

# Analyzing the Noise Robustness of Deep Neural Networks

Kelei Cao, Mengchen Liu, Hang Su, Jing Wu, Jun Zhu, Shixia Liu

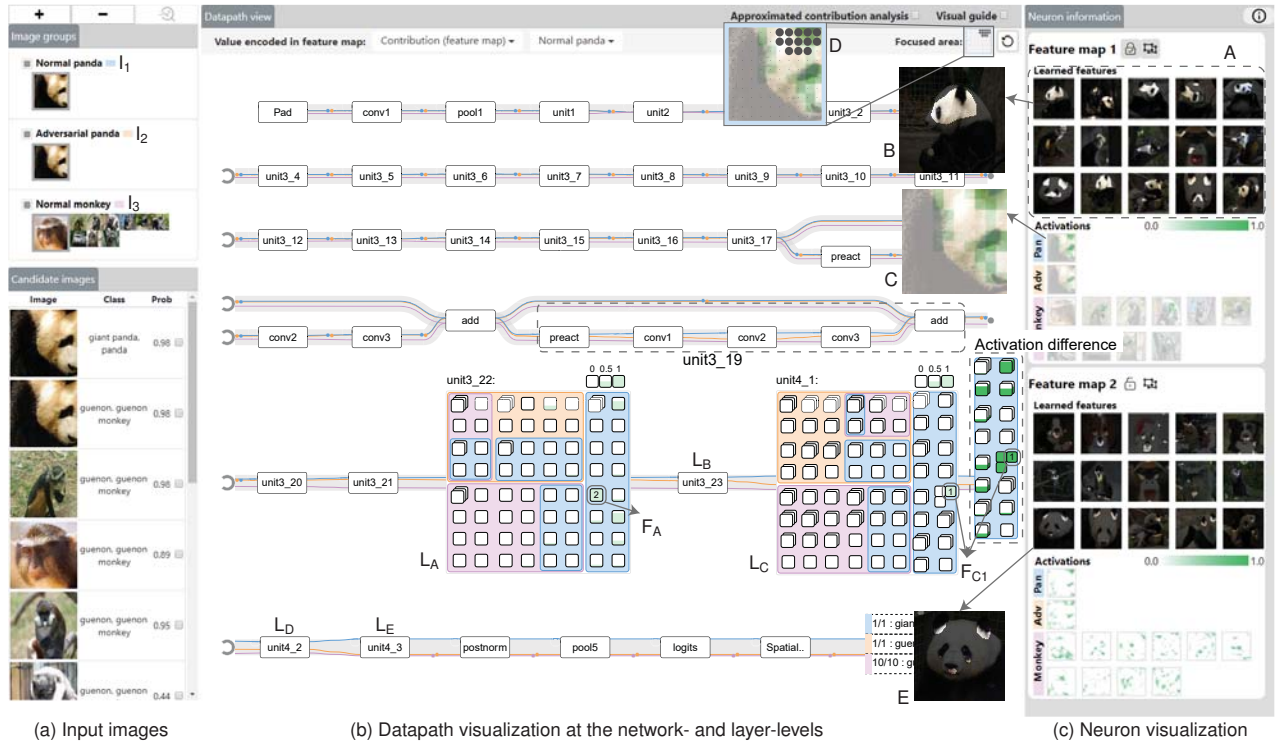


Figure 1: Explanation why an adversarial panda image is not classified as a panda. The root cause is identified as the neurons in the feature map  $F_A$  failing to detect the outline of the panda's ear (E) in the adversarial example, which further leads to the failure of detecting the panda's ear (B) in  $F_{C1}$ .

**Abstract**—Adversarial examples, generated by adding small but intentionally imperceptible perturbations to normal examples, can mislead deep neural networks (DNNs) to make incorrect predictions. Although much work has been done on both adversarial attack and defense, a fine-grained understanding of adversarial examples is still lacking. To address this issue, we present a visual analysis method to explain why adversarial examples are misclassified. The key is to compare and analyze the datapaths of both the adversarial and normal examples. A datapath is a group of critical neurons along with their connections. We formulate the datapath extraction as a subset selection problem and solve it by constructing and training a neural network. A multi-level visualization consisting of a network-level visualization of data flows, a layer-level visualization of feature maps, and a neuron-level visualization of learned features, has been designed to help investigate how datapaths of adversarial and normal examples diverge and merge in the prediction process. A quantitative evaluation and a case study were conducted to demonstrate the promise of our method to explain the misclassification of adversarial examples.

**Index Terms**—Robustness, deep neural networks, adversarial examples, explainable machine learning.

## 1 INTRODUCTION

Deep neural networks (DNNs) have demonstrated superior performance in many artificial intelligence applications, such as pattern recognition and natural language processing [1], [2]. However, researchers have recently found that even a highly accurate DNN

can be vulnerable to carefully-crafted adversarial examples that are intentionally designed to mislead a DNN into making incorrect predictions [3], [4], [5], [6]. For example, an attacker can make imperceptible modifications to a panda image (from  $I_1$  to  $I_2$  in Fig. 1) to mislead a state-of-the-art DNN model [7] to classify it as a monkey. This phenomenon creates high risk when applying DNNs to safety- and security-critical applications, such as driverless cars, face recognition ATMs, and Face ID security on mobile phones [8]. For example, researchers have recently shown that even the state-of-the-art public Face ID system can be fooled by

- K. Cao, H. Su, J. Zhu, and S. Liu, are with Tsinghua University.
- M. Liu is with Microsoft.
- J. Wu is with Cardiff University.

using a carefully-crafted sticker on a hat [9]. Thus, there is an urgent need to understand the prediction process of adversarial examples and identify the root cause of incorrect predictions [8], [10]. Such an understanding is valuable for developing adversarially robust solutions [11], [12], [13]. A recent survey identifies two important questions that require analysis [8]: (1) why similar images (e.g., adversarial and normal panda images) **diverge** into different predictions, and (2) why images from different classes (e.g., adversarial panda images and normal monkey images) **merge** into the same prediction.

To give analytical answers to the questions, we need to solve two technical challenges. The first is to disclose the prediction process of a DNN. To this end, we need to extract the critical neurons and their connections that are responsible for the predictions of examples (Fig. 2 (b)). Such neurons and their connections form the datapaths of examples [10]. However, in a DNN, the neurons have complex interactions with each other [14]. Thus, it is technically demanding to disentangle the roles of these neurons within the entire network and extract the critical neurons to form the datapath. The second challenge is to effectively illustrate and compare the prediction processes of adversarial and normal examples based on the extracted datapaths. A state-of-the-art DNN usually contains hundreds of layers, with millions of neurons in each layer [7]. Thus, an extracted datapath potentially contains millions of neurons and even more connections. Directly visualizing all the neurons and connections in the extracted datapath will result in excessive visual clutter.

To tackle these challenges, we have developed a visual analysis tool, AEVis, to help identify the root cause of misclassification of adversarial examples. Fig. 1 shows an example of using AEVis to analyze why an adversarial panda image is misclassified. On the one hand, we find that the extracted datapaths of the adversarial and normal panda images start to diverge at layer  $L_A$  (Fig. 1) and eventually lead to different predictions. On the other hand, merging starts at layer  $L_C$  (Fig. 1) in the datapaths of the adversarial panda and monkey images. With the use of the developed multi-level visualization, we identify the root cause of this misclassification as both a failed detection of the outline of one of the panda's ears and a faulty detection of a monkey face in the adversarial panda image using the target DNN.

