inputs. Experimental results show that our approach outperforms random selection by 35% on average on killable mutants; by 61% on non-killable ones.

Keywords Software testing · Deep learning testing · Test input clustering

1 Introduction

As Deep Learning (DL)-based software solutions are spreading across a plethora of different application domains, testing of DL systems is progressively fostering related research [44] to ensure the integrity and reliability of those systems. In particular, one of the main goals is to identify the most effective test cases, i.e., inputs which bear features that are more likely to expose failures.

Existing approaches to DL input generation [12, 40, 55, 56] often apply small input perturbations until a misbehaviour is exposed. However, these approaches may focus the test generation effort in a specific promising region of the input space, neglecting other, possibly less easily reachable regions, where different misbehaviours might be exposed. To the best of our knowledge, the only approach that explores the input feature space during test generation is DeepHyperion [62, 63]. However, this approach requires substantial human effort for the definition and quantification of the feature space and

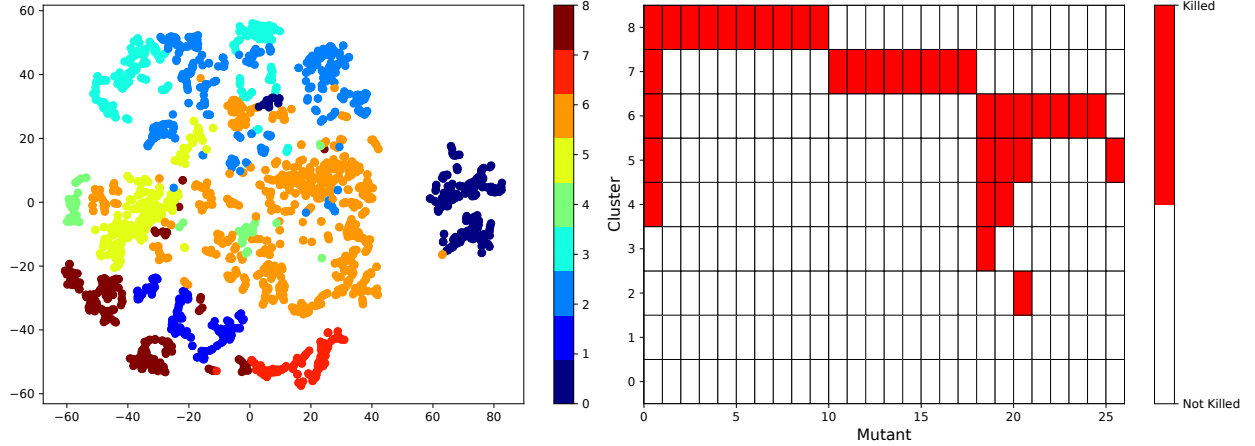


Figure 1: Topographical map of aggregations of test input clusters killing various mutants.

requires a parameterised, model-based input generator that can be steered toward the exhaustive exploration of the input space.

In this paper, we propose TOPOMAP, a novel approach for fully automated creation of a topographical map of inputs, which are grouped into regions according to their features. We validated the resulting map through a human study and showed that it can be used to expose different DL failures.

The input to a DL system is usually complex, in particular due to the high dimensionality of the data being processed [53]. Hence, the first step of our approach is the automated extraction of an input embedding that has lower dimensionality than the input space, while still retaining the key features that characterise each input and trigger a specific behaviour of the DL system under test. Second, we group the inputs according to the features represented in their embeddings using unsupervised learning methods, which infer the membership of the instances to a certain region (i.e., a cluster), based on the characteristics of the lower-dimensional data. As multiple embeddings and clustering algorithms/configurations are possible, we automate the selection of the optimal set of clusters by mimicking human assessment using a Deep Neural Network (DNN). A human may approach the quality assessment process for the output of a clustering by comparing pairs of clusters to decide if they contain discriminative features that make them meaningful and distinguishable groups of inputs. This task can be automated using a DNN, which is trained to classify the inputs according to the pseudo-labels (regions) produced by clustering. We assume that a high accuracy of this DNN represents a reliable indicator of the presence of meaningful, discriminative features in the two clusters being compared. The output of this phase is a clustering of the inputs, represented as a topographical map (see Fig. 1, left). It should be noticed that construction of the topographical map is black-box and model-agnostic, in that it is based uniquely on the analysis of the input domain, without considering any internal property of the DL model under test. Actually, the resulting topographical map is DL model-independent and can be reused across different model architectures operating on the same data.

The results indicate that the topographical maps constructed by TOPOMAP contain regions that are clearly discernible and grounded in meaningful feature selection. Furthermore, the human study suggests that participants’ judgements are consistent with those of TOPOMAP when separating inputs into map regions, particularly in cases where the dataset images are easy to interpret.

We empirically evaluated whether aggregations of TOPOMAP’s clusters are indicative of the presence of DL failures. To do this, we applied mutation analysis and aggregated the clusters based on their ability to kill a mutated DL system. Fig. 1 (right) shows an example in which aggregations made up of only two clusters or fewer are capable of killing the available mutants. The plot shows how the inputs responsible for triggering a misbehaviour in the tested DL system can be found to be concentrated in a few clusters, sometimes even just one. Finally, comparing side by side the two panels of Fig. 1, we can notice how the killing clusters exhibit some spatial proximity, e.g. in the case of the pairs of clusters (6, 7) and (5, 8), which supports our conjecture on the existence of a topographical structure behind the failure-inducing inputs. Experimental results show that our approach can identify mutation killing inputs very efficiently achieving a killing probability of 100% with just 9% of the inputs, while random selection of the same number of inputs leads to a much lower killing (hence failure) probability of 65%. This work aims to be a first stepping stone for a model-agnostic segmentation of the input space and the subsequent identification of critical regions where to locate and generate new test input instances that have a high chance of exposing deficiencies of the DL model under test.

The paper is organised as follows: Sec. 2 illustrates the recent advances in research on related fields. Sec. 3 introduces some background on embedding computation and clustering. Sec. 4 describes our approach, followed by its empirical assessment in Sec. 5. Sec. 6 analyses the results of the experiments, and finally, we draw our conclusions in Sec. 7.

2 Related Work

In this section, we consider recent developments in research areas related to TOPOMAP: black-box testing of DNNs, prioritisation of the inputs exposing misbehaviours in a DL system, and analysis of faults in DL systems. Finally, we compare our work with the most related one, DeepHyperion [62, 63].

2.1 Black-box Testing of DNNs

Aghababaeyan et al. [6] propose a black-box approach to DNN testing that focuses on measuring diversity in test inputs. The authors selected 3 different diversity metrics for image-based datasets and evaluated their correlation with the fault detection capability of test inputs. The fault types are defined by applying a clustering technique to mispredicted inputs and assuming that each discovered cluster characterises a distinct DNN fault type. Their results show a positive and statistically significant correlation between geometric diversity metric [34] and faulty regions in the DNN input space. The existence of such correlation supports our hypothesis on the effectiveness of input selection based on specific combinations of input features. However, in our work we do not rely on diversity metrics, which are necessarily domain-specific (Aghababaeyan et al. [6] investigated only image-based datasets). Our approach is more general, as it aims to separate test inputs into discriminative regions that emulate human-grade assessment without making any assumption on the existence of domain-specific diversity metrics. As a consequence, our experimental evaluation spans multiple domains, beyond image classification.

2.2 Input Prioritisation

DL input prioritisation is a widely studied area, due to the high cost of labelling new inputs, which is often done manually and requires deep domain knowledge. Thus, when prioritising the inputs to test a DL system, it is important to automatically identify those that have high misbehaviour-revealing potential.

Uncertainty quantifiers such as DeepGini [17], Vanilla Softmax [21, 41], Prediction-Confidence Score (PCS) [59], and Entropy [48], estimate the uncertainty of a DL model based on the outputs of the *Softmax* layer, which makes these approaches to be only applicable to classification systems. Monte-Carlo Dropout [19], however, can be used both for classification and regression tasks, as it uses the dropout layers of a model during the prediction phase to obtain multiple results for the same input and it evaluates uncertainty based on the discrepancies in the predicted values. Inputs with higher uncertainty are deemed as likely to be mispredicted.

Another family of approaches are based on surprise adequacy (SA) [30]. These techniques focus on measuring how surprising an input is for a DNN, i.e., how different the DNN’s activations produced on this input are, with respect to the activations observed for the training dataset. Among those are Likelihood-based SA [30], Distance-based SA [30], and Mahalanobis-Distance Based SA [31]. The intuition behind the concept of SA is that the more surprising the input to the DNN, the higher its ability to trigger failures.

The last group of approaches relies on the neural coverage. Similar to traditional code coverage measures, metrics such as Neuron Coverage (NC) [43], k-Multisection NC [37], Neuron Boundary Coverage [37], Strong Neuron Activation Coverage [37] and Top-k NC [37] prioritise inputs based on their capability of triggering different DNN neuron activation patterns.

All these techniques require access to the internal components of a DNN, such as the output of the softmax layer or neural activations. Our approach, however, is black-box and is practically model-agnostic as it operates on the input features only. We create a discriminative topographical map of the inputs and we conjecture that different regions in such map can be aggregated to expose different DL failures.

2.3 DL Faults

There is a number of studies that investigated faults in DL systems. In particular, Humatova et al. performed an analysis of different sources of DL faults (StackOverflow, GitHub, expert interviews) to produce a comprehensive taxonomy of real faults in DL systems [23]. They define a DL fault as an inadequacy of a DL component to complete the task at hand that was caused by a human mistake during DL development or training.

The works by Zhang et al. [60] and Islam et al. [25] also mined faults from StackOverflow and GitHub platforms. Zhang et al. [60] focused on a set of applications that were developed using the Tensorflow platform and produced a set of 175 generic programming and DL-specific faults. Islam et al. [25] considered a larger set of frameworks such as Theano, Caffe, Keras, TensorFlow, and PyTorch and collected 447 different faults.

DeepCrime [24] is a mutation testing tool designed for automated seeding of artificial faults (mutations) into DL systems. In DeepCrime the authors propose 35 and implement 24 *source level* mutation operators that target different aspects of the development and training of DL systems. This set of operators was extracted from an existing taxonomy of real faults in deep learning systems [23] and was complemented with the issues found in the replication packages for the studies by Islam et al. [25] and Zhang et al. [60].

None of the existing works on the nature of DL faults attempted to characterise the input space regions that are capable of triggering a DL fault. The empirical evaluation of our approach relies on DeepCrime’s mutants to draw such a connection.

2.4 DeepHyperion

DeepHyperion [62, 63] is the most similar approach to TOPOMAP available in the literature. DeepHyperion creates a feature map of the inputs based on a quantification of the structural and behavioural features that characterise the application domain. For instance, the structural features of the image of a digit may include its luminosity and orientation, while the behavioural (resp. structural) features of a self-driving car may include its maximum distance from the centre of the lane (resp. the road curvature). The definition of such features and their quantification require profound expert domain knowledge. The authors of DeepHyperion propose a methodology for the manual identification of relevant features and definition of metrics that quantify them. While we share with them the goal of finding discriminative features in the input space, our approach is fully automated. Indeed, one key design decision behind our approach is the usage of a DNN to mimic the human’s discrimination capabilities. As a consequence, our approach has a wide range of applicability, which goes beyond the subjects considered in the evaluation of DeepHyperion.

3 Background

Our approach builds upon two well-established techniques: embedding computation and clustering. Below, we briefly summarise each to clarify their role in our pipeline.

3.1 Embedding Computation

Embedding models are important tools for transforming high-dimensional inputs such as images or text into more compact vector representations that capture meaningful semantic or structural information. Transforming inputs into an embedding representation enables downstream models, such as clustering algorithms, to process and compare inputs efficiently. In our approach, we explore the potential of a range of different embedding models in shaping topographical maps of input space. Among the existing embedding procedures, we consider: (i) Principal Component Analysis (PCA) [22] that returns a set of components obtained as linear combinations of the original data set features; (ii) Uniform Manifold Approximation and Projection (UMAP) [46], which performs a non-linear dimension reduction on the data set; (iii) t-distributed Stochastic Neighbour Embedding (t-SNE) [52] - a non-linear dimension reduction method optimised for two or three dimensions; (iv) Singular Value Decomposition (SVD) [54], which projects data onto a lower-dimensional space while retaining the most significant singular values; (v) Linear Discriminant Analysis (LDA) [61], which identifies the feature subspace by maximising the separation between classes.

In particular, PCA produces a new space of orthogonal variables through a linear combination of potentially complementary features. This new space is usually the same size as the original feature space and consists of a coordinate system orientated in the directions expressing the largest variance of the data. The variables that explain the highest proportion of variance are selected (the common rule of thumb is 90% of variance [20, 28]). This reduces the dimensionality of the data while retaining most of the variability. On the other hand, UMAP captures both local and global non-linear data relationships by providing a topological representation of the data. This is achieved by connecting each data point to its nearest neighbours and identifying connectivity patterns across these local subsets. The ensemble of these subsets is then mapped onto a lower-dimensional approximation aiming at preserving the detected data relationships. Similarly, t-SNE is a non-linear dimensionality reduction technique that aims to map data onto a very low-dimensional space, usually two or three dimensions, so that data points can be visualised. Another approach, named SVD, decomposes the original data matrix into three matrices without relying on a covariance matrix. Like PCA, SVD involves eigen-decomposition. In fact, PCA can be considered a special case of SVD. Finally, LDA is a linear embedding procedure which assumes that all the classes in the dataset share the same covariance matrix.