Technically, AEVis aims to disclose the prediction process of a DNN by extracting and visualizing the datapaths for adversarial and normal examples, especially focusing on illustrating how these datapaths diverge and merge.

To achieve this aim, we first formulate the datapath extraction as a subset selection problem, which aims to select a minimum set of neurons that can maintain the predictions of a set of examples. As neurons in a DNN sometimes have similar roles, there is randomness in selecting neurons in the datapath extraction process. As a result, the uniqueness of an example's extracted datapath cannot be guaranteed. Moreover, the randomness hinders the detection of the diverging and merging patterns in the extracted datapaths. To reduce the randomness, we introduce the constraint that it is desirable for the datapaths of adversarial and normal examples to share common feature maps (a set of neurons that share the same weights in a DNN). To extract the datapaths for large DNNs, we approximate the subset selection problem as a continuous optimization that can be efficiently solved by constructing and training a neural network [10].

Second, we have developed a multi-level visualization that illustrates how the extracted datapaths diverge and merge in the prediction process. In particular, at the network-level, we have

created a river-based visualization to provide an overview of the diverging and merging patterns of datapaths. At a detected diverging/merging point (layer level), we employ a treemap-based set visualization to illustrate the neuron groups at this layer and their belonging to different datapaths. This helps experts determine the critical neurons that cause the diverging and merging patterns. In addition, we have enhanced the multi-level visualization with a set of rich interactions that enable experts to effectively analyze the cause of diverging/merging of datapaths. For example, we allow experts to interactively analyze the contribution of neurons in one layer to those of another deeper layer in order to disclose the root cause of a diverging/merging pattern in the compared datapaths.

The paper is an extension of our previous work [15], in which datapaths of examples are extracted and illustrated. In this paper, we address the problem of merging patterns that were not detected in our previous method. We provide a better overview of diverging and merging between the datapaths of examples. In addition, the root cause of such patterns is analyzed more deeply with our refined analysis workflow. To evaluate the usefulness of the new system, we re-invited one expert in our previous work and conducted a deeper analysis of the same two adversarial image pairs, panda-monkey (Sec. 6.2.1) and cannon-racket (Sec. 6.2.2). These improvements come from the following technical contributions:

- **A constrained datapath extraction algorithm** to extract datapaths while preserving their diverging and merging patterns.
- **A river-based visualization** to provide an overview of how datapaths diverge and merge at the network level and **a refined layer-level visualization** to reveal the feature maps of interest.
- **A contribution analysis method** to iteratively investigate the contribution of neurons between two layers and help experts analyze the root cause of diverging/merging in certain layers.

In this paper, we focus on analyzing adversarial examples generated for convolutional neural networks (CNNs), because CNNs are among the most widely-used networks, and most of the current adversarial example generation methods focus on attacking CNNs [8]. Our method can also be used to analyze adversarial examples for other deep networks that use CNNs as the key components.

## 2 RELATED WORK

In the field of visual analytics, a number of methods have been developed to illustrate the working mechanism of a variety of DNNs, such as CNN [16], [17], [18], RNN [19], [20], [21], [22], deep generative models [23], [24], [25], and deep reinforcement learning models [26]. Hohman et al. [27] presented a comprehensive survey to summarize the state-of-the-art visual analysis methods for explainable deep learning. Existing methods can be categorized into three classes: network-centric [28], [29], [30], instance-centric [18], [31], [32], [33], and hybrid [34], [35]. **Network-centric methods.** Network-centric methods help explore the entire network structure of a DNN, illustrating the roles of neurons/neuron connections/layers in the training/test process. In the pioneering work, Tzeng et al. [29] employed a DAG visualization to illustrate the neurons and their connections. This method can illustrate the structure of a small neural network but suffers from severe visual clutter when visualizing state-of-the-art DNNs. To solve this problem, Liu et al. [28] developed a scalable visual analysis tool, CNNVis, based on clustering techniques. It helps explore the roles of neurons in a deep CNN and diagnose

failed training processes. Wongsuphasawat et al. [30] developed a tool with a scalable graph visualization to present the dataflow of a DNN. To produce a legible graph visualization, they applied a set of graph transformations that converts the low-level graph of dataflow to the high-level structure of a DNN.

The aforementioned methods help experts better understand the network structure, but they are less capable of explaining the predictions of individual examples.

**Instance-centric methods.** To address the aforementioned issue, researchers made several recent attempts that focus on instances. These attempts aim at analyzing the learning behavior of a DNN revealed by the instances. A widely-used method is feeding a set of instances into a DNN and visualizing the corresponding log data, such as the activation or the final predictions.

For example, Rauber et al. [31] designed a compact visualization to reveal how the internal activation of training examples evolves during a training process. They used t-SNE [36] to project the high-dimensional activation maps of training examples in each snapshot onto a 2D plane. The projected points are connected by 2D trails to provide an overview of the activation during the whole training process. The method successfully demonstrated how different classes of instances are gradually distinguished by the target DNN. In addition to internal activation, the final predictions of instances can also help experts analyze the instance relationships. For example, the tool Blocks [18] utilizes a confusion matrix to visualize the final predictions of a large number of instances. To reduce the visual clutter caused by a large number of instances and classes, researchers enhanced the confusion matrix using techniques such as non-linear color mapping and halo-based visual boosting. The enhanced confusion matrix was able to disclose the confusion pattern among different classes of instances and further indicated the learning behavior of a target CNN.

The above methods can provide an overview of a large number of instances and help experts analyze their relationships. However, the prediction process of individual instances is less considered. Compared with these macro-level methods, our method focuses on the micro-level and targets the prediction processes of a set of instances (usually a few to dozens). The prediction processes of these instances are visualized using a multi-level datapath visualization. Revealing the prediction processes enables experts to analyze the root cause of the misclassification of adversarial examples.

**Hybrid methods.** The hybrid methods combine the advantages of network-centric and instance-centric methods. Like instance-centric methods, the hybrid methods also feed the target instances into the network and extract log data such as activation maps. The extracted log data is often visualized in the context of the network structure, which provides visual hints to select and explore the data of interest, e.g., the activation in a specific layer. Visualizing the log data in the context of network structure also helps experts explore the data flow from the network input to the output.

There are several papers making progress in this direction. For example, Hartley et al. [34] developed an interactive node-link visualization to show the activation in a DNN. Although this method is able to illustrate detailed activation on feature maps, it suffers from severe visual clutter when dealing with large CNNs. To solve this problem, Kahng et al. [35] developed ActiVis to interpret large-scale DNNs and their results. They employed a multiple coordinated visualization to facilitate experts in comparing activation among examples. The above works mainly focus on exploring the prediction process of **normal** examples. Recently, there is an emerging need in safety-critical fields to

analyze **adversarial** examples of DNNs. While machine learning researchers have developed some holistic views on understanding the existence of adversarial examples [11], [12], there is still a lack of visualization tools to analyze the details. In response to this need, we developed AEVis [15] to analyze the root cause of misclassifications produced by malicious **adversarial** examples. In particular, we developed a datapath extraction method to extract critical neurons and their connections in the prediction process. To enable experts to explore the extracted datapaths, we designed a multi-level visualization that presented datapaths from the high-level network structure to the detailed neuron activation.

As an extension of our previous work [15], this paper re-identifies the central analytical task as analyzing the diverging and merging patterns of normal and adversarial examples. Based on this task, we developed a constrained datapath extraction method that better preserves the diverging and merging patterns of normal and adversarial examples. We also enhanced the whole analysis workflow by introducing several useful interactions, such as activation analysis and contribution analysis. These interactions enable the experts to gradually investigate the major reason for this diverging/merging pattern and thus help them analyze the misclassification of adversarial examples.

### 3 THE DESIGN OF AEVIS

The development of AEVis was in collaboration with the machine learning team that won first place in the NIPS 2017 non-targeted adversarial attack and targeted adversarial attack competitions, which aimed at attacking CNNs [37], [38]. Despite the promising results they achieved, the experts found the research process inefficient and inconvenient, especially in terms of the explanation of the model outputs. In their research process, a key step was to explain the misclassification introduced by adversarial examples. Understanding why an error has been made helps the experts identify the model weakness and further design a more effective attack/defense method. The experts thus desire a tool that can assist them in understanding the prediction process of the target CNN.

#### 3.1 Requirement Analysis

We have identified the following high-level requirements based on previous research and discussions with two experts ( $E_1$  and  $E_2$ ) from the winning team of the NIPS 2017 competition.

**R1 - Extracting the datapaths for adversarial and normal examples.** Both experts expressed the need for extracting the datapaths of adversarial examples, which can disclose the prediction process of adversarial examples and thus serves as the basis for analyzing why the adversarial examples were misclassified. In a CNN, different neurons learn to detect different features [39], and play different roles for the prediction of an example.  $E_1$  said that only analyzing the datapath can greatly reduce their effort by allowing them to only focus on critical neurons rather than having to examine all of them. In addition to the datapaths for adversarial examples,  $E_1$  emphasized the need for extracting datapaths for normal examples simultaneously. He commented that as an adversarial example is often generated by slightly perturbing the pixel values of a normal image, there must be similarities between the two extracted datapaths. Considering the similarity during the datapath extraction process will help extract more meaningful datapaths for comparison during the analysis.

**R2 - Comparing the datapaths of adversarial and normal examples.** As mentioned before, an adversarial example is often



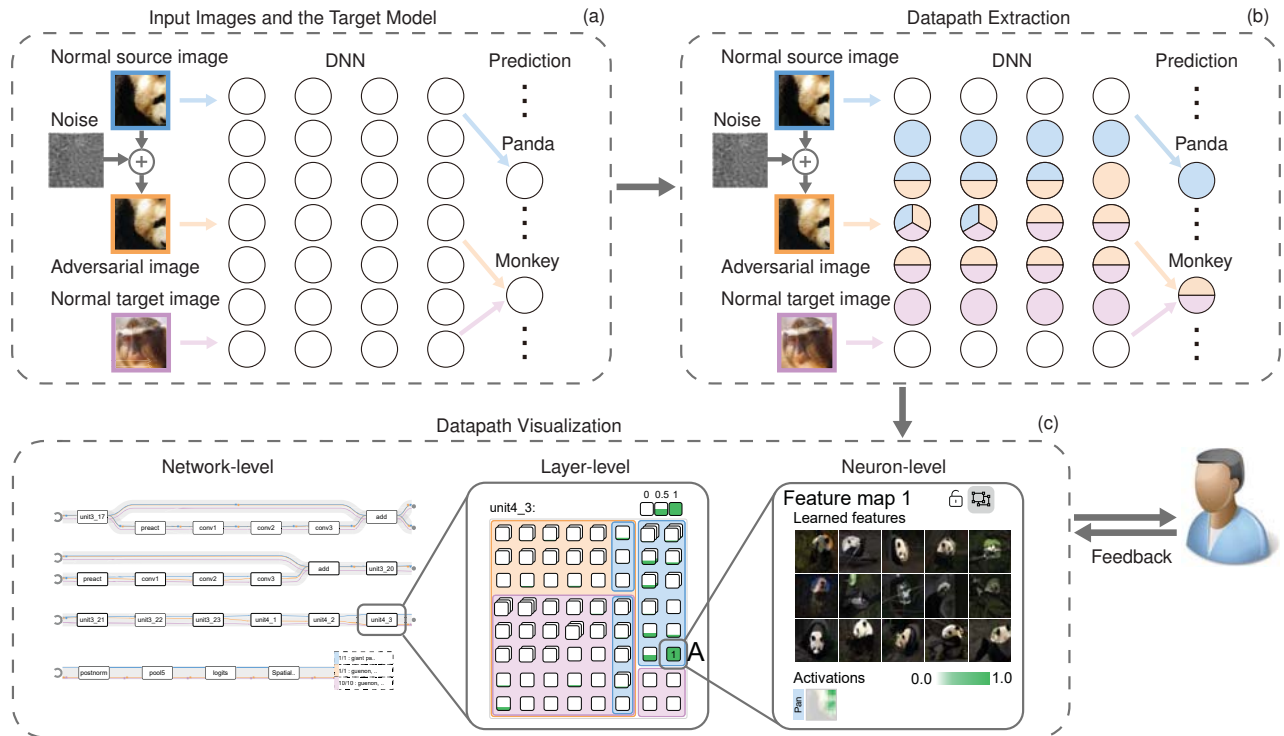


Figure 2: AEVis system overview. (a) Input of the AEVis system; (b) the datapath extraction module; (c) the datapath visualization module that illustrates the extracted datapaths at the network-, layer-, and neuron-level.

generated by adding unperceivable noise to a normal example, and thus there is little difference from the normal image in the input space. However, their prediction results are different. The experts are interested in how they diverge into different predictions. For example,  $E_2$  commented, “I want to know whether there are some critical ‘diverging points’ for the different predictions or they accumulate gradually layer by layer through the network.” To this end,  $E_2$  wanted to compare the datapaths of normal source examples and adversarial examples. Triggered by  $E_2$ ,  $E_1$  added that it was interesting to compare the datapath of an adversarial example (e.g., a panda image that is misclassified as a monkey) with that of normal target examples (e.g., normal monkey images). Such comparisons help understand how these very different images “merge” into the same prediction (e.g., the monkey). The need for visual comparison is consistent with the findings of previous research [40], [41].

**R3 - Exploring datapaths at different levels.** In a large CNN, a datapath often contains millions of neurons and connections. Directly presenting all neurons in a datapath will induce severe visual clutter.  $E_1$  commented, “I cannot examine all the neurons in a datapath because there are too many of them. Instead, I often start by selecting an important layer based on my knowledge and examine the neurons in that layer to analyze the learned features and the activation of these neurons. The problem is that when dealing with a new architecture, I may not know which layer to start with. Thus, I have to examine a bunch of layers, which is very tedious.” He advocated for the idea of providing an overview of the datapath with visual guidance to facilitate experts in selecting the layer of interest. The requirement of providing an overview of a CNN aligns well with previous research [15], [30], [35]. Although the overview of a datapath facilitates experts in finding the layer of interest, it is not enough to diagnose the root cause of the wrong prediction. The experts said that a link between the overview of a datapath and the detailed neuron activation is required, which helps them

identify the most important neurons that lead to misclassification. To summarize, it is desirable to provide a multi-level exploration mechanism that allows experts to zoom into the neurons of interest gradually. Previous research also indicates that visual analytics for deep learning benefits from multi-level visualization [15], [35].