We also considered other embedding techniques, but these were only employed on a limited subset of the available data scenarios due to certain input-specific characteristics that could not be generalised to all scenarios. In particular, we considered (vi) Isomap [51], a non-linear dimensionality reduction technique aiming at preserving the global geometry of the data; (vii) Autoencoder (AE), a type of a DNN that learns a dimensionally efficient (encoded) representation of the data, which can be subsequently used by other methods or eventually decoded back into a form that is roughly similar to the original; (viii) Word2vec [39], a model that produces vector representations of words that preserve their proximity to neighbours within a corpus of text.

Isomap first creates a neighbourhood graph for all the points in the dataset, computing the shortest path between each pair of points to approximate the actual geodesic distances, i.e. the number of edges connecting two nodes, in the manifold. The data is then mapped onto a lower-dimensional coordinate space, ensuring that the distances remain unchanged. Although effective, this technique does not scale for very large datasets [38]. Alternatively, AEs are often used to encode inputs into a lower-dimensional representation. To integrate them into our empirical settings, we use domain-specific implementations for our subjects [5, 13, 50]. Finally, Word2vec maps words that appear in similar contexts or are semantically related to each other as relatively close vectors, with high cosine similarity. This embedding is specifically designed to be applied to textual input.

Other potential embedding procedures for the input data include extracting the activation vectors of the last hidden layers of the DNN under test (which would render the approach non-black box) or a feature extractor DNN. However, these approaches would introduce an additional layer of complexity, as the obtained embedding would depend not only on the data itself, but also on the proper training of the model from which the new reduced dimensions are extracted. Thus, we excluded them from our study.

3.2 Clustering

Clustering is a widely used unsupervised learning approach that identifies structure and groups similar data points based on feature similarity. In our study, we analyse the structure of the obtained embeddings using several clustering algorithms to eventually group semantically related inputs together. In particular, we consider the K-means [49], BIRCH [58], HDBSCAN [11] and affinity propagation [18] algorithms. The first two of them (i.e., K-means and BIRCH) are parametric on the number k of clusters. For a given value of k , the K-means algorithm selects k data points to represent the centres of each cluster (centroids) and populates the clusters by assigning input vectors to their closest centroid. Then, the algorithm computes the new centroids by taking the average of the points within each cluster. It then re-computes the cluster membership of the data points iteratively until convergence. BIRCH belongs to the class of hierarchical methods. It arranges the data in a tree-like structure called a *dendrogram*, where each inner node can be expanded into its sub-clusters. The algorithm iteratively merges and refines the dendrogram nodes until the desired number of clusters is reached. The optimal selection of k is task-specific. Therefore, as part of our approach, we propose our own k selection step, which is presented in the next section.

In contrast, clustering methods such as HDBSCAN [11] and affinity propagation [18] do not require the specification of k as an input parameter, since the optimal value is determined internally as part of the clustering algorithm. HDBSCAN is a density-based clustering algorithm that creates a hierarchy of clusters according to varying density levels, before extracting the most stable clusters from this hierarchy. It can effectively identify clusters of varying shapes (including non-convex shapes) and sizes, while also automatically classifying outliers as noise. Affinity propagation, in its turn, is based on the exchange of messages between data points to identify *exemplars*, data points that best represent clusters. This approach automatically determines the number of clusters based on a similarity matrix and a *preference parameter* that decides how likely a point is to become a cluster center.

4 Approach

In this section, we describe TOPOMAP, our approach for creating a topographical input map based on the input space features. The map can then be used to explore different regions of the input space and identify areas of interest, such as those that cause DL failures.

The main steps of TOPOMAP are shown in Fig. 2. Given a set of inputs, split into training set and test set (Step 1: Input Dataset Collection), we compute a numerical embedding from the raw inputs (Step 2: Dimensionality Reduction). This is done for both the training set, which possibly includes a validation set (blue, in Fig. 2), and the test set (purple, in Fig. 2). This step should ideally preserve the relevant features while reducing the input dimensionality, making it suitable for the next step — clustering. More specifically, we take all instances of a given dataset as input, to ensure that the embedding is applied consistently to both the training (including the validation) and the test data. Then, we apply a clustering algorithm that separates the inputs according to their features and assigns each input vector a

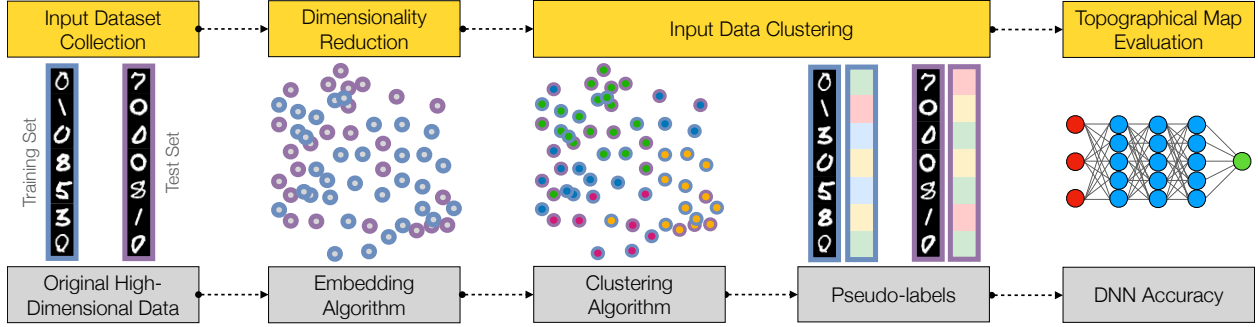


Figure 2: Cluster generation and evaluation pipeline.

pseudo-label representing its cluster (Step 3: Input Data Clustering). The clusters identified by the clustering algorithm outline distinct regions in the topographical map of the input space. After clustering the data, we divide it back into training (including validation) and test sets, according to the original split. In the final step (Step 4: Topographical Map Evaluation), the training (including validation) set is then used to train a DNN classifier that predicts the clusters' pseudo-labels. Finally, we use the trained DNN to automatically assess the quality of the map, by measuring its accuracy on the test set.

It is important to note that the described pipeline contains multiple configurable parameters. The two most important ones are the choice of the embedding algorithm and the choice of the clustering algorithm. In addition, each specific embedding technique and clustering algorithm may have its own hyper-parameters that the user has to select. We refer to each combination of embedding and clustering algorithm, including the specific selection of their hyper-parameter values, as a *clustering configuration*.

The aim of the topographical map quality evaluation step is to select the optimal map for a given input domain. Following the evaluation, the clustering configuration that produces the highest DNN prediction accuracy on the test set is reported to the user and the clusters it produces determine the resulting topographical map of the inputs.

Among the hyper-parameters to be chosen for each clustering algorithm, the number of clusters k , which is required by some, but not all, clustering algorithms, deserves special treatment. In fact, running the entire pipeline depicted in Fig. 2 for all possible values of k is usually not affordable from a computational point of view, due to the typically large range of possible values of k . Hence, only for this hyper-parameter we designed an ad-hoc selection procedure, detailed in the next sub-section. Then, in the following sub-section we provide details about the last step of the pipeline in Fig. 2, automated assessment of the quality of the map associated with each considered clustering configuration (choice of k excluded).

4.1 Automated Selection of the Number of Clusters k

Some clustering algorithms, such as K-means, require the number of clusters, k , to be specified as an additional input. Finding an optimal number of clusters is crucial for our approach as this parameter directly influences the construction of the map and its correctness. However, it would be computationally too expensive to evaluate a large number of values for k for use in our pipeline (see Fig. 2). To address this, we designed an ad hoc selection procedure for this specific hyperparameter only. The main idea is that clustering can be considered as a classification algorithm that operates by majority vote within each cluster, with the categorical labels (or numerical buckets, in case of a regression problem) of the original dataset serving as the ground truth to be predicted by the clusters. Specifically, each cluster is assigned a label determined as the majority label of the train set inputs in the cluster. Such a label is used as the cluster's prediction for all the test set inputs assigned to the such a cluster. The predicted label is then compared to the ground truth label, to determine the accuracy of the cluster's classification. We expect the clustering classifier to become more accurate as k increases, because clusters become smaller and more cohesive. However, at some point, further subdivision of the data into smaller clusters becomes either unnecessary or only marginally beneficial. This happens because the common features that characterise each cluster have already been identified in previous steps (i.e., at lower k), and the new partitions are based on irrelevant and noisy features. Therefore, we can stop increasing k when the plot of accuracy over k flattens, as this indicates that the cluster classification capability has reached the saturation point. When this happens, we expect that all the features relevant for an optimal classification have been discovered in the current set of clusters. When dealing with regression problems, the ground truth labels are represented by continuous numerical values rather than classes. Thus, to be able to apply the proposed automatic selection of k , we split the continuous

Algorithm 1: Selection of number of clusters k

Input : $\mathcal{E}_{train}, \mathcal{E}_{test}$, embedded training and test set
Input : y_{train}, y_{test} , train and test set labels
Input : n , number of classes
Output : k^* , optimal number of clusters

```

1  $\delta_0 \leftarrow 0, \alpha_0 \leftarrow 0, k_0 \leftarrow 0, k^* \leftarrow n;$ 
2 while True do
3    $m_X \leftarrow \text{CLUSTERINGALGORITHM}(\mathcal{E}_{train}, k^*);$ 
4    $m_C \leftarrow \text{ASSIGNCLUSTERLABELSBYMAJORITY}(m_X, \mathcal{E}_{train}, y_{train});$ 
5    $y_{test}^C \leftarrow m_C(\mathcal{E}_{test});$ 
6    $\alpha_1 \leftarrow \text{COMPUTEACCURACY}(y_{test}, y_{test}^C);$ 
7    $\delta_1 \leftarrow \frac{\alpha_1 - \alpha_0}{k^* - k_0};$ 
8   if  $(\delta_1 + \delta_0)/2 < 0.001$  then
9     break;
10  end
11   $\delta_0 \leftarrow \delta_1, \alpha_0 \leftarrow \alpha_1, k_0 \leftarrow k^*, k^* \leftarrow k^* + n;$ 
12 end
13 return  $k^*;$ 
    
```

values associated with the input data into buckets and treat each bucket as a class. We use the standard deviation of the continuous training labels to determine the binning ranges, which define the boundaries of the buckets.

Algorithm 1 presents the pseudo-code of our k selection procedure. We begin (Line 1) with an initial value of k equal to the number of classes. Choosing a smaller value is not appropriate, as our goal is to further partition each class based on the distinguishing features exhibited by its inputs, rather than simply rediscovering the original class labels through clustering. Next, we perform clustering on the training set \mathcal{E}_{train} using the current value for the number of clusters k^* (Line 3). Cluster labels are then assigned by majority vote to the unlabelled clusters of m_X using the train set and its ground truth labels (Line 4). Since the original labels of the test set inputs (y_{test}) are known, we treat the clustering output as a form of classification. We use the labelled clusters of m_C to predict the labels of the test set inputs (Line 5) and compute the corresponding accuracy (Line 6), as well as its derivative (Line 7). If the derivative indicates that accuracy gains have plateaued (Line 8), the loop terminates and the current value of k^* is returned. Specifically, we stop increasing k when the average of the last two derivatives drops below 0.001, suggesting that the accuracy curve has likely flattened. If this condition is not met, k^* is incremented by a multiple of the number of classes (Line 11), and the loop repeats (Line 2).

There exist alternative approaches for the selection of the optimal number of clusters k^* , such as the Silhouette score [45], the Davies-Bouldin index [16], the Calinski–Harabasz index [10] and the elbow rule [29, 47]. However, these methods have well-documented limitations and rely on assumptions that may not hold for our datasets (e.g., convexity of the ground truth clusters). Since we operate in a supervised deep learning context, where the model is tasked with predicting a label, we can leverage the available test set labels to assess the clustering quality for different values of k . Specifically, we evaluate how well each cluster groups vectors that share the same ground truth label. In preliminary experiments, we explored the use of standard clustering quality metrics to guide the selection of k . However, none of the tested metrics consistently yielded reliable choices of k across all datasets. This is likely due to violations of the underlying assumptions these metrics make, which are not satisfied by some of our datasets. In contrast, the procedure described in Algorithm 1 consistently produced meaningful values of k across all datasets, as confirmed by our manual inspection of representative clusters.

4.2 Automated Assessment of the Quality of Clusters

To identify the clustering configuration that produces the most meaningful and discriminative clusters, we evaluate combinations of embedding models, clustering algorithms, and their corresponding hyperparameters. We simulate the perspective of a human user who is asked to manually distinguish clusters based on the shared features of the elements within each group. From this standpoint, we consider clusters to be of high quality if they are internally cohesive and clearly distinguishable from one another. We assume that, similarly to a human evaluator, a deep neural network would be able to differentiate between vectors from different clusters and recognise vectors within the same cluster when the clustering is of high quality. In contrast, such discrimination would be more difficult when the clusters are less coherent. This perspective can be viewed as an automated counterpart of human-subject evaluations, which have been adopted in related work focused on assessing the interpretability and consistency of input features in domains similar to ours [63].

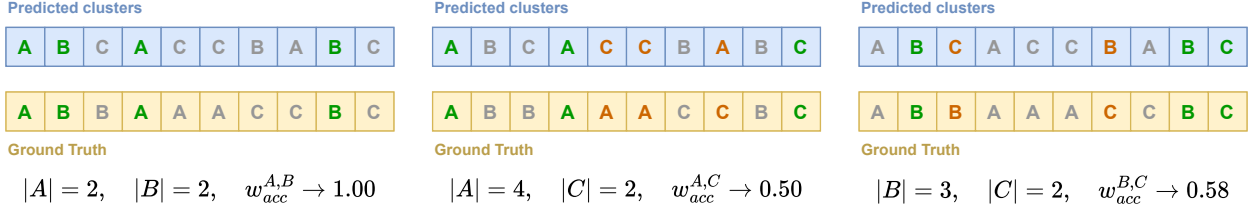


Figure 3: Computation of the *weighted pairwise accuracy* on a three-cluster prediction example.