**R4 - Examining how neurons contribute to each other in a datapath.** Finding a diverging or merging point is not the end of the analysis. To develop effective defense methods, we must disclose how such divergence or merging happens. As the data flows from previous layers to the current diverging or merging point, a practical method of finding the root cause is tracing back to the previous layers and examining how the neurons there contribute to the neurons at the diverging or merging point.  $E_1$  commented, “When I find a neuron or feature map that performs very differently for an adversarial and a normal example, I’m interested in the cause of this difference. For example, it is useful to know whether it was caused by the neurons in the previous layer or even the neurons in a far-away layer due to the skip-connections [7] in modern CNNs.” Therefore, we need to analyze how neurons contribute to each other in a DNN. Previous research also indicates that presenting the contributions among neurons is important for understanding the outputs and roles of neurons [15].

### 3.2 System Overview

Driven by the requirements suggested by these experts, we have developed a visual analysis tool, AEVis, to help experts analyze the root cause of the robustness issues arising from adversarial examples. It consists of the following two parts.

- A **datapath extraction module** that extracts the critical neurons and their connections for the predictions of adversarial and normal examples (**R1**).
- A **datapath visualization module** that enables a multi-level (**R3**) visual comparison (**R2**) of the extracted datapaths and

provides rich interactions (**R4**) to analyze the root cause of a misclassification.

As shown in Fig. 2 (a), AEVis takes a trained CNN and the examples to be analyzed as its input. The examples usually include the adversarial examples, normal source examples, and normal target examples. Given the examples and the CNN, the datapath extraction module extracts the critical neurons and their connections that are responsible for the predictions of the examples (Fig. 2 (b)). The extracted datapaths are then fed into the visualization module (Fig. 2 (c)), which supports the navigation and comparison of the datapaths from the high-level layers to the detailed neuron activation.

## 4 DATAPATH EXTRACTION

### 4.1 Basic Problem Formulation

Extracting datapaths of adversarial and normal examples is the basis for analyzing why an adversarial example is misclassified (**R1**). The key challenge is to identify the critical neurons in the prediction process. Once the critical neurons have been identified, selecting the corresponding connections to form the datapath is straightforward. Critical neurons are those that highly contribute to the final prediction. In other words, by only combining the critical neurons and corresponding connections, the prediction of an example will not be changed. Therefore, we aim to select a minimized subset of neurons that can maintain the original prediction. Accordingly, we formulate critical neurons extraction as a subset selection problem:

$$N^{opt} = \arg \min_{N_s \subseteq N} (p(x) - p(x; N_s))^2 + \lambda |N_s|. \quad (1)$$

The first term is to keep the original prediction, and the second term ensures the selection of a minimized subset of neurons. Specifically,  $N$  is the set of neurons in a CNN,  $N_s$  is a subset of  $N$ ,  $N^{opt}$  is the optimized subset consisting of critical neurons,  $p(x)$  is the prediction of example  $x$ , and  $p(x; N_s)$  is the prediction if we only consider the neuron subset  $N_s$ . To measure the difference between two predictions, we adopt the widely used  $\ell_2$  norm.  $|N_s|$  is the size of  $N_s$  and  $\lambda$  is used to balance the two terms. Compared with our previous work, we change the second term from  $|N_s|^2$  to  $|N_s|$ . With this change, we are able to accelerate the entire optimization process by obtaining a minimized subset of neurons more easily according to the Lasso algorithm [42].

The large search space in Eq. (1) hinders a direct solution, which is mainly due to a large number of neurons in a CNN (usually millions). To reduce the search space, we utilize the weight-sharing property in CNNs [1] and group neurons into a set of feature maps. Specifically, in a CNN, neurons in a feature map share the same weights, and thus learn to detect the same feature. Making use of this characteristic, we replace the problem of critical neuron selection with feature map selection and reformulate the problem as:

$$F_{opt} = \arg \min_{F_s \subseteq F} (p(x) - p(x; F_s))^2 + \lambda |F_s|, \quad (2)$$

where  $F$  is the set of feature maps in a CNN, and  $F_s$  is a subset of  $F$ .

### 4.2 Constrained Datapath Extraction

The above method is successful in extracting the critical feature maps for **one** example but sometimes creates difficulty when **comparing** datapaths of adversarial and normal examples, especially in the detection of merging patterns [15]. After discussions with the

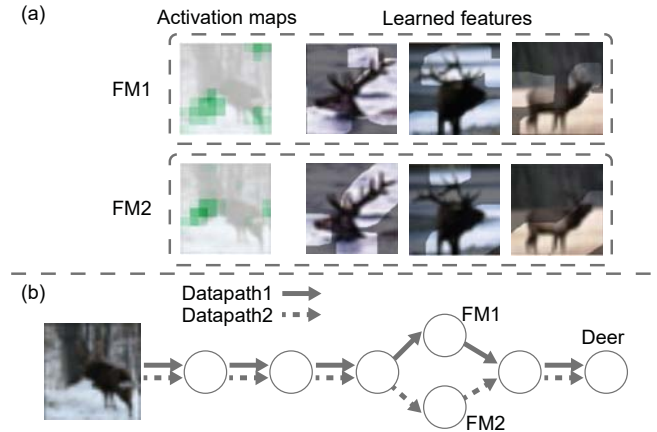


Figure 3: The cause of randomness in datapath extraction. (a) two feature maps detect the same feature (a deer head); (b) there are two equivalent candidate datapaths for the deer image.

domain experts ( $E_1$  and  $E_2$ ) and conducting several experiments, we find that the difficulty is mainly due to the randomness in datapath extraction. Specifically, different feature maps in a CNN may have very similar roles, i.e., detecting nearly the same features [2]. It means that the optimized datapath for an individual adversarial or normal example may not be unique, given the many feature maps with equivalent roles. Thus, extracting a datapath can be treated as sampling one from equivalently good candidate datapaths, which introduces randomness into the datapath extraction. Extracting datapaths that share common feature maps may lead to two feature map subsets that lack common feature maps. As a result, we may over-estimate the difference between two extracted datapaths, which hinders the detection of the diverging and especially the merging patterns.

To illustrate the above analysis, we trained a 6-layer CNN on the CIFAR10 dataset [43]. The network contains 5 convolutional layers and 1 fully connected layer. After training, two equivalently good datapaths for a deer image (the difference between values of Eq. (2) is less than 0.127) were extracted. By examining the feature maps in the two datapaths, we found two feature maps that detected the same feature (a deer head, Fig. 3 (a)) but belonged to different datapaths (Fig. 3 (b)). The above experiment illustrates that feature maps in a CNN may detect the same feature and thus perform similar roles. As a result, there is randomness in the datapath extraction process. This randomness hinders the detection of diverging and merging patterns when comparing datapaths.

To faithfully disclose the diverging and merging patterns, we need to reduce the randomness in datapath extraction. The randomness is mainly caused by the lack of preference among the feature maps that detect the same features. To tackle this issue, we introduce additional constraints into the datapath extraction process, to prioritize the extraction of datapaths that share common feature maps.

Accordingly, the datapaths  $F_{opt}^1, \dots, F_{opt}^i, \dots, F_{opt}^n$  for examples  $x^1, \dots, x^i, \dots, x^n$  are extracted by optimizing:

$$F_{opt}^1, \dots, F_{opt}^i, \dots, F_{opt}^n = \arg \min_{F_s^i \subseteq F} \sum_i L^i + \gamma \sum_{i,j} dis(F_s^i, F_s^j), \quad (3)$$

where the first term  $L^i = (p(x^i) - p(x^i; F_s^i))^2 + \lambda |F_s^i|$  measures how good the datapath is for the  $i$ -th example  $x^i$ . The second term  $dis(F_s^i, F_s^j)$  is the distance to measure the difference

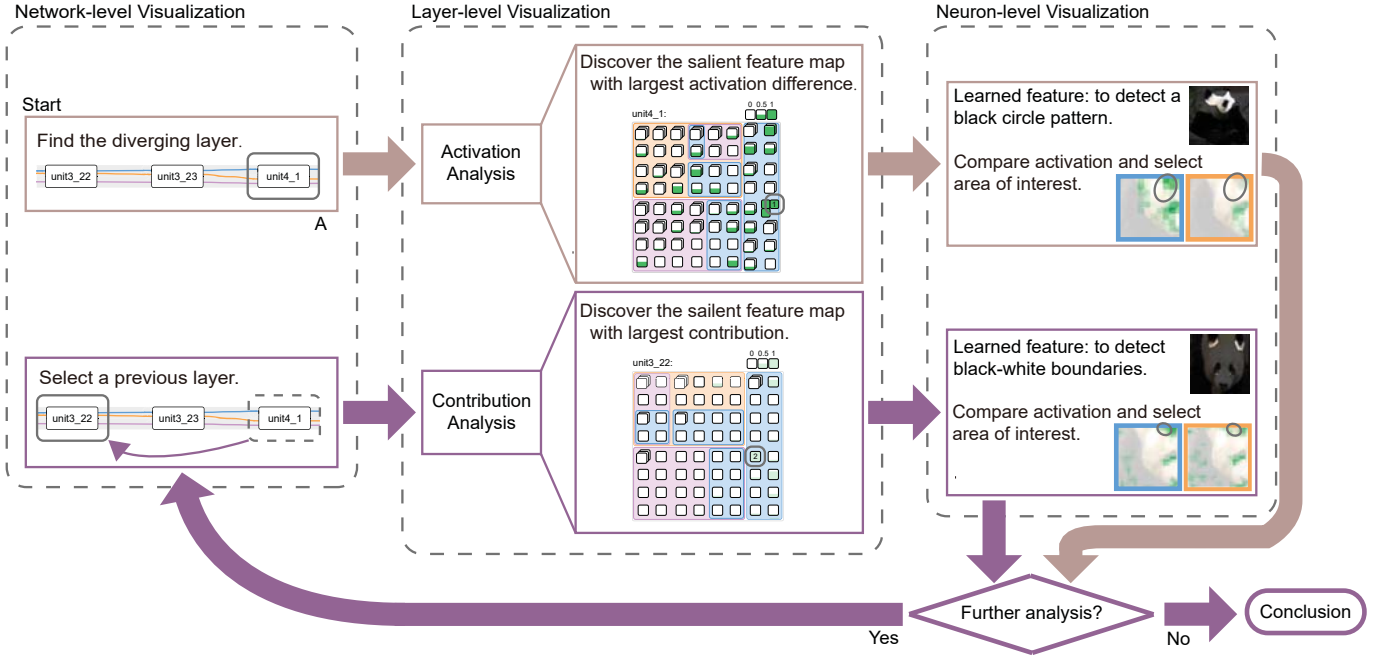


Figure 4: The analysis workflow of the diverging pattern. The brown color represents the first step analysis, and the purple color indicates the subsequent analysis, which is iterative.

between the datapaths for the  $i$ -th and  $j$ -th example, which is defined in Eq. (5). Adding the constraint helps extract datapaths that share common feature maps.  $\gamma$  is used to balance the two terms. Although this method is theoretically sound, in practice, we find that it is difficult to maintain all the predicted labels of the examples, a fundamental requirement for explaining the prediction process. The root cause of the problem is the complexity in jointly finding optimized datapaths for different examples. To solve this problem, we instead approximate the joint optimization into a chain of simpler conditional optimizations. We first obtain the datapath for one example (e.g., the adversarial example to analyze) and iteratively obtain others by treating the previously calculated datapaths as constraints. In particular, for the  $i$ -th example, we solve:

$$F_{opt}^i = \arg \min_{F_s^i \subseteq F} L_i + \gamma \sum_{j=1}^{i-1} dis(F_s^i, F_{opt}^j). \quad (4)$$

To efficiently solve the above subset selection problem, we approximate this NP-hard discrete optimization [44] with a continuous optimization:

$$\begin{aligned} \mathbf{z}_{opt}^i &= \arg \min_{\mathbf{z}^i \in [0,1]^n} L(x^i, \mathbf{z}^i) + \gamma \sum_{j=1}^{i-1} dis(\mathbf{z}_{opt}^j, \mathbf{z}^i), \\ L(x^i, \mathbf{z}^i) &= (p(x^i) - p(x^i; \mathbf{z}^i))^2 + \lambda \|\mathbf{z}^i\|, \\ dis(\mathbf{z}_{opt}^j, \mathbf{z}^i) &= \|\mathbf{z}_{opt}^j - \mathbf{z}^i\|_2, \end{aligned} \quad (5)$$

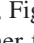
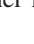
where  $\mathbf{z}^i = [z_1^i, \dots, z_n^i]$  and  $z_k^i \in [0, 1]$  is the contribution of the  $k$ -th feature map in the datapath of the  $i$ -th example  $x^i$ . We apply the commonly-used  $\ell_2$  norm to measure the difference between two datapaths (the second term in Eq. (4)). Eq. (5) is further solved by constructing and training a DNN as in [10]. In particular, we embed the variable  $\mathbf{z}^i$  into the target DNN and train the network on the target adversarial/normal examples by stochastic gradient descent (SGD).

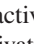
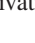
## 5 DATAPATH VISUALIZATION

### 5.1 Overview

An extracted datapath usually contains millions of neurons and even more connections, which prohibits efficient examination of the datapath or discovery of the merging-diverging patterns. To help experts systematically investigate the extracted datapaths, we have designed a multi-level visualization to facilitate the datapath analysis from the high-level network structure to the detailed neuron activation (R3). Accordingly, it consists of three major visualization components at the network-, layer-, and neuron-levels.

**Network-level visualization of data flows.** As shown in Fig. 2 (c), the network-level visualization provides an overview of the extracted datapaths, discloses the potential diverging and merging points, and further guides experts in selecting a layer of interest for examination (R2). Compared with our previous work, we replace the dot-plot-based network-level visualization with a river-based visual metaphor, which has better scalability and is more effective in depicting diverging and merging patterns.

**Layer-level visualization of feature maps.** When an expert identifies a layer of interest (e.g., a diverging or merging point), s/he then zooms in to examine the critical feature maps in that layer (Fig. 2 (c)). For a diverging point, the unique feature maps of each datapath lie in the center of experts' analysis. For a merging point, the shared feature maps between/among datapaths are critical to the analysis. To help experts more quickly find the important and informative feature maps, we use two types of filling styles to encode the activation difference (solid filling , Fig. 2A) and contribution (dotted filling , Fig. 1FA). The higher filling represents a larger value.