For each clustering configuration, we apply the clustering model to the full dataset (both train and test set), and get a pseudo-label (i.e., a cluster identifier) for each input. This pseudo-label is assigned to the original, non-embedded data to serve as ground truth label to be predicted by training a deep neural network (DNN) classifier. We use the relabelled training (and possibly validation) set to train the DNN, and then we evaluate its performance on the corresponding relabelled test set. For each input domain, we choose a well-established and documented DNN architecture that proved to be effective in that particular domain, with adaptations applied as necessary to suit our specific classification task.

Since we cannot determine in advance whether the resulting clusters are balanced, and we actually expect them to be imbalanced in most cases, we apply class weighting during training. Specifically, we adopt a weighted training procedure to ensure that the loss function is not disproportionately influenced by the most populated classes. Furthermore, our focus is on evaluating the worst-case performance of the DNN, particularly its ability to distinguish between the pair of clusters that are most difficult to separate. In this context, the overall weighted accuracy may not provide sufficient insight, as it reflects an aggregate measure across all clusters rather than highlighting the most challenging distinctions. To address this, we introduce a novel metric, referred to as *weighted pairwise accuracy*, which effectively projects standard accuracy onto each possible pair of clusters. We then consider the minimum value across all pairs as a representative indicator of the DNN’s worst-case performance.

We evaluated the DNN’s ability to distinguish between each pair of clusters using the *weighted pairwise accuracy metric*. For a pair $\langle A, B \rangle$, we compute its weighted pairwise accuracy as:

$$w_{acc}^{A,B} = \begin{cases} \frac{w \cdot n_{AA} + n_{BB}}{w \cdot n_{AA} + n_{BB} + w \cdot n_{BA} + n_{AB}}, & w = \frac{|B|}{|A|}, \quad \text{if } |A| < |B| \\ \frac{n_{AA} + w \cdot n_{BB}}{n_{AA} + w \cdot n_{BB} + n_{BA} + w \cdot n_{AB}}, & w = \frac{|A|}{|B|}, \quad \text{otherwise.} \end{cases} \quad (1)$$

where n_{AA} (resp. n_{BB}) indicates the total number of instances belonging to cluster A (resp. B) that are correctly classified, while n_{BA} (resp. n_{AB}) indicates the number of elements of cluster A (resp. B) being predicted as instances belonging to cluster B (resp. A). As this computation is conducted pairwise, any further instance of the two clusters that are labelled neither as A or B is discarded when considering the pair $\langle A, B \rangle$. For the purpose of computing $w_{acc}^{A,B}$, the size of cluster A (resp. B) is determined as $|A| = n_{AA} + n_{BA}$ (resp. $|B| = n_{BB} + n_{AB}$). It can be noticed that the definition of weighted accuracy $w_{acc}^{A,B}$ is symmetric w.r.t. the pair $\langle A, B \rangle$, i.e., $w_{acc}^{A,B} = w_{acc}^{B,A}$. Consequently, we only need to compute it for half of the pairs $\langle A, B \rangle$ with $A \neq B$, and the remaining pairs can be obtained by symmetry.

Fig. 3 shows an example of computation of weighted accuracy for all pairs of 3 clusters, namely $\langle A, B \rangle$, $\langle A, C \rangle$, $\langle B, C \rangle$. For the first pair, $\langle A, B \rangle$, 4 entries (in green) are taken into account, i.e., only those that contain either of the two clusters in both the ground truth and the prediction. It should be noticed that we discard entries that contain cluster B in the ground truth, but C in the prediction (i.e., the 3rd entry) or A in the prediction, but C in the ground truth (i.e., the 8th entry). While these entries do not contribute to the weighted accuracy for the pair $\langle A, B \rangle$, they will contribute to the other pairs (specifically, the 3rd entry to $\langle B, C \rangle$, the 8th to $\langle A, C \rangle$). So, they are taken into account later. On the remaining pairs, we measure the DNN weighted accuracy, which is 100% as all predictions are correct (all 4 kept entries are green). In this case, weighted accuracy and normal accuracy coincide, as $|A|$ is the same as $|B|$. When considering the pair $\langle A, C \rangle$, 6 entries are kept (3 correct/green; 3 incorrect/red), while 5 entries are kept for the pair $\langle B, C \rangle$ (3 correct/green; 2 incorrect/red). In both cases, $|A|$ and $|B|$ differ, so we need to adjust for the class imbalance by introducing a weight w , equal to 2 in the first case and to $3/2$ in the second. This weight multiplies the number of ground truth entries of the minority class, to account for the imbalance. Hence, for $\langle A, C \rangle$, $w_{acc}^{A,C} = (2 + 2 \cdot 1) / (2 + 2 \cdot 1 + 2 + 2 \cdot 1) = 1/2$. For $\langle B, C \rangle$, $w_{acc}^{B,C} = (2 + 3/2 \cdot 1) / (2 + 3/2 \cdot 1 + 1 + 3/2 \cdot 1) = 7/12$.

The minimum of such pairwise accuracies, i.e., 0.5 for the pair $\langle A, C \rangle$, is an indicator of the worst case DNN’s ability to identify discriminative features in the clusters. We use such minimum pairwise accuracy as selection criterion for choosing the best clustering configuration.

Algorithm 2: Selection of clustering configuration

Input : $CC = \{\langle E_1, C_1 \rangle, \langle E_2, C_2 \rangle, \dots\}$, candidate clustering configurations
Input : $X_{train}, X_{valid}, X_{test}$, original training, validation and test set
Output : $\langle E^*, C^* \rangle$, optimal clustering configuration

```

1  $Min \leftarrow \emptyset$ ;
2 foreach  $\langle E, C \rangle \in CC$  do
3    $m_C \leftarrow C(E(X_{train}))$ ;
4    $y_{train}^C, y_{valid}^C, y_{test}^C \leftarrow m_C(E(X_{train})), m_C(E(X_{valid})), m_C(E(X_{test}))$ ;
5    $m_{DNN} \leftarrow \text{TRAINDNN}(X_{train}, y_{train}^C, X_{valid}, y_{valid}^C)$ ;
6    $pred_{test}^C \leftarrow m_{DNN}(X_{test})$ ;
7    $min\_acc \leftarrow 1$ ;
8   foreach  $\langle A, B \rangle \in m_C.clusterIds \times m_C.clusterIds, A \neq B$  do
9      $acc \leftarrow \text{COMPUTEPAIRWISEACCURACY}(pred_{test}^C, y_{test}^C, A, B)$ ;
10    if  $acc < min\_acc$  then  $min\_acc \leftarrow acc$ ;
11  end
12   $Min \leftarrow Min \cup \{\langle E, C, min\_acc \rangle\}$ ;
13 end
14 return  $\langle E^*, C^* \rangle$  such that  $\langle E^*, C^*, min\_acc \rangle \in Min$  and  $min\_acc = \max\{Min.min\_acc\}$ ;
```

Algorithm 2 shows the pseudo-code of our procedure for the selection of the most discriminative clustering configuration. The input CC is a set of clustering configurations $\langle E_i, C_i \rangle$, where E_i represents an embedding and C_i a clustering algorithm, both instantiated with a concrete choice of hyper-parameters. Algorithm 2 takes also in input the original (pre-embedding) input vectors, split into training, validation and test set.

We use set Min , initialized at Line 1, to collect the minimum weighted pairwise accuracy of each configuration. For each configuration $\langle E, C \rangle$ (Line 2), we apply the clustering algorithm C to the training set transformed through the embedding E (Line 3). We use the resulting clustering model m_C to assign a pseudo-label to the input vectors in training, validation, and test set (Line 4) and we train a neural network m_{DNN} using training and validation set (Line 5). We predict the labels of the test set vectors at Line 6, and then we iterate over all pairs of clusters to determine the pair with minimum weighted accuracy (loop at Line 8). In particular, at Line 9 we compute the pairwise weighted accuracy for each pair $\langle A, B \rangle$ and at Line 12 we assign the minimum value to the current configuration $\langle E, C \rangle$, by storing the triple $\langle E, C, min_acc \rangle$ into set Min . Finally, we return as optimal configuration $\langle E^*, C^* \rangle$ the one that maximizes the worst case pairwise accuracy (Line 14).

5 Experimental Procedure

In this section we describe the evaluation procedure used for the assessment of the topographical maps produced by TOPOMAP. We also define a set of research questions that form the basis of our empirical assessment.

5.1 Evaluation Based on Manual Human Assessment

The primary objective of the pipeline that generates a topographical map is to identify input space regions that are both recognisable and distinguishable. To evaluate this, we investigate whether the pipeline focusses on features that are interpretable by humans. For example, in the domain of digit classification, such features might include the orientation or boldness of the digits. To this end, we performed a user study to understand whether manual human assessment agrees with the clustering produced by our approach. When assigned the task of separating distinguishable inputs into two groups, human assessors would ideally produce a clustering consistent with that of TOPOMAP. On the other hand, when the input is hard to distinguish, both a human and TOPOMAP are expected to struggle to achieve high separation accuracy.

The idea behind the design of the human study is presented in Fig. 4. Given a topographical map, we identify the pairs of clusters and rank them by $w_{acc}^{A,B}$. To allow a thorough evaluation, we use a selection of inputs from pairs of clusters with different values of $w_{acc}^{A,B}$. In other words, we carefully select diverse pairs of clusters that are easy, moderate, or difficult to distinguish. For each cluster in a pair, we sample a fixed number (15) of train set inputs and present them to human assessors to give them an idea of how clusters were formed. Then, for each pair, we pick 10 test set inputs from each cluster and ask a human assessor to assign each input to one of the two clusters. We then compare the accuracy of the human assessor with that of the DNN used by TOPOMAP. Furthermore, we also compare it with the accuracy of a binary DNN trained only on the sampled cluster pair, which represents an empirical upper bound for accuracy.

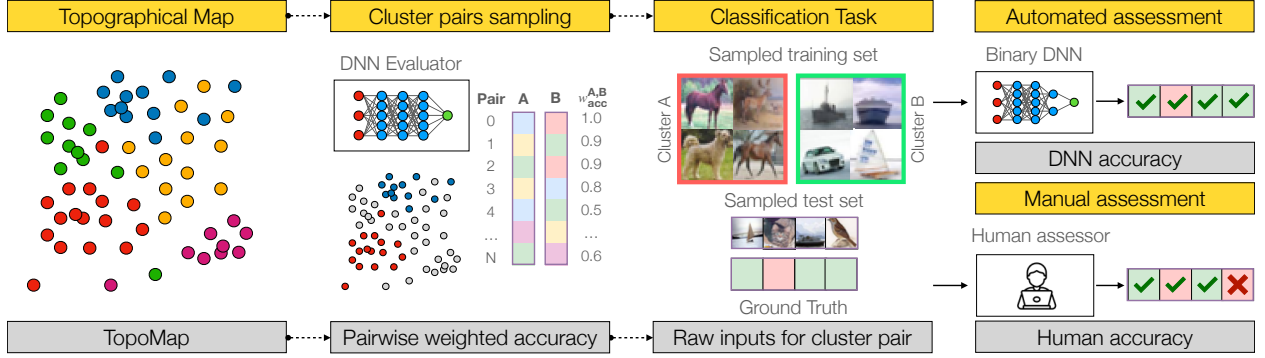


Figure 4: Empirical design of the human evaluation study.

We also ask the human assessor for each pair to indicate how difficult the task was, on a scale from 1 to 5, with the difficulty being explained as the arbitrariness of the choice when assigning an image to a cluster (the more the choices are arbitrary, the more the classification task is difficult). Fig. 4 illustrates the procedure. In the first step, we construct a topographical map of the inputs using the TOPOMAP approach described in the previous section. In Step 2, we consider all possible pairs A, B of clusters and rank them by the TOPOMAP’s weighted accuracy, obtained from the DNN evaluator of TOPOMAP. Among them, we take a subset so as to cover pairs that are easy, moderate and difficult to separate according to the DNN evaluator of TOPOMAP. In practice, we sample these pairs from the initial, middle and final parts of the list ranked by weighted accuracy. In Step 3, for each sampled pair of clusters A and B (resp. red and green), we create a classification task, consisting of examples of inputs selected from the training set and belonging to both clusters, along with a set of inputs from the test set belonging to those clusters and to be classified automatically or manually. The ground truth for the sampled test inputs is the cluster (by construction, A or B) they belong to. Then, in Step 4, both the DNN and the human assessor classify the test inputs into clusters A or B . These classifications are then checked against the ground-truth labels, enabling a comparison of the two accuracies.

The study was structured in the form of a questionnaire consisting of 15 distinct tasks. In addition, we included an attention check to ensure that responses are not given at random. The attention check is not disclosed to the human assessor and consists of a single trivial task, such as separating multiple copies of the same handwritten bold sample of digit ‘2’ from narrow and tilted samples of digit ‘0’. To guarantee a range of diverse and significant responses, a total of four questionnaires have been designed, two based on handwritten digit recognition (MNIST) and the other two on image classification (CIFAR-10). We designed our questionnaires using Qualtrics [2], an online platform commonly used to conduct surveys, evaluations, and other forms of data collection. Participants were recruited through the Amazon Mechanical Turk [1] platform. Amazon Mechanical Turk allows users to filter participants by the level of performance measured on a previously performed task. We chose an approval rate threshold of 95% in accordance with previous studies [42].