**Neuron-level visualization of learned features.** When an expert finds a feature map of interest, AEVis helps him/her understand what features the neurons of interest have learned in the prediction process. Following previous research [2], [18], we employ the learned features of the neurons  (Fig. 1A) and their activation maps  (Fig. 1C) to facilitate the understanding. The activation of



a neuron in a feature map is encoded by the color. Darker green indicates a higher value.

**Analysis workflows.** These three visualizations work together to support a progressive analysis of adversarial examples, which helps experts understand the root cause of the divergence between normal source examples and the corresponding adversarial examples, as well as the merging between the adversarial examples and the normal target examples. Fig. 4 shows the typical workflow for analyzing a diverging pattern. It starts from the network-level visualization where a diverging pattern (Fig. 4A) with several layer groups is identified first. Then with the layer-level visualization and activation analysis, the salient feature map is discovered. Next, by analyzing the learned features and activation in the neuron-level visualization, the expert can identify an area of interest in the focused feature map, which is sent to the contribution analysis module. This module computes the contribution to the activation of the selected neurons from corresponding neurons in previous feature maps. Finally, by examining the contribution of the feature maps in the diverging pattern, the expert gradually investigates the major reason for this divergence. In the merging pattern analysis, instead of using activation analysis, we use contribution analysis as the first step. This is because the merging point is usually followed by the prediction. Contribution analysis helps identify the most important learned feature for the final prediction.

With this exploratory analysis, the potential cause for the wrong predictions is disclosed to facilitate experts in their task of noise robustness analysis. In the below sections, we focus on introducing the network-level visualization, the layer-level visualization, and contribution analysis.

## 5.2 Network-level visualization

In our previous work [15], we employed a dot plot to visualize the difference between **two** datapaths. As shown in Fig. 5, each rectangle represents a layer group, where layers are hierarchically grouped according to the hierarchical computation graph defined in the widely-used TensorFlow Graph Visualization [30]. Each dot in the plot represents the activation similarity between two datapaths of a layer. As a result, the dot plot is combined with the layer group to illustrate the similarities between the extracted datapaths of each layer in each layer group. The position of a dot on the x-axis denotes the similarity value, from 0 (left) to 1 (right). The method has been demonstrated to be useful in detecting the diverging/merging point of two datapaths (e.g., the datapaths of adversarial panda and normal panda images). However, we have received feedback from the experts that the dot-plot-based visualization is less intuitive in revealing the overall evolution pattern of datapath merging and diverging as well as the transition between them. Moreover, it cannot compare three datapaths, which is specifically requested by the experts. The experts said that in analysis, they often needed to examine the adversarial examples in the context of both normal source examples (e.g., panda) and normal target examples (e.g., monkey) to identify the critical diverging/merging points.

To tackle these issues, we have developed a river-based visual metaphor [45], which is inspired by the natural phenomenon of a river merging and diverging along a riverbed. The river-based visualization has been proven effective at depicting diverging and

merging patterns over time [45]. As shown in Fig. 6, we use a curve to represent a datapath. Considering the complexity of the current system, we do not use the curve width to encode extra information, and thus the width is always the same. The distance between the curves represents the similarity between two datapaths. The smaller the distance, the more similar the two datapaths. When comparing three datapaths (adversarial, normal source, and normal target examples), we employ a rule-based method to highlight diverging and merging patterns. In particular, the datapath of the adversarial example stays in the middle with the other two (source and target) on either side (Fig. 6). The screen distance is proportional to the datapath distance  $d_1$  (source-target). The position of the datapath of the adversarial example is determined by retaining the ratio of  $d_2$  (adversarial - source) and  $d_3$  (adversarial - target). To better reveal how data flows in the network, we embed the river-based visualization into the DAG visualization representing the network structure (Fig. 6). With this combination, the **merging** (Fig. 6 (a)), **diverging** (Fig. 6 (b)), and **transition** between datapaths can be easily recognized by examining the distance changes. (Fig. 6 (c)). For example, in Fig. 6 (a), the distance between the blue curve (normal panda) and the orange curve (adversarial example) increases. This indicates that the critical neurons of these two datapaths are gradually becoming less similar to each other, creating a **diverging pattern**. While in Fig. 6 (b), the distance between the orange curve (adversarial example) and the purple curve (normal monkey) decreases. This indicates that the critical neurons of these two datapaths are gradually becoming more similar to each other, creating a **merging pattern**. Fig. 6 (c) shows a transition process from the diverging between a normal panda and an adversarial panda to the merging of an adversarial panda and a normal monkey. Revealing these patterns helps experts quickly locate the layer of interest for further investigation.

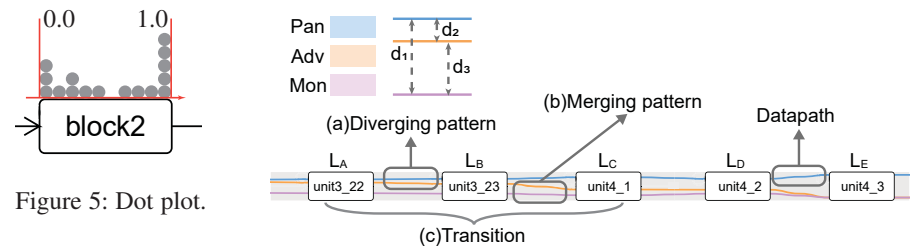


Figure 5: Dot plot.

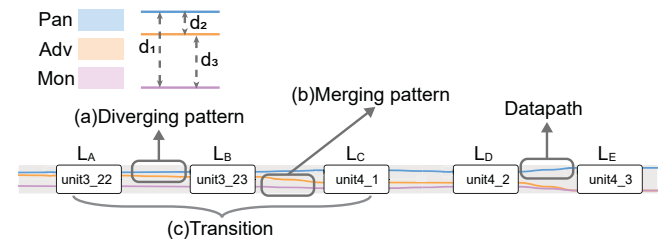


Figure 6: Visualization of three datapaths, with the illustration of (a) the diverging pattern, (b) the merging pattern, and (c) the transition from diverging to merging.

## 5.3 Layer-Level Visualization

When examining a layer of interest, such as a layer with a diverging or merging pattern, the critical feature maps of that layer (Fig. 2 (c)) are important for understanding the key features learned in that layer. These feature maps and corresponding learned features are generally useful for understanding why an adversarial example diverges from its original category and merges into another category. As a result, the unique and shared feature maps between/among datapaths are expected to be encoded and visualized clearly. To this end, we employ a treemap-based visualization to describe the set relationships among the feature maps of different datapaths. The basic idea is illustrated in Fig. 7. In the figure, (a) shows three sets of feature maps belonging to three datapaths (normal panda, adversarial panda, and normal monkey), and their set relations. We

first compute the shared (intersection) and unique parts of the three sets (Fig. 7 (b)). Then a hierarchy is built based on the set inclusion relationships (Fig. 7 (c)). To have more space for displaying the shared and unique parts and distinguishing them clearly, we put the shared parts into the largest set with more feature maps. Finally, the feature map sets and their intersection relationships are visualized with a squarified treemap layout [46] (Fig. 7 (d)).