5.2 Evaluation Based on Mutation Analysis

Our goal is to determine whether the clusters identified by TOPOMAP contain inputs that can effectively discriminate between faulty (in our experiments, mutant) and correct (in our experiments, original) DNN models, i.e. whether TOPOMAP identifies fault-revealing features. The structure of the evaluation pipeline is shown in Fig. 5.

After computing a topographical map of the inputs with TOPOMAP (Step 1), we assign each of the available test inputs to a TOPOMAP cluster (Step 2). The inputs used for the evaluation consist of the test dataset available in each application domain (e.g., digit recognition). They are considered in raw, pre-embedding form (e.g., pixels of an image), together with their original ground truth labels (e.g., the ground truth digit to be recognised).

In Step 3 (deep learning mutants), inputs are fed to a set of pretrained original DL models and their mutants obtained by applying different mutation operators (MO) and mutation configurations. We need multiple instances of pre-trained original DL models to support statistical mutation killing computation in the next step. To obtain mutants, we use the DeepCrime [24] pre-training DL mutation tool, which implements 24 MOs, each configurable by setting the values of its parameters (e.g., the MO “change learning rate” is configurable by setting a specific mutated learning rate).

In Step 4 (mutation killing analysis), to account for nondeterminism, which affects the training process of the original and mutant models, we compare the predictions performed by the mutants with the ones obtained from the original DL model on N (20 in our studies) experimental runs, each associated with a separate retraining of both the original model

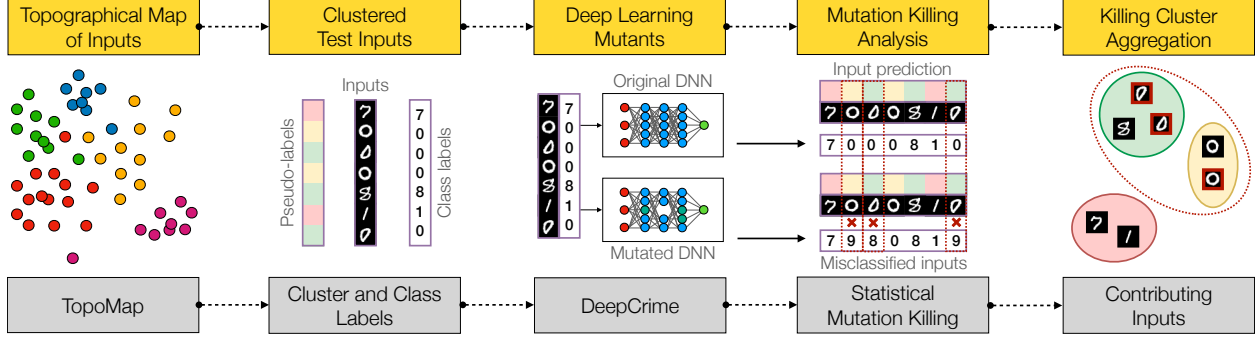


Figure 5: Structure of the mutation-based evaluation pipeline.

and the mutants. Specifically, we implemented the statistical notion of mutation killing proposed by Jahangirova & Tonella [24, 26]: a mutant is *killed* if there is a statistically significant difference between the N performance metric values (i.e., accuracy in the case of classification systems or mean squared error in case of regression problems) obtained for the original model and the N performance metric values of the mutant, with a non-negligible and non-small Cohen’s d effect size.

In the last step, if a mutant is killed, we determine the *contributors*, i.e., the input vectors that potentially contribute to killing the mutant by triggering a misprediction of the mutant. Specifically, a contributor is any input test vector for which the probability of being mispredicted (P^K) by the mutant instances is higher than that of the original baseline DL model instances f_{DL} (see Eq. (2)). The probability of being mispredicted is calculated across N retrainings (instances) of both the original and the mutated model. For regression problems, we introduce an application domain specific threshold τ on the prediction error to decide if an input is a mispredicted, as done in previous work [24]: $\|f(x) - y\| \geq \tau$, where $f(x)$ is the model prediction and y is the ground truth prediction for x . In Fig. 5 (Step 4), the three misclassified inputs indicated with a red cross (2nd, 3rd, and 7th entries) are contributors.

$$\text{Contrib}(M) = \{x \mid x \in X_{\text{test}} \wedge P_{DL}^K(x) < P_M^K(x)\} \quad (2)$$

We sort the clusters produced by our approach in descending order of contributor density and incrementally aggregate them into a candidate killing set, until the obtained aggregation of the clusters results in the mutant being statistically significantly killed [26]. We consider our approach to be effective in killing a mutant if a random selection of inputs with the same size as our killing aggregation has a substantially lower probability of killing the same mutant. In fact, this would indicate that the features characterising each of our clusters are useful to identify a high concentration of killing inputs, which in turn would allow developers to kill a mutant with a small number of instances. We repeated all our experiments 10 times to account for the randomness associated with clustering and random selection of inputs.

5.3 Research Questions

We have performed a set of experiments to answer the following research questions:

RQ1 (Discrimination): *How useful is the proposed approach to select the optimal clustering configuration, i.e., the clustering configuration that better discriminates the test input instances?*

We proposed a DNN-based cluster evaluation pipeline (see Fig. 2) for the selection of the optimal clustering configuration under the assumptions that (1) not all clustering configurations are equivalent; and (2) the optimal clustering configuration might be different across different DL problems (regression vs classification)/domains/subjects.

The goal of this question is to determine whether different clustering configurations are chosen by our evaluation pipeline when applied to different subjects, and whether there is a substantial difference between the optimal choice and the other, discarded options.

Metrics: We use the minimum pairwise weighted accuracy (see Eq. (1)) to characterise the ability of each clustering configuration to produce clusters that a DNN classifier can discriminate.

RQ2 (Explainability): *Are the TOPOMAP regions of the optimal configuration selected by the DNN assessor distinguishable according to human assessors?*

Our intuition behind a DNN selecting the best topographical configuration lies in the idea that it may simulate a human assessing whether two clusters are distinguishable or not. We want to determine whether human assessors are able to

identify patterns that differentiate TOPOMAP’s clusters in a way that correlates with the automated assessment made by TOPOMAP’s DNN assessor.

Metrics: To answer this question, we compute the Pearson correlation coefficient (PCC) and the associated p -value between the human accuracy ($h_{acc}^{A,B}$) and (i) the weighted pairwise accuracy of TOPOMAP’s DNN assessor ($w_{acc}^{A,B}$), (ii) the accuracy of a DNN trained and tested only on the pair of clusters under consideration ($t_{acc}^{A,B}$), (iii) the accuracy of the latter DNN on the same exact 10-input sample analysed by the user ($t'_{acc}^{A,B}$). We compute the same correlation statistic also for the human perceived difficulty ($h_{diff}^{A,B}$).

Human assessment: We conducted a human study (detailed in Section 5.1) to assess whether the DNN classifier can provide a reliable approximation of the human’s ability to recognise and distinguish clusters. We base our evaluation on a sample of cluster pairs with different pairwise-accuracies. In total, this assessment involved 80 different human assessors, 20 for each questionnaire.

RQ3 (Fault Detection): *Are clusters indicative of failure regions in DL systems?*

This question aims to assess whether the clusters obtained from TOPOMAP group together instances that share features making them capable of triggering misbehaviours in the tested DL systems. In particular, we are interested in determining whether there exist small cluster aggregations that can kill each mutant of the original DL system thanks to a high density of killing contributors.

We deem a mutant killing cluster aggregation as *failure-discriminative* if only a small number of clusters is sufficient to create it and if it contains a high density of mutant killing contributors. In fact, the presence of only a few contributors in a cluster may not necessarily result in statistically killing the mutant, as their influence may be diluted by the other non-contributing inputs, which means a larger and denser cluster aggregation should be formed to possibly kill the mutant.

We compare the killing capability of each cluster aggregation with a random selection of a set of inputs having the same size as the aggregation. The cluster aggregation is regarded as a high density failure region if the corresponding random selection has a substantially lower probability of killing each mutant. In fact, this would indicate that developers choosing the test inputs based on our topographical map have a higher chance of exposing failures of the DL system than developers who choose the test inputs randomly.

We perform this experiment on a set of killable mutants, namely, mutants that can be killed by the test data at hand. Additionally, we checked whether our approach can locate a region in the map that successfully kills mutant that are not killable by the whole test set.

Metrics: To answer this question, we introduce the following three metrics computed for each mutant M and then aggregated across MOs. The percentage of input instances gathered within the clusters K in the killing aggregation \mathcal{A} over the size of the entire input test set.

$$\rho_k = \frac{1}{|X_{test}|} \sum_{K \in \mathcal{A}} |K| \quad (3)$$

The proportion of contributors in the aggregation, computed as the ratio between the total number of contributors $c(K)$ in the aggregated clusters and the size of the aggregation \mathcal{A} .

$$\rho_c = \frac{\sum_{K \in \mathcal{A}} c(K)}{|\mathcal{A}|} \quad (4)$$

We check whether a random selection of a number of inputs equal to the size of the killing aggregation \mathcal{A} is capable of killing the mutant. For that purpose, we compute the probability ρ_d that a mutant M is killed across R random input selections of size $|\mathcal{A}|$.

$$\rho_d = \frac{1}{R} \sum_{i=1}^R d_i, \quad \text{with } d_i = \begin{cases} 1, & M \text{ is killed by } |\mathcal{A}| \text{ randomly selected inputs} \\ 0, & \text{otherwise} \end{cases} \quad (5)$$

RQ4 (Killing Propagation): *Do different aggregations, which kill different MOs and mutation configurations, include repeatedly the same clusters?*

The goal of this question is to assess whether the obtained killing clusters persist in triggering misbehaviours across different MOs and MO configurations. This would indicate that the failure regions of our topographical map are useful across different types of faults possibly affecting the DL system under test.

Metrics: We measure the *cluster killing strength* as the percentage of aggregations (ρ_a) that contain a given cluster when considering different mutants. We also measure whether the same aggregations consistently kill at least half or all the configurations of analysed MOs (fault types). The existence of clusters with high killing strength would indicate that their failure exposure capability is not confined to a single, specific type of fault or its configuration.

5.4 Subjects

In order to attain a comprehensive evaluation of our approach, we considered data sets pertaining to different domains, representative of different types of DL applications, spanning from image classification to the prediction of angles through regression. As currently there is no comprehensive and realistic benchmark of DL faults available, we adopt faulty models generated by the DeepCrime mutation tool and available in its replication package [24]. We consider only killable mutants, as their performance can be distinguished from that of the original model by a set of test inputs.

Four of the data sets considered solve classification problems: (i) the *MNIST* (MN) data set [35], consisting of 28×28 greyscale images of ten handwritten digits. The digit classifier is implemented as an eight-layered convolutional neural network (CNN) [14]; (ii) the *CIFAR-10* (CF) data set [33], consisting of 32×32 colour images belonging to ten different classes. The used DL system is a 2D CNN; (iii) the Kaggle *Speaker Recognition* (SR) data set [32], consisting of audio speech samples from five different politicians, labelled with their respective names. The classification task is carried out by deriving a frequency domain representation of the input samples through the Fast Fourier Transform and then training a 1D CNN model [7]; (iv) the *Reuters* (RE) data set from Keras [15]—a preprocessed version of the original data collected by Lewis [36]—consisting of a selection of news articles labelled by topic. Each entry is modelled as a high-dimensional binary vector marking with a 1 the words present in the document (one-hot vocabulary encoding).

The other two data sets pertain to regression problems. In these two cases, only for the automated selection of the optimal number of clusters k^* , we map each numeric output to a set of discrete buckets, determined based on the standard deviation σ_y of the output y across the entire data set. Specifically, for the bucketing of each output y we use the following three ranges: $(-\infty, -\sigma_y)$, $[-\sigma_y, +\sigma_y]$, $(+\sigma_y, +\infty)$. In the case of UE, where the output labels y_1, y_2 belong to \mathcal{R}^2 , we apply these three ranges to each dimension, thus obtaining nine distinct buckets.

The regression problems considered are: (v) the *Udacity* (UD) data set [27], consisting of images of urban roads labelled with the steering angle used by a self-driving car model to learn how to keep the car in the lane. The DL model we employ is the one implemented by NVIDIA [9]; (vi) the *Unity Eyes* (UE) data set [3], consisting of synthetic eye images rendered by the Namesake framework [57], with the target learning model, implemented as a CNN based on the LeNet architecture [3], being an estimator of the 2D eye gaze direction angle. For what concerns the assessment of the correctness of a prediction, required to determine the contributing inputs for the UD data set, we adopt the threshold $\tau = 0.3$, following previous empirical work [24]. For UE we instead set a threshold $\tau = 5^\circ$ which is the step for both yaw and pitch angles used during the generation of the data set.

For each data set, we determine the ability of test set instances to kill mutants of the original high-performance DL model. The MOs used come from previous research [24], to which we refer for further details on the MOs as well as their respective configuration parameters. Of these MOs, we selected the instances that are killable by the entire test set. In particular, we consider the following operators: add activation function to layer (AAL), add noise to training data (TAN), add weights regularisation (RAW), change activation function (ACH), change batch size (HBS), change early-stopping patience (RCP), change number of epochs (HNE), change labels of training data (TCL), change learning rate (HLR), change loss function (LCH), change optimisation function (OCH), change weights initialisation (CWI), remove portion of training data (TRD), disable data batching (HDB), make output classes overlap (TCO), remove activation function (ARM), unbalance train data (TUD).

6 Results

RQ1. Instance Discrimination: *How useful is the proposed pipeline to select the optimal clustering configuration, i.e., the clustering configuration that better discriminates the input instances?*

Table 1 shows the size of the embedded feature space for each considered data scenario. It can be seen that the data dimensionality is significantly reduced: more than 80% of the original size in all cases. Besides the expected outcome of enhancing significant, discriminative features, the reduced size of the feature space has a positive impact on the computational workload of the pipeline generating the topographical map.

Table 2 shows the results of the application of different combinations of embeddings and clustering algorithms. The column highlighted in gray reports the minimum weighted accuracy of the best-performing topographical map, chosen among all considered configurations of embeddings and clustering algorithms (column ‘ $\min w_{acc}^{A,B}$ ’). For each subject

Table 1: Size of the test set and number of components in the feature space in their original raw form compared to the reduced dimensionality in the embedded form for each data scenario. Although both UD and UE data sets concern regression problems, we can group instances into *buckets* for the sake of having a starting point with respect to the number of clusters for the generation of the topographical map.

	MN	CF	SR	RE	UD	UE
Size of the test set	10000	10000	1350	2246	2432	25857
Number of classes (<i>buckets</i>)	10	10	5	46	(3)	(9)
Raw data dimensionality	784	3072	8000	10000	51200	2162
Embedding dimensionality	86 (10.97%)	98 (3.19%)	565 (7.06%)	1910 (19.10%)	274 (0.53%)	35 (1.62%)

Table 2: Best topographical map for each data scenario by embedding/clustering algorithm configuration with respect to the accuracy of the cluster pair most difficult to discriminate by a DNN. Widely used metrics for cluster quality assessment (Silhouette score, DB/CH index) are included for comparison.

Subject	Embedding	Clustering	k^*	Sil. score ^(↑)	DB index ^(↓)	CH index ^(↑)	$\min w_{acc}^{A,B}$	$w_{acc}^{A,B}$
MN	SVD	K-means	90	0.065	2.453	120.077	0.941	0.999
CF	PCA	K-means	40	0.023	2.806	246.947	0.919	0.992
SR	t-SNE	K-means	40	0.376	0.806	2555.596	0.439	0.982
RE	PCA	BIRCH	138	-0.088	3.607	6.164	0.679	0.996
UD	t-SNE	BIRCH	9	0.380	0.820	2422.404	0.994	0.998
UE	LDA	K-means	27	0.120	1.677	2057.630	0.920	0.986

system (column ‘Subject’), the best configuration is reported in the two columns ‘Embedding’ and ‘Clustering’, while column ‘ k^* ’ shows the optimal number of clusters identified for each corresponding configuration.

For completeness, we also include a number of well established clustering quality metrics, such as Silhouette score [45] (‘Sil. score’ column), the Davies-Bouldin index [16] (‘DB index’ column), the Calinski–Harabasz index [10]. Results indicate that none of these commonly employed metrics for cluster evaluation is consistently correlated with the selection of the best clustering configuration performed by the proposed approach (column ‘ $\min w_{acc}^{A,B}$ ’), based on a DNN classifier.

The high values (higher than 90%) of the cumulative average weighted pairwise accuracy (column ‘ $w_{acc}^{A,B}$ ’) suggest that the trained DNN model is always able to successfully recognise the generated clusters. However, we can also see that the lowest pairwise weighted accuracy (column ‘ $\min w_{acc}^{A,B}$ ’) changes significantly across the subjects and configurations (see replication package [4] for full results). In fact, it can be observed that the choice of the best configuration is dependent on the application scenario, with no specific configuration clearly dominating the others across all considered subjects. Moreover, the choice of a specific embedding (e.g., for UD) or clustering algorithm (e.g., for MN) can have a significant impact on the accuracy (see replication package [4] for full results) and, as a result, produce very different topographical maps. In particular, we can see that the combination of a linear embedding with K-means is the best choice in 3 cases out of 6, while for the remaining 3 subjects, the best topographical map is produced either by BIRCH or by a non-linear embedding. The optimal number of clusters also varies from dataset to dataset, ranging from 9 for UD to 138 for the RE dataset.

Answer to RQ1: Our results show that, for each subject, TOPOMAP can generate a map with meaningful and distinguishable clusters, and that the adopted DNN-based evaluation pipeline is important to discriminate maps of different quality, because different subjects require different clustering configurations.

RQ2. Explainability: *Are the TOPOMAP regions of the optimal configuration selected by the DNN assessor distinguishable according to human assessors?*

For each of the four questionnaires, we recruited 20 participants. For each questionnaire, we considered only the responses that passed the attention check, republishing the questionnaire until the expected number of responses was obtained. Hence, we have 20 respondents for each of the two questionnaires on the MN subject and 20 respondents for each of the two questionnaires on the CF subject. The results, aggregated by subject, are reported in Table 3.

Table 3: PCC, mean and standard deviation computed between human-based metrics and DNN-based metrics with respect to classification according to the cluster labels. Highlighted correlation coefficient indicates statistical significance ($p < 0.05$).

Subject	Metric	Mean \pm StDev	$w_{acc}^{A,B}$		$t_{acc}^{A,B}$		$t'_{acc}^{A,B}$	
			PCC	Mean \pm StDev	PCC	Mean \pm StDev	PCC	Mean \pm StDev
MN	$h_{acc}^{A,B}$	0.671 ± 0.198	0.735	0.915 ± 0.080	0.736	0.919 ± 0.076	0.622	0.913 ± 0.099
	$h_{diff}^{A,B}$	0.147 ± 0.036	-0.691		-0.631		-0.585	
CF	$h_{acc}^{A,B}$	0.624 ± 0.130	0.351	0.856 ± 0.147	0.401	0.882 ± 0.107	0.347	0.866 ± 0.149
	$h_{diff}^{A,B}$	0.159 ± 0.047	0.086		0.050		-0.050	

Within the MN data scenario, empirical evidence suggests a strong correlation for both $h_{acc}^{A,B}$ and $h_{diff}^{A,B}$ with all the DNN-based metrics. In particular, the negative correlation between $h_{diff}^{A,B}$ and the other metrics suggests that the difficulty for a human assessor to distinguish a pair of clusters is aligned with the inaccuracies of the DNN model. However, on the CF dataset, results show a less marked relationship between the metrics, with a moderate correlation between $h_{acc}^{A,B}$ and the DNN-based metrics. In the MN scenario, the measured correlation between $w_{acc}^{A,B}$ and $h_{acc}^{A,B}$ achieves a low p -value ($p < 0.05$) thus suggesting statistical significance. The same also occurs with the $h_{diff}^{A,B}$ metric. On the other hand, in the CF scenario, we observe a moderate correlation between the DNN-based metrics and the $h_{acc}^{A,B}$, however, only the correlation with $t_{acc}^{A,B}$ achieves statistical significance, while the p -value of PCC between $h_{acc}^{A,B}$ and $w_{acc}^{A,B}$ is 0.057, and 0.061 with respect to $t'_{acc}^{A,B}$. Conversely, none of the correlations computed with respect to $h_{diff}^{A,B}$ achieves statistical significance.

In addition to the classification tasks, assessors had also the possibility to indicate which strategy they followed for assigning images to a cluster. In the case of the MN questionnaires, the *digit value* was taken into account by 73% of the participants, *digit orientation* by 58%, while *digit boldness* by 65%. On the other hand, for the CF questionnaires, *image subject* was considered by 83% of the participants, while *image color* was deemed relevant by 60%.

Despite meeting our expectations, the human study comes with some limitations that should be taken into account. First of all, as for the CF dataset, images depict real world subjects with a very low resolution, making it hard even for a human to understand what is represented (this observation came out during a pilot of our study). This is backed by the fact that, by considering solely the human-based metrics, we get that $h_{acc}^{A,B}$ and $h_{diff}^{A,B}$ are moderately negatively correlated (-0.634) in a statistical significant way (p -value: 1.7×10^{-4}) in the MN scenario, while for the questionnaires involving CF, the assessment performed by the participants results in a not discernible negative correlation (-0.288) which does not achieve statistical significance (p -value: 0.12). Second, human-based metrics are compared with DNN-inferred metrics even though the two evaluators (human assessor and DNN model) have been supplied different data samples for training. Indeed, for practical reasons, the human assessor has been shown only a random selection of inputs of a cluster to give them an idea of the features represented by such cluster. Instead, the DNN model was provided with the entire training set for that cluster. This might have affected to some degree the human assessors' performance in the tasks related to the CF scenario, whose clusters are generally heterogeneous with respect to the original class label, the underlying features being not necessarily easily recognizable by humans.

Answer to RQ2: Empirical results show correlation between human-measured accuracy and accuracy measured through DNNs, as well as some degree of negative correlation between the human-perceived task difficulty and the DNN accuracies, particularly on the MN dataset. However, the dataset with more complex features (CF) shows that humans and DNN might consider different features when discriminating input clusters.

RQ3. Fault Detection: Are clusters indicative of failure regions in DL systems?

The results of the mutation analysis for MN and CF can be found in Table 4, while Table 5 shows the results obtained for SR and RE, and Table 6 shows the results for UD and UE.

Results reported for the ρ_k metric show that across subjects on average 9.25% of the test data is enough to kill a mutant. This allows us to pinpoint a specific region of the input space that can quickly lead to fault detection. This value goes to as low as 0.1% for the RAW operator and as high as to 99.5% for TAN operator of the RE subject. For the UD dataset, the average ρ_k across all MOs for killable mutants reaches 30.3%, which is the highest across all subjects.

Table 4: Proportion ρ_k of input instances included into the killing cluster aggregation for two image classification data scenarios (MN and CF) averaged over ten empirical runs. The density of contributing inputs ρ_c quantifies the concentration of killing inputs in the considered cluster aggregations. Results are compared with the concentration of contributors ρ_c and the probability ρ_d of the mutant being killed by the aggregation obtained from the feature regions identified by a random selection.

	MO	#	t_c	MN					CF						
				ρ_k	ρ_c	ρ_d	$\rho_c^{(R)}$	$\rho_d^{(R)}$	#	t_c	ρ_k	ρ_c	ρ_d	$\rho_c^{(R)}$	$\rho_d^{(R)}$
Killable	TAN	2	0.014	0.025	0.070	1.000	0.016	0.300	4	0.462	0.027	0.541	1.000	0.465	0.575
	RAW	3	0.013	0.030	0.060	1.000	0.013	0.233	12	0.433	0.053	0.511	1.000	0.427	0.750
	ACH	9	0.019	0.032	0.096	1.000	0.018	0.322	5	0.573	0.019	0.653	1.000	0.576	1.000
	HBS	3	0.015	0.022	0.075	1.000	0.016	0.033	1	0.391	0.033	0.466	1.000	0.375	0.100
	HNE	3	0.023	0.028	0.107	1.000	0.028	0.433	3	0.459	0.208	0.522	1.000	0.453	0.867
	TCL	6	0.043	0.031	0.332	1.000	0.043	0.367	4	0.416	0.035	0.556	1.000	0.416	0.750
	HLR	3	0.098	0.021	0.308	1.000	0.089	0.367	—	—	—	—	—	—	—
	LCH	10	0.028	0.017	0.153	1.000	0.027	0.640	—	—	—	—	—	—	—
	OCH	3	0.028	0.019	0.132	1.000	0.026	0.333	4	0.347	0.095	0.415	1.000	0.346	0.625
	CWI	4	0.453	0.024	0.534	1.000	0.444	0.625	3	0.882	0.021	0.913	1.000	0.880	1.000
	TRD	4	0.114	0.022	0.318	1.000	0.106	0.575	5	0.512	0.025	0.579	1.000	0.507	0.560
	HDB	1	0.015	0.020	0.094	1.000	0.014	0.200	1	0.428	0.022	0.513	1.000	0.432	1.000
	TCO	1	0.122	0.023	0.972	1.000	0.130	1.000	1	0.419	0.037	0.559	1.000	0.414	1.000
	ARM	4	0.260	0.024	0.304	1.000	0.261	0.425	—	—	—	—	—	—	—
	TUD	6	0.037	0.044	0.229	1.000	0.040	0.450	3	0.429	0.021	0.612	1.000	0.431	0.600
Mean				0.025	0.252	1.000	0.085	0.420			0.050	0.570	1.000	0.477	0.736
Non-Killable	TAN	2	0.011	0.081	0.051	1.000	0.011	0.050	2	0.361	0.045	0.435	1.000	0.355	0.050
	RAW	—	—	—	—	—	—	—	1	0.356	0.040	0.439	1.000	0.349	0.100
	HNE	2	0.011	0.176	0.040	0.950	0.011	0.100	1	0.336	1.000	0.343	0.000	0.336	0.000
	TCL	1	0.010	0.022	0.081	1.000	0.010	0.000	1	0.363	0.141	0.415	1.000	0.358	0.200
	HLR	2	0.010	0.317	0.032	0.750	0.011	0.050	2	0.337	0.282	0.393	0.900	0.338	0.150
	TRD	3	0.012	0.317	0.033	0.833	0.012	0.000	3	0.360	0.060	0.434	1.000	0.368	0.167
	TCO	2	0.011	0.093	0.047	1.000	0.012	0.100	1	0.379	0.047	0.465	1.000	0.390	0.200
Mean				0.157	0.048	0.933	0.011	0.043			0.231	0.418	0.837	0.355	0.113

The average proportion of contributors (ρ_c) in the killing aggregation reaches 51.7%. The smallest proportion is observed for UD (22.2%), while the highest is observed for UE (80.9%). Across all the subjects and MOs, in 59% of the subject/MO pairs, the proportion of contributors is 50% or more, in 25% it is at least 75%, and in 11%, the percentage ranges between 90% and 100%. This further supports the claim that the obtained killing aggregations are *failure-discriminative*: TOPOMAP’s killing aggregations are very small compared to their respective parent test set and, at the same time, the percentage of contributors within them is high enough to kill the mutants. Furthermore, across the subjects, the average ρ_c of the aggregations produced with TOPOMAP is always greater than the one measured on the random selection. In particular, in the case of MN, the concentration of contributors is three times bigger in the TOPOMAP aggregations than in the random ones. Similar results are obtained for RE and SR.