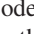
To better reveal the relationships between the shared and unique parts of different datapaths, the treemap cells of the shared parts are placed in a position that is as close as possible to every treemap cell of related feature map sets. For example, the shared part,  $S$ , of feature map sets  $A$ ,  $B$ , and  $C$  are placed near the centers of the three related treemap cells representing  $A$ ,  $B$ , and  $C$  (Fig. 7 (d)). Accordingly, the treemap layout is formulated as an optimization problem with the goal of placing the treemap cells of shared sets close to the center of the cells of related sets:


$$\min \sum_{s_i \in \{0,1\}, \sum_i s_i \geq 2} (f_e(\bigcap_{i:s_i=1} A_i) - f_m(f_e(A_i)))^2, \quad (6)$$

where  $\bigcap_{i:s_i=1} A_i$  is the set that contains all feature maps shared by  $A_i$ , for all  $s_i = 1$ .  $s_i$  is the status variable of  $A_i$ .  $s_i = 1$  indicates that  $A_i$  contains the shared part  $\bigcap_{i:s_i=1} A_i$ , while  $s_i = 0$  means that  $A_i$  does not contain this set.  $f_e(\cdot)$  denotes the center of the treemap cell representing a set, while  $f_m(\cdot)$  is the mean of the centers. Accordingly, the first term represents the center of the shared feature map set  $\bigcap_{i:s_i=1} A_i$ , and the second term represents the mean of the centers of the feature map sets that share  $\bigcap_{i:s_i=1} A_i$ .

This optimized treemap-based visualization can clearly reveal the shared and unique feature maps on the datapaths of interest, which is useful for investigating the roles of different types of feature maps (e.g., unique or shared feature maps) in the prediction. For example, the experts are interested in examining the unique feature maps on each datapath for a diverging point. While for a layer with a merging pattern, the shared feature maps among datapaths are critical for the prediction analysis.

To facilitate the identification of salient feature maps, two types of encoding are employed for activation and contribution, respectively.

**Encoding the activation difference.** We select the maximum neuron activation in a feature map to represent its activation, with the aim of emphasizing the most salient feature detected by the feature map. The solid filling style  is used to encode the activation difference between two datapaths. Taking datapaths  $A$  and  $B$  as an example, their activation difference is  $acti(A) - acti(B)$ . The larger the value, the higher the filling, which indicates that the learned feature is more salient in  $A$  than in  $B$ .

**Encoding the contribution.** A subset of neurons in a specific feature map can be selected as an area of interest (Fig. 1D). Experts can trace how corresponding neurons in previous layers contribute to the activations of the selected neurons. This is useful for identifying the key learned features that lead to diverging or merging of datapaths. In addition, to facilitate the analysis of the diverging/merging patterns for adversarial examples, it is essential to understand how feature maps at each layer contribute to the final prediction. In our visualization, the dotted filling style  is employed to encode the contribution from corresponding neurons in previous layers to the activation of the neurons in the focused area or to the final prediction. The higher filling represents a larger contribution value.

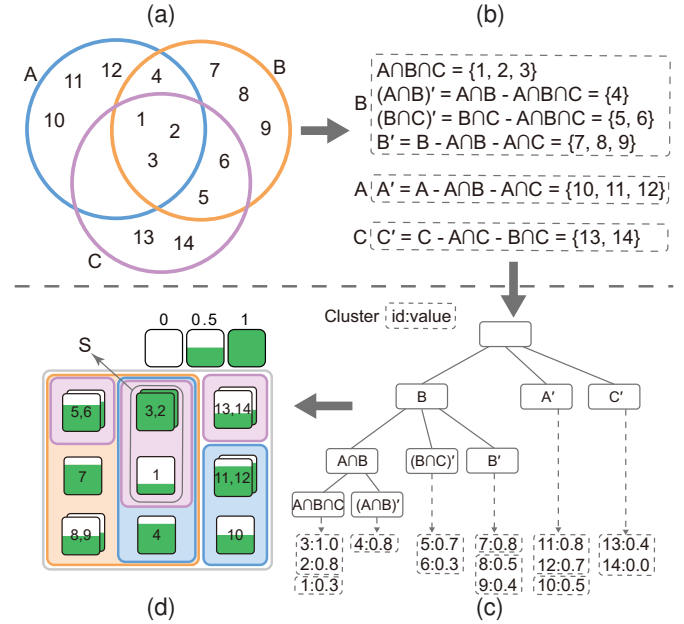


Figure 7: Illustration on how to create a feature map visualization for three datapaths. (a) The three sets of feature maps at a selected layer; (b) The intersection relationships among the sets; (c) A hierarchy based on the set inclusion relationships; (d) The treemap-based visualization of feature maps.

## 5.4 Contribution Analysis

When an expert finds a pattern of interest (e.g., a critical feature map in a diverging point or merging point), s/he often wants to analyze the major cause that leads to the pattern. To this end, we have developed a contribution analysis method to compute the contribution of the previous feature maps to the neuron activation of the feature map of interest (target feature map).

Initially, the contribution analysis is performed based on the whole target feature map. We formulate this problem as a subset selection problem. It aims to select a minimum number of feature maps that can maximally preserve the activation of the target feature map. This formulation is similar to the datapath extraction discussed in Sec. 4. As a result, we also employ the continuous optimization method to select the feature maps and compute the corresponding contribution. In particular, we replace the first term in Eq. (5) with the preservation of the activation of the target feature map:

$$\begin{aligned} z_{opt}^i = \arg \min_{z_{prev}^i \in [0,1]^n} & (f(x^i) - f(x^i; z_{prev}^i))^2 \\ & + \lambda |z_{prev}^i| + \gamma \sum_{j \in [1,m], j \neq i} \|z_{prev}^j - z_{prev}^i\|_2, \end{aligned} \quad (7)$$

where  $f(x^i)$  is the neuron activation of the target feature map on example  $x^i$  and  $f(x^i; z_{prev}^j)$  is the corresponding neuron activation in consideration of the previous feature map contributions  $z_{prev}^j$ .  $m$  is the number of datapaths being analyzed. This optimization problem can be solved similarly with our proposed algorithm in Sec. 4.

When using this contribution analysis method to analyze the adversarial noise, we discover an issue. As shown in Fig. 8 (a), if all the neurons of the target feature map are considered, some irrelevant feature maps, such as FM1 and FM2, are ranked highest, while the relevant one, FM3, is ranked third. This is because, in



addition to the feature that is misled by the adversarial noise, other irrelevant features with high neuron activation, are also considered. These irrelevant features may trigger several irrelevant feature maps and rank them higher.

To tackle this issue, we allow experts to only select neurons (Fig.8A) that are highly activated on the adversarial noise and examine the influence of other feature maps on these selected neurons. For example, when an expert examines a feature map of interest (e.g., Fig. 1C), s/he finds that a certain area of the example is identified as a panda's ear. S/he then checks the previous layers to investigate the reason why this area is activated by the neurons. Therefore, it is useful for the expert to focus on the neurons in this area and check the contribution of the previous feature maps to these selected neurons in the prediction. The key challenge of this problem is to calculate the contribution of the corresponding neurons in each of the previous feature maps. Here the corresponding neurons correspond to the selected neurons in the target feature map. Similarly, we also aim to select a minimum number of feature maps that can maximally preserve the activation of the selected neurons in the target feature map. Accordingly, we change the optimization variable  $z_s$  in Eq. (7) from the whole feature map to the corresponding neurons:

$$z_{opt,p}^i = \arg \min_{z_{prev,p}^i \in [0,1]^n} (f(x^i) - f(x^i; z_{prev,p}^i))^2 + \lambda |z_{prev,p}^i| + \gamma \sum_{j \in [1,m], j \neq i} \|z_{prev,p}^j - z_{prev,p}^i\|_2, \quad (8)$$

where  $z_{prev,p} = [z_{prev,p}^1, \dots, z_{prev,p}^n]$ ,  $z_{prev,p}^k \in [0,1]$  approximates the contribution of the neurons in the focused area of  $k$ -th feature map to the activation of the selected neurons in the target feature map.

The top 3 most contributed feature maps identified by the new method are shown in (Fig.8 (b)), where the most relevant one, FM3, ranks first.

**Time complexity.** The contribution analysis is the only part in AEVis that cannot be pre-computed (the offline pre-computation usually takes a few minutes) because the contribution is calculated based on the selected neurons. It usually takes about 5 seconds to calculate the contribution using SGD to solve Eq. (8). To accelerate the process, we can use the quadratic approximation from our previous work [15], which is faster (computation time  $< 1s$ ) but less accurate. Since our target users (machine learning experts) focus more on analysis accuracy, the SGD-based solution is set as default. Users can switch to the approximated contribution analysis in the interface.

## 6 EVALUATION

We first quantitatively evaluated the effectiveness of the proposed constrained datapath extraction method in comparison with a previous state-of-the-art method, the DGR method [10]. We then demonstrated through a case study how AEVis helped the analysis of the root cause for misclassification of adversarial examples. Expert E<sub>1</sub>, one of the two experts who participated in the evaluation of the previous version of AEVis [15], was invited again to evaluate the usefulness of the new system. The CNN used for evaluation is a pretrained **ResNet-101** [7], which contains 101 layers and is a state-of-the-art CNN for image classification.

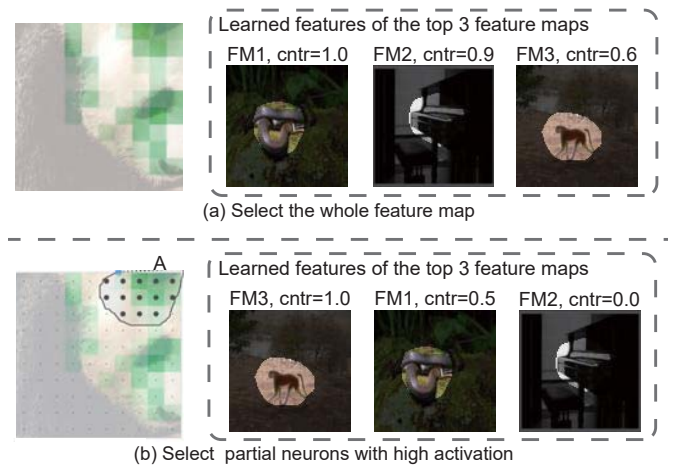


Figure 8: The top 3 most contributed feature maps identified by considering a) the whole feature map of interest; b) part of the feature map of interest.

### 6.1 Quantitative Analysis

As there is no ground-truth for datapaths, the effectiveness of the datapath extraction method is measured by the ability to detect the diverging-merging patterns between the extracted datapaths. Two datasets with different scales were used for this evaluation. One dataset contains all the images of 10 randomly selected classes (shown at the top of Table 1) from ImageNet ILSVRC 2012 [47]. The other contains 100 classes, with 10 randomly selected images in each class. We used a state-of-the-art attacking method, namely, the momentum iterative fast gradient sign method [37], [48], to generate an adversarial example for each image in the datasets. From the classification results of these adversarial images, for each class that was mistakenly classified into, we further sampled 20 target images from the original ImageNet ILSVRC 2012 dataset. Then for each misclassified adversarial image, we constructed 20 triplets of normal source/adversarial/normal target images and extracted datapaths for each triplet.

As shown in Fig. 6, there is a diverging point ( $L_A$ ) where the datapath of the misclassified adversarial image gradually deviates from the datapath of the normal source image, and gets closer to that of the normal target image and merges with it at a point ( $L_E$ ), resulting in the misclassification. Such a diverging followed by a merging pattern (simplified as a **diverging-merging pattern**) is an important characteristic indicating the misclassification of adversarial images. The ability of the extracted datapaths to reflect such patterns is thus used to evaluate the effectiveness of the datapath extraction method.

To determine the occurrence of a diverging-merging pattern, we calculated the difference in the distances between 1) the datapaths of the adversarial and normal source images and 2) the datapaths of the adversarial and normal target images. The difference at layer  $i$  is calculated as:

$$diff(i) = \|z_{opt}^{adv}(i) - z_{opt}^{src}(i)\|_2 - \|z_{opt}^{adv}(i) - z_{opt}^{tar}(i)\|_2 \quad (9)$$

where  $z_{opt}^{adv}$ ,  $z_{opt}^{src}$ , and  $z_{opt}^{tar}$  denote datapaths of the adversarial example, the normal source image that corresponds to the adversarial, and the normal target image, respectively.

If the distance difference of the last  $r$  layers continues to increase towards the end of the model, a diverging-merging

		Dataset 1										Dataset 2
		jeep	schooner	banana	pizza	panda	goldfish	rosehip	snake	tusker	sunglass	100 class (average)
<b>Top-1 Score</b>	DGR	0.000	0.038	0.039	0.059	0.051	0.006	0.000	0.025	0.000	0.029	0.011
	Ours	0.017	0.076	0.058	0.078	0.061	0.029	0.045	0.025	0.021	0.057	0.042
<b>Top-3 Score</b>	DGR	0.033	0.114	0.107	0.176	0.111	0.041	0.061	0.101	0.042	0.095	0.064
	Ours	0.083	0.253	0.165	0.333	0.222	0.110	0.091	0.139	0.104	0.152	0.123
<b>Top-5 Score</b>	DGR	0.067	0.177	0.117	0.275	0.172	0.076	0.091	0.215	0.083	0.162	0.123
	Ours	0.133	0.468	0.204	0.647	0.404	0.174	0.136	0.291	0.167	0.248	0.209

Table 1: The top -1, -3, -5 scores on the 10-class and the 100-class datasets using DGR [10] and our method for datapath extraction.

pattern is detected. Thus, we count the number of layers ( $n_l$ ) that continuously increasing in the last  $r$  layers:

$$n_l = \sum_{i=m-r+1}^m \mathbb{I}(\text{diff}(i) > \text{diff}(i-1)) \quad (10)$$

where  $\mathbb{I}(\cdot)$  is an indicator function. It equals 1 if the predicate is true, and 0 otherwise.  $m$  is the number of layers in the model, and  $r$  is a number recommended by experts to set the minimum length of the diverging-merging pattern ( $r = 8$  in our experiments). Then such a pattern is detected when  $n_l = r$  and  $\text{diff}(m)$  reaches the maximum.

To evaluate the effectiveness of the datapath extraction method, we then defined the following score:

**Top-K score.** For each misclassified adversarial example, we sorted the corresponding target images in descending order according to their datapath similarity (based on the distance defined in Eq. (5),  $1/\text{dis}$ ) with the adversarial image. The top-K score is calculated as the average number of diverging-merging patterns in the top-K target-images of each adversarial example. The higher the score, the better detection of the diverging-merging pattern in the extracted