Column $\rho_d^{(R)}$ in Table 4, Table 5 and Table 6 reports the killing probability of a set of randomly selected inputs having the same size as the killing aggregation identified by our approach. It is worth mentioning that the probability of killing for our approach is 100%, as by construction we aggregate regions in the input space that are together capable of killing a mutant. The comparison shows that the gain in killing probability is about 58% for the RE and MN subjects, is more than 50% for SR, and is 26% and 12% for CF and UD, respectively. For the UE dataset, the difference in probability is around 1%, which is due to the mutants of this subject being killable by a diverse set of several unrelated inputs, so that a random selection is able to achieve good performance. The average gain across subjects when using our approach against random input selection is 35%.

Furthermore, in Table 4, Table 5 and Table 6 we present the average killing probability results of TOPOMAP vs random selection on the set of mutants not killed by the original test set (see the *Non-Killable* rows in the tables). For UD, neither TOPOMAP nor random selection could kill any additional mutant, confirming the hypothesis that these mutants are likely equivalent (i.e. non-distinguishable from the original program) [24]. For the other subjects, however, the picture is completely different. The killing probability on *Non-Killable* mutants for MN by TOPOMAP reaches 94%, with a gain of 85% over the random approach. For the rest of the subjects, the killing probability ranges from 61% to 80% with the average gain over random selection of 74%.

Table 5: Proportion ρ_k of input instances included into the killing cluster aggregation for two classification data scenarios (SR and RE) averaged over ten empirical runs. The density of contributing inputs ρ_c quantifies the concentration of killing inputs in the considered cluster aggregations. Results are compared with the concentration of contributors ρ_c and the probability ρ_d of the mutant being killed by the aggregation obtained from the feature regions identified by a random selection.

	MO	SR							RE						
		#	t_c	ρ_k	ρ_c	ρ_d	$\rho_c^{(R)}$	$\rho_d^{(R)}$	#	t_c	ρ_k	ρ_c	ρ_d	$\rho_c^{(R)}$	$\rho_d^{(R)}$
Killable	TAN	—	—	—	—	—	—	—	1	0.089	0.995	0.089	1.000	0.089	1.000
	RAW	—	—	—	—	—	—	—	6	0.230	0.001	0.742	1.000	0.204	0.317
	ACH	1	0.005	0.034	0.311	1.000	0.029	0.000	14	0.431	0.005	0.757	1.000	0.438	0.493
	HBS	—	—	—	—	—	—	—	3	0.098	0.013	0.545	1.000	0.085	0.167
	RCP	4	0.015	0.039	0.548	1.000	0.113	0.450	—	—	—	—	—	—	—
	HNE	4	0.024	0.033	0.557	1.000	0.162	0.525	2	0.097	0.125	0.472	1.000	0.081	0.250
	TCL	6	0.027	0.033	0.637	1.000	0.190	0.667	3	0.217	0.023	0.849	1.000	0.219	0.433
	HLR	—	—	—	—	—	—	—	8	0.607	0.002	0.944	1.000	0.609	0.588
	LCH	12	0.135	0.014	1.000	1.000	1.000	1.000	9	0.146	0.002	0.954	1.000	0.180	0.289
	OCH	4	0.011	0.029	0.439	1.000	0.069	0.425	3	0.150	0.015	0.584	1.000	0.189	0.167
	CWI	1	0.006	0.037	0.338	1.000	0.042	0.100	5	0.366	0.188	0.734	1.000	0.319	0.580
	TRD	6	0.025	0.075	0.566	1.000	0.202	0.517	5	0.171	0.004	0.555	1.000	0.175	0.260
	TCO	6	0.029	0.030	0.703	1.000	0.217	0.733	2	0.285	0.034	0.813	1.000	0.296	0.750
	ARM	—	—	—	—	—	—	—	2	0.484	0.002	0.875	1.000	0.400	0.400
	TUD	—	—	—	—	—	—	—	6	0.099	0.003	0.617	1.000	0.109	0.083
Mean				0.036	0.567	1.000	0.225	0.491			0.101	0.681	1.000	0.242	0.413
Non-Killable	TAN	—	—	—	—	—	—	—	2	0.088	0.004	0.619	1.000	0.092	0.100
	AAL	3	0.004	0.092	0.216	0.967	0.035	0.067	—	—	—	—	—	—	—
	ACH	3	0.004	0.058	0.175	1.000	0.025	0.067	3	0.079	0.159	0.259	1.000	0.082	0.133
	HBS	—	—	—	—	—	—	—	1	0.066	1.000	0.066	0.000	0.066	0.000
	RCP	2	0.006	0.324	0.183	0.800	0.042	0.000	—	—	—	—	—	—	—
	TCL	—	—	—	—	—	—	—	4	0.088	0.359	0.243	0.750	0.090	0.000
	LCH	—	—	—	—	—	—	—	1	0.068	1.000	0.068	0.000	0.068	0.000
	OCH	—	—	—	—	—	—	—	3	0.086	0.004	0.535	1.000	0.086	0.233
	CWI	1	0.003	1.000	0.023	0.000	0.022	0.000	9	0.075	0.228	0.426	0.778	0.068	0.056
	TRD	—	—	—	—	—	—	—	1	0.084	0.004	0.750	1.000	0.100	0.000
	TCO	—	—	—	—	—	—	—	4	0.084	0.515	0.228	0.500	0.084	0.000
Mean				0.369	0.150	0.692	0.031	0.034			0.364	0.355	0.670	0.082	0.058

Answer to RQ3: Our results show that TOPOMAP is highly effective in locating input space regions with high density of inputs sharing features that expose faults in DNN models. The average advantage across subjects of TOPOMAP compared to a random sampling is 35% on killable mutants and 74% on non-killable ones, using on average just 9% of the input data.

RQ4. Killing Propagation: Do different aggregations, which kill different MOs and mutation configurations, include repeatedly the same clusters?

Results presented in Table 7 show that the killing capability of the first most selected cluster (row ‘First’) reaches 67% of mutants for MN, 65% for RE, 45% for SR, and drops to 38%, 33% and 32% for UD, CF, and UE, respectively. The highest participation in killing aggregations is observed for the second (35%) and third MSC (31%) for UD, and the lowest for UE (18% and 10%).

Additionally, we analysed whether different configurations of the same MO tend to be killed by the same aggregations. In Table 8, we present the results at both the mutation operator and the subject levels. For each mutation operator in column ‘MO’, column ‘#Subj’ shows the number of subjects to which it was applied, while columns with subject names show the number of configurations this operator has for each subject. Column ‘Mean#Confs’ show the average number of configurations across all subjects for each operator. Column ‘KillHalf’ shows the probability that the killing aggregation selected by TOPOMAP kills at least half of the MO’s configurations, and column ‘KillFull’ shows the probability that the aggregation kills all the configurations of the MO. At the bottom of the table, row ‘#MOs’ shows the number of MOs applied per subject and rows ‘KillHalf’ and ‘KillFull’ show the average probabilities per subject calculated across all applicable operators.

Table 6: Proportion ρ_k of input instances included into the killing cluster aggregation for two regression data scenarios (UD and UE) averaged over ten empirical runs. The density of contributing inputs ρ_c quantifies the concentration of killing inputs in the considered cluster aggregations. Results are compared with the concentration of contributors ρ_c and the probability ρ_d of the mutant being killed by the aggregation obtained from the feature regions identified by a random selection.

	MO	UD							UE						
		#	t_c	ρ_k	ρ_c	ρ_d	$\rho_c^{(R)}$	$\rho_d^{(R)}$	#	t_c	ρ_k	ρ_c	ρ_d	$\rho_c^{(R)}$	$\rho_d^{(R)}$
Killable	TAN	—	—	—	—	—	—	—	3	0.713	0.035	0.801	1.000	0.716	1.000
	AAL	—	—	—	—	—	—	—	7	0.671	0.035	0.911	1.000	0.674	0.943
	RAW	—	—	—	—	—	—	—	11	0.669	0.034	0.795	1.000	0.670	1.000
	ACH	—	—	—	—	—	—	—	5	0.793	0.035	0.917	1.000	0.791	1.000
	HBS	—	—	—	—	—	—	—	2	0.687	0.036	0.879	1.000	0.684	1.000
	HNE	2	0.059	0.531	0.083	1.000	0.059	0.900	4	0.571	0.033	0.699	1.000	0.571	0.900
	TCL	1	0.151	0.234	0.258	1.000	0.154	1.000	4	0.679	0.035	0.785	1.000	0.679	1.000
	HLR	1	0.026	0.567	0.034	1.000	0.027	0.900	3	0.505	0.033	0.648	1.000	0.505	1.000
	LCH	10	0.805	0.121	0.814	1.000	0.805	0.990	9	0.873	0.032	0.934	1.000	0.873	1.000
	OCH	3	0.064	0.234	0.105	1.000	0.064	0.800	2	0.672	0.034	0.838	1.000	0.672	1.000
	CWI	—	—	—	—	—	—	—	7	0.816	0.036	0.897	1.000	0.816	1.000
	TRD	3	0.106	0.187	0.210	1.000	0.107	0.567	3	0.598	0.033	0.684	1.000	0.597	1.000
	TCO	4	0.109	0.318	0.159	1.000	0.109	0.950	1	0.591	0.104	0.728	1.000	0.590	1.000
	TUD	2	0.067	0.234	0.118	1.000	0.066	0.900	—	—	—	—	—	—	—
Mean				0.303	0.222	1.000	0.174	0.876			0.040	0.809	1.000	0.680	0.988
Non-Killable	RAW	—	—	—	—	—	—	—	1	0.448	0.541	0.502	0.600	0.448	0.000
	HNE	5	0.038	1.000	0.041	0.000	0.038	0.000	—	—	—	—	—	—	—
	TCL	3	0.048	1.000	0.053	0.000	0.048	0.000	—	—	—	—	—	—	—
	HLR	2	0.046	1.000	0.051	0.000	0.046	0.000	—	—	—	—	—	—	—
	TRD	2	0.043	1.000	0.047	0.000	0.043	0.000	—	—	—	—	—	—	—
	TCO	2	0.056	1.000	0.070	0.000	0.056	0.000	1	0.208	0.053	0.364	1.000	0.205	0.000
Mean				1.000	0.052	0.000	0.046	0.000			0.297	0.433	0.800	0.327	0.000

Table 7: Proportion ρ_a of aggregations that contain the 1st, 2nd, 3rd most selected cluster (MSC) across the killable mutants per subject.

MSC	MN	CF	SR	RE	UD	UE
First	0.677	0.272	0.452	0.391	0.654	0.572
Second	0.352	0.198	0.318	0.333	0.346	0.267
Third	0.245	0.163	0.180	0.217	0.115	0.051

Results indicate that ‘KillHalf’ probability is higher than 80% for 71% of the operators. The highest probability of 100% is observed for the AAL, RCP, HDB, TCO, ARM operators, and the lowest probability is reached by the RAW operator (55%). Moreover, this probability reaches 100% on at least one subject for 12 operators out of 17 (TAN, ACH, HNE, HBS, HDB, TCL, HLR, LCH, OCH, CWI, TCO, TUD). When it comes to ‘KillFull’ probability, the results range from 0% (AAL, ARM) to 100% (HDB) with an average of 34%.

At the subject level, results show that the probability that TOPOMAP’s input aggregation kills no less than half of the MO configurations (‘KillHalf’ row) is higher than 60% for all subjects. In particular, the median probability is 92.5% with the lowest achieved on the RE subject (63%) and the highest on the UD and SR subjects (100%). When it comes to the probability of killing all the configurations of an MO by the same killing aggregation (‘KillFull’ column), the lowest results are observed for the RE (7%) and UE (22%) subjects. For the remaining subjects, it ranges between 30% (CF) and 75% (UD). These results indicate that our approach can successfully detect the intrinsic characteristics of the faults and identify discriminative killing input regions in the topographical map.

Answer to RQ4: Results confirm our hypothesis that there exist regions in the topographical map that consistently contribute to fault exposure when different faults are injected by mutation analysis.

Table 8: Probability that an input cluster aggregation kills at least half (*KillHalf*) or all (*KillFull*) of the MO configurations across different MO categories and subjects.

MO	#Configurations						#Subj	Mean#Confs	KillHalf	KillFull
	MN	CF	SR	RE	UD	UE				
AAL	—	—	—	—	—	7	1	7.00	100%	0%
TAN	2	4	—	2	—	3	4	2.75	78%	50%
RAW	3	12	—	12	—	11	4	9.50	55%	12%
ACH	9	5	1	28	—	5	5	9.60	72%	36%
HBS	3	1	—	6	—	2	4	3.00	95%	62%
RCP	—	—	4	—	—	—	1	4.00	100%	60%
HNE	3	3	4	4	2	4	6	3.33	98%	28%
TCL	6	4	6	6	1	4	6	4.50	93%	23%
HLR	3	—	—	16	1	3	4	5.75	95%	40%
LCH	10	—	12	18	10	9	5	11.80	80%	28%
OCH	3	4	4	6	3	2	6	3.67	85%	22%
CWI	4	3	1	10	—	7	5	5.00	88%	24%
TRD	4	5	6	10	3	3	6	5.17	72%	3%
HDB	1	1	—	—	—	—	2	1.00	100%	100%
TCO	1	1	6	4	4	1	6	2.83	100%	67%
ARM	4	—	—	4	—	—	2	4.00	100%	0%
TUD	6	3	—	12	2	—	4	5.75	60%	25%
#MOs	15	12	9	14	8	13	—	—	—	—
KillHalf	97%	70%	100%	63%	100%	88%	—	—	—	—
KillFull	39%	30%	42%	7%	75%	22%	—	—	—	—

6.1 Visualisation

To better understand the structure and properties of the generated clusters, we provide a visual representation of the topographical map. Visualisation helps reveal how the clusters relate to each other in the embedding space and highlights patterns that may not be evident from numerical results alone.

Fig. 6 presents a selection of clusters from the topographical map generated by TOPOMAP for the MN dataset, using SVD as the embedding method and K-means with 90 clusters as the clustering algorithm. The clusters are ranked by the ratio of mutants killed to the total number of mutants considered, and only those with a ratio greater than 0.05 are shown. This results in 23 clusters out of 90, which together kill 92% of all available mutants. To improve interpretability, in Fig. 6 we display the medoid of each cluster. The figure illustrates how groups of similar digits are positioned close to each other in the space and how stylistic features of the digits remain consistent across nearby clusters.

To visualise a complete topographical map, we adopted a graph-based representation. Specifically, each map is converted to a fully connected, weighted, undirected graph. In this graph, the vertices represent clusters and the edges encode the connections between clusters. The weight of each edge corresponds to the Euclidean distance between the centroids of the connected clusters. After computing all distances, we construct a minimum spanning tree of the graph and remove all edges with distances greater than the maximum distance in the tree. This pruning step ensures that the visualisation remains clear and interpretable.

For rendering, we use Gephi [8]. Fig. 7 shows graph representations of the maps generated by TOPOMAP for the MN and SR datasets. The size of each vertex reflects the number of inputs falling into a cluster, while the colour intensity indicates the mutation-killing capability of that cluster. Edge thickness and colour intensity correspond to the distances between clusters. The visualisation highlights strong spatial relationships, showing that clusters with higher fault-revealing power tend to be closely connected.

6.2 Threats to Validity

Construct Validity. The choice of evaluation metrics could threaten our conclusions. To mitigate this risk, we did not limit ourselves to the most popular metrics such as the Silhouette score, but investigated additional measures that better capture the discriminative quality of clusters.

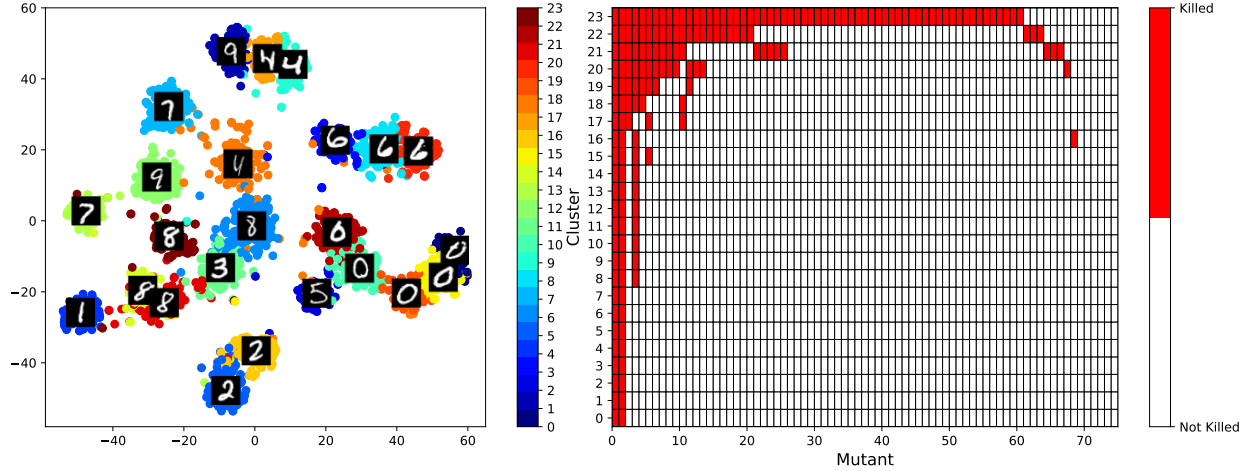
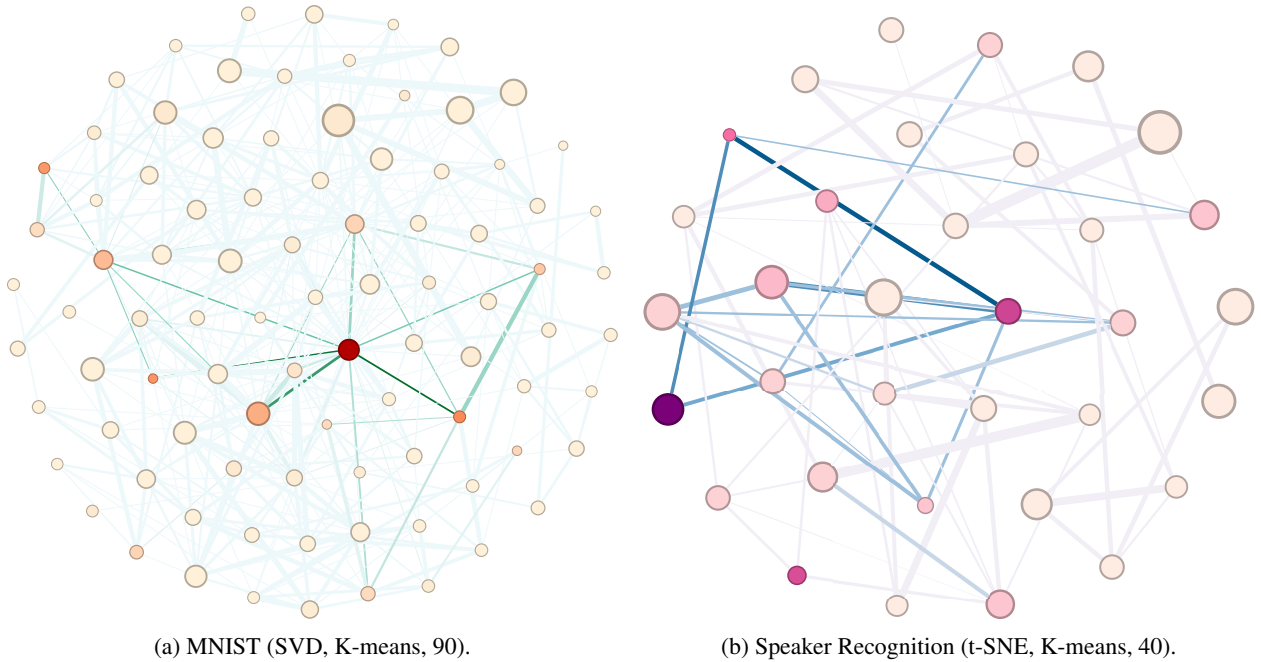


Figure 6: Topographical map computed on the MNIST dataset (SVD, K-means, 90).



(a) MNIST (SVD, K-means, 90).

(b) Speaker Recognition (t-SNE, K-means, 40).

Figure 7: Graph-based visualisation of the topographical maps for MN and SR subjects.

Internal Validity. The main internal threat comes from the selection of dimensionality reduction techniques and clustering algorithms. We experimented with the most frequently used approaches and performed a sophisticated evaluation to choose the best performing configuration for each subject.

External Validity. We have evaluated TOPOMAP on 6 subjects, and therefore there is no guarantee of generalisation. We made a careful choice to represent different tasks (classification and regression) from different application domains.

7 Conclusion

In this work, we propose TOPOMAP, a black-box model-agnostic approach that generates a topographical map of the DNN input feature space. We apply this approach to explore and identify the map regions that share features that are likely to expose faults in a DL model. We also propose an automated technique to configure the dimensionality reduction and clustering methods used in our approach, to optimise the discriminative power of the resulting map regions. We evaluated the effectiveness of TOPOMAP in selecting failure-exposing inputs. The results indicate that in

terms of killing probability, our approach outperforms the random baseline by 35% on average on killable mutants and by 61% on non-killable ones.

Beyond quantitative performance, our analysis shows that the maps produced by TOPOMAP highlight regions that are both meaningful and easily distinguishable. Importantly, our human study indicates that participants’ judgments tend to align with the input grouping of TOPOMAP, particularly when the images in the data set are straightforward to interpret. This supports the interpretability and practical utility of the generated maps.

By focussing their testing effort on map regions that expose misbehaviours of the model under test, developers can efficiently create a challenging test set for each fault possibly affecting their model. In our future work, our aim is to rely on the features exposed by TOPOMAP to generate new input instances tailored to target specific faults when testing DL systems. Furthermore, it may be interesting to investigate the spatial relations between killing clusters, to identify potential macro-regions of contributing inputs.

Data Availability

All data, scripts, and results of this study are publicly available [4].

Ethics approval

The human study was approved by the Ethics Committee of the Università della Svizzera italiana (Decision CE-2025-14).

Funding

This work is funded by the Swiss National Science Foundation (SNSF) programme under the project Toposcope, grant agreement n. 214989.

References

- [1] Amazon Mechanical Turk, 2005. URL <https://www.mturk.com/>.
- [2] Qualtrics, 2019. URL <https://www.qualtrics.com>.
- [3] An implementation of a multimodal CNN for appearance-based gaze estimation, 2020. URL <https://github.com/dlsuroviki/UnityEyesModel>.
- [4] TopoMap: A Feature-based Semantic Discriminator of the Topographical Regions in the Test Input Space, 2025. URL <https://github.com/GianmarcoDeVita/topomap-replication-package>.
- [5] Uzair Adamjee. Cifar10-Autoencoder with CNN. <https://www.kaggle.com/code/uzairadamjee/cifar10-autoencoder-with-cnn>, 2016.
- [6] Zohreh Aghababaeian, Manel Abdellatif, Lionel Briand, Ramesh S, and Mojtaba Bagherzadeh. Black-Box Testing of Deep Neural Networks through Test Case Diversity. *IEEE Trans. Softw. Eng.*, 49(5):3182–3204, may 2023. ISSN 0098-5589. doi:10.1109/TSE.2023.3243522. URL <https://doi.org/10.1109/TSE.2023.3243522>.
- [7] Fadi Badine. Speaker Recognition. https://keras.io/examples/audio/speaker_recognition_using_cnn/, 2020.
- [8] Mathieu Bastian, Sebastien Heymann, and Mathieu Jacomy. Gephi: An open source software for exploring and manipulating networks. In *Proceedings of the 3d International AAAI Conference on Web and Social Media*, 2009.
- [9] Mariusz Bojarski, Davide Del Testa, Daniel Dworakowski, Bernhard Firner, Beat Flepp, Praseen Goyal, Lawrence D. Jackel, Mathew Monfort, Urs Muller, Jiakai Zhang, Xin Zhang, Jake Zhao, and Karol Zieba. End to End Learning for Self-Driving Cars, 2016.
- [10] T. Caliński and J Harabasz. A dendrite method for cluster analysis. *Communications in Statistics*, 3(1):1–27, 1974. doi:10.1080/03610927408827101.
- [11] Ricardo J. G. B. Campello, Davoud Moulavi, and Joerg Sander. Density-Based Clustering Based on Hierarchical Density Estimates. In Jian Pei, Vincent S. Tseng, Longbing Cao, Hiroshi Motoda, and Guandong Xu, editors, *Advances in Knowledge Discovery and Data Mining*, pages 160–172, Berlin, Heidelberg, 2013. Springer Berlin Heidelberg. ISBN 978-3-642-37456-2. doi:10.1007/978-3-642-37456-2_14.

- [12] Nicholas Carlini and David Wagner. Towards Evaluating the Robustness of Neural Networks. In *2017 IEEE Symposium on Security and Privacy (SP)*, pages 39–57, 2017. doi:10.1109/SP.2017.49.
- [13] François Chollet. Building Autoencoders in Keras. <https://blog.keras.io/building-autoencoders-in-keras.html>, 2016.
- [14] François Chollet. Simple MNIST ConvNet. https://keras.io/examples/vision/mnist_convnet/, 2020.
- [15] François Chollet et al. Keras, 2015. URL <https://github.com/fchollet/keras>.
- [16] David L. Davies and Donald W. Bouldin. A Cluster Separation Measure. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, PAMI-1(2):224–227, 1979. doi:10.1109/TPAMI.1979.4766909.
- [17] Yang Feng, Qingkai Shi, Xinyu Gao, Jun Wan, Chunrong Fang, and Zhenyu Chen. DeepGini: prioritizing massive tests to enhance the robustness of deep neural networks. In *Proceedings of the 29th ACM SIGSOFT International Symposium on Software Testing and Analysis, ISSTA 2020*, page 177–188, New York, NY, USA, 2020. Association for Computing Machinery. ISBN 9781450380089. doi:10.1145/3395363.3397357. URL <https://doi.org/10.1145/3395363.3397357>.
- [18] Brendan J. Frey and Delbert Dueck. Clustering by Passing Messages Between Data Points. *Science*, 315(5814): 972–976, 2007. doi:10.1126/science.1136800.
- [19] Yarin Gal and Zoubin Ghahramani. Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning. In Maria Florina Balcan and Kilian Q. Weinberger, editors, *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pages 1050–1059, New York, New York, USA, 20–22 Jun 2016. PMLR. URL <https://proceedings.mlr.press/v48/gal16.html>.
- [20] Fouzi Harrou, Mohamed N. Nounou, Hazem N. Nounou, and Muddu Madakyaru. Statistical fault detection using PCA-based GLR hypothesis testing. *Journal of Loss Prevention in the Process Industries*, 26(1):129–139, 2013. ISSN 0950-4230. doi:<https://doi.org/10.1016/j.jlp.2012.10.003>.
- [21] Dan Hendrycks and Kevin Gimpel. A Baseline for Detecting Misclassified and Out-of-Distribution Examples in Neural Networks. In *Proceedings of International Conference on Learning Representations*, 2017. doi:10.48550/arXiv.1610.02136.
- [22] Harold Hotelling. Relations Between Two Sets of Variates. *Biometrika*, 28(3/4):321–377, 1936. ISSN 00063444. doi:<https://doi.org/10.2307/2333955>.
- [23] 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, New York, NY, USA, 2020. Association for Computing Machinery. ISBN 9781450371216. doi:10.1145/3377811.3380395. URL <https://doi.org/10.1145/3377811.3380395>.
- [24] Nargiz Humbatova, Gunel Jahangirova, and Paolo Tonella. DeepCrime: mutation testing of deep learning systems based on real faults. In *Proceedings of the 30th ACM SIGSOFT International Symposium on Software Testing and Analysis, ISSTA 2021*, page 67–78, New York, NY, USA, 2021. Association for Computing Machinery. ISBN 9781450384599. doi:10.1145/3460319.3464825. URL <https://doi.org/10.1145/3460319.3464825>.
- [25] Md Johirul Islam, Giang Nguyen, Rangeet Pan, and Hridesh Rajan. A Comprehensive Study on Deep Learning Bug Characteristics. In *Proceedings of the 2019 27th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering, ESEC/FSE 2019*, pages 510–520, New York, NY, USA, 2019. ACM. ISBN 978-1-4503-5572-8. doi:10.1145/3338906.3338955. URL <http://doi.acm.org/10.1145/3338906.3338955>.
- [26] Gunel Jahangirova and Paolo Tonella. An Empirical Evaluation of Mutation Operators for Deep Learning Systems. In *2020 IEEE 13th International Conference on Software Testing, Validation and Verification (ICST)*, pages 74–84, 2020. doi:10.1109/ICST46399.2020.00018.
- [27] Gunel Jahangirova, Andrea Stocco, and Paolo Tonella. Quality Metrics and Oracles for Autonomous Vehicles Testing. In *2021 14th IEEE Conference on Software Testing, Verification and Validation (ICST)*, pages 194–204, 2021. doi:10.1109/ICST49551.2021.00030.
- [28] Ian T. Jolliffe. *Choosing a Subset of Principal Components or Variables*, pages 111–149. Springer New York, New York, NY, 2002. ISBN 978-0-387-22440-4. doi:10.1007/0-387-22440-8_6.
- [29] David J. Ketchen and Christopher L. Shook. The application of cluster analysis in strategic management research: An analysis and critique. *Strategic Management Journal*, 17(6):441–458, 1996. doi:[https://doi.org/10.1002/\(SICI\)1097-0266\(199606\)17:6<441::AID-SMJ819>3.0.CO;2-G](https://doi.org/10.1002/(SICI)1097-0266(199606)17:6<441::AID-SMJ819>3.0.CO;2-G).

- [30] Jinhan Kim, Robert Feldt, and Shin Yoo. Guiding deep learning system testing using surprise adequacy. In *Proceedings of the 41st International Conference on Software Engineering, ICSE '19*, page 1039–1049. IEEE Press, 2019. doi:10.1109/ICSE.2019.00108. URL <https://doi.org/10.1109/ICSE.2019.00108>.
- [31] Jinhan Kim, Jeongil Ju, Robert Feldt, and Shin Yoo. Reducing DNN labelling cost using surprise adequacy: an industrial case study for autonomous driving. In *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering, ESEC/FSE 2020*, page 1466–1476, New York, NY, USA, 2020. Association for Computing Machinery. ISBN 9781450370431. doi:10.1145/3368089.3417065. URL <https://doi.org/10.1145/3368089.3417065>.
- [32] Evans Kiplagat. Speaker Recognition Dataset. <https://www.kaggle.com/kongaevans/speaker-recognition-dataset>, 2020.
- [33] Alex Krizhevsky. Learning Multiple Layers of Features from Tiny Images. <https://www.cs.toronto.edu/~kriz/cifar.html>, 2009.
- [34] Alex Kulesza and Ben Taskar. *Determinantal Point Processes for Machine Learning*. Now Publishers Inc., Hanover, MA, USA, 2012. ISBN 1601986289. doi:10.5555/2481023.
- [35] Yann LeCun, Corinna Cortes, and Christopher J.C. Burges. The MNIST Database of Handwritten Digits. <http://yann.lecun.com/exdb/mnist/>, 1998.
- [36] David Lewis. Reuters-21578 Text Categorization Collection. UCI Machine Learning Repository, 1997.
- [37] Lei Ma, Felix Juefei-Xu, Fuyuan Zhang, Jiyuan Sun, Minhui Xue, Bo Li, Chunyang Chen, Ting Su, Li Li, Yang Liu, Jianjun Zhao, and Yadong Wang. DeepGauge: multi-granularity testing criteria for deep learning systems. In *Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering, ASE '18*, page 120–131, New York, NY, USA, 2018. Association for Computing Machinery. ISBN 9781450359375. doi:10.1145/3238147.3238202. URL <https://doi.org/10.1145/3238147.3238202>.
- [38] Eysan Mehrbani and Mohammad Hossein Kahaei. Low-rank isomap algorithm. *IET Signal Processing*, 16(5): 528–545, 2022. doi:<https://doi.org/10.1049/sil2.12124>.
- [39] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient Estimation of Word Representations in Vector Space, 2013.
- [40] Anh Nguyen, Jason Yosinski, and Jeff Clune. Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 427–436, 2015. doi:10.1109/CVPR.2015.7298640.
- [41] Yaniv Ovadia, Emily Fertig, Jie Ren, Zachary Nado, D. Sculley, Sebastian Nowozin, Joshua Dillon, Balaji Lakshminarayanan, and Jasper Snoek. Can you trust your model's uncertainty? evaluating predictive uncertainty under dataset shift. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper_files/paper/2019/file/8558cb408c1d76621371888657d2eb1d-Paper.pdf.
- [42] Eyal Peer, Joachim Vosgerau, and Alessandro Acquisti. Reputation as a sufficient condition for data quality on amazon mechanical turk. *Behavior Research Methods*, 46(4):1023–1031, Dec 2014. ISSN 1554-3528. doi:10.3758/s13428-013-0434-y. URL <https://doi.org/10.3758/s13428-013-0434-y>.
- [43] Kexin Pei, Yinzhi Cao, Junfeng Yang, and Suman Jana. DeepXplore: automated whitebox testing of deep learning systems. *Commun. ACM*, 62(11):137–145, oct 2019. ISSN 0001-0782. doi:10.1145/3361566. URL <https://doi.org/10.1145/3361566>.
- [44] Vincenzo Riccio, Gunel Jahangirova, Andrea Stocco, Nargiz Humbatova, Michael Weiss, and Paolo Tonella. Testing machine learning based systems: a systematic mapping. *Empirical Software Engineering*, 25(6):5193–5254, Nov 2020. ISSN 1573-7616. doi:10.1007/s10664-020-09881-0. URL <https://doi.org/10.1007/s10664-020-09881-0>.
- [45] Peter J. Rousseeuw. Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. *Journal of Computational and Applied Mathematics*, 20:53–65, 1987. ISSN 0377-0427. doi:[https://doi.org/10.1016/0377-0427\(87\)90125-7](https://doi.org/10.1016/0377-0427(87)90125-7).
- [46] Tim Sainburg, Leland McInnes, and Timothy Q. Gentner. Parametric UMAP Embeddings for Representation and Semisupervised Learning. *Neural Computation*, 33(11):2881–2907, 10 2021. ISSN 0899-7667. doi:10.1162/neco_a_01434. URL https://doi.org/10.1162/neco_a_01434.
- [47] Erich Schubert. Stop using the elbow criterion for K-means and how to choose the number of clusters instead. *SIGKDD Explor. Newsl.*, 25(1):36–42, jul 2023. ISSN 1931-0145. doi:10.1145/3606274.3606278. URL <https://doi.org/10.1145/3606274.3606278>.

- [48] Claude E. Shannon. A mathematical theory of communication. *The Bell System Technical Journal*, 27(3):379–423, 1948. doi:10.1002/j.1538-7305.1948.tb01338.x.
- [49] Douglas. Steinley. K-means clustering: A half-century synthesis. *British Journal of Mathematical and Statistical Psychology*, 59(1):1–34, 2006. doi:https://doi.org/10.1348/000711005X48266.
- [50] Andrea Stocco, Michael Weiss, Marco Calzana, and Paolo Tonella. Misbehaviour Prediction for Autonomous Driving Systems. In *Proceedings of 42nd International Conference on Software Engineering, ICSE '20*, page 12 pages. ACM, 2020.
- [51] Joshua B. Tenenbaum, Vin de Silva, and John C. Langford. A Global Geometric Framework for Nonlinear Dimensionality Reduction. *Science*, 290(5500):2319–2323, 2000. doi:10.1126/science.290.5500.2319. URL <https://www.science.org/doi/abs/10.1126/science.290.5500.2319>.
- [52] Laurens van der Maaten and Geoffrey Hinton. Visualizing Data using t-SNE. *Journal of Machine Learning Research*, 9(86):2579–2605, 2008. URL <http://jmlr.org/papers/v9/vandermaaten08a.html>.
- [53] Edoardo Vecchi, Lukáš Pospíšil, Steffen Albrecht, Terence J. O’Kane, and Illia Horenko. eSPA+: Scalable Entropy-Optimal Machine Learning Classification for Small Data Problems. *Neural Computation*, 34(5):1220–1255, 04 2022. ISSN 0899-7667. doi:10.1162/neco_a_01490. URL https://doi.org/10.1162/neco_a_01490.
- [54] Michael E. Wall, Andreas Rechtsteiner, and Luis M. Rocha. *Singular Value Decomposition and Principal Component Analysis*, pages 91–109. Springer US, Boston, MA, 2003. ISBN 978-0-306-47815-4. doi:10.1007/0-306-47815-3_5. URL https://doi.org/10.1007/0-306-47815-3_5.
- [55] Jingyi Wang, Guoliang Dong, Jun Sun, Xinyu Wang, and Peixin Zhang. Adversarial Sample Detection for Deep Neural Network through Model Mutation Testing. In *2019 IEEE/ACM 41st International Conference on Software Engineering (ICSE)*, pages 1245–1256, 2019. doi:10.1109/ICSE.2019.00126.
- [56] Matthew Wicker, Xiaowei Huang, and Marta Kwiatkowska. Feature-Guided Black-Box Safety Testing of Deep Neural Networks. In Dirk Beyer and Marieke Huisman, editors, *Tools and Algorithms for the Construction and Analysis of Systems*, pages 408–426, Cham, 2018. Springer International Publishing. ISBN 978-3-319-89960-2. doi:10.1007/978-3-319-89960-2_22.
- [57] Erroll Wood, Tadas Baltrušaitis, Louis-Philippe Morency, Peter Robinson, and Andreas Bulling. Learning an Appearance-Based Gaze Estimator from One Million Synthesised Images. In *Proceedings of the Ninth Biennial ACM Symposium on Eye Tracking Research & Applications, ETRA '16*, page 131–138, New York, NY, USA, 2016. Association for Computing Machinery. ISBN 9781450341257. doi:10.1145/2857491.2857492. URL <https://doi.org/10.1145/2857491.2857492>.
- [58] Tian Zhang, Raghu Ramakrishnan, and Miron Livny. BIRCH: an efficient data clustering method for very large databases. In *Proceedings of the 1996 ACM SIGMOD International Conference on Management of Data, SIGMOD '96*, page 103–114, New York, NY, USA, 1996. Association for Computing Machinery. ISBN 0897917944. doi:10.1145/233269.233324. URL <https://doi.org/10.1145/233269.233324>.
- [59] Xiyue Zhang, Xiaofei Xie, Lei Ma, Xiaoning Du, Qiang Hu, Yang Liu, Jianjun Zhao, and Meng Sun. Towards characterizing adversarial defects of deep learning software from the lens of uncertainty. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering, ICSE '20*, page 739–751, New York, NY, USA, 2020. Association for Computing Machinery. ISBN 9781450371216. doi:10.1145/3377811.3380368. URL <https://doi.org/10.1145/3377811.3380368>.
- [60] Yuhao Zhang, Yifan Chen, Shing-Chi Cheung, Yingfei Xiong, and Lu Zhang. An Empirical Study on TensorFlow Program Bugs. In *Proceedings of the 27th ACM SIGSOFT International Symposium on Software Testing and Analysis, ISSTA 2018*, pages 129–140, New York, NY, USA, 2018. ACM. ISBN 978-1-4503-5699-2. doi:10.1145/3213846.3213866. URL <http://doi.acm.org/10.1145/3213846.3213866>.
- [61] Shuping Zhao, Bob Zhang, Jian Yang, Jianhang Zhou, and Yong Xu. Linear discriminant analysis. *Nature Reviews Methods Primers*, 4(1):70, Sep 2024. ISSN 2662-8449. doi:10.1038/s43586-024-00346-y. URL <https://doi.org/10.1038/s43586-024-00346-y>.
- [62] 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, ISSTA 2021*, page 79–90, New York, NY, USA, 2021. Association for Computing Machinery. ISBN 9781450384599. doi:10.1145/3460319.3464811. URL <https://doi.org/10.1145/3460319.3464811>.
- [63] Tahereh Zohdinasab, Vincenzo Riccio, Alessio Gambi, and Paolo Tonella. Efficient and Effective Feature Space Exploration for Testing Deep Learning Systems. *ACM Trans. Softw. Eng. Methodol.*, 32(2), mar 2023. ISSN 1049-331X. doi:10.1145/3544792. URL <https://doi.org/10.1145/3544792>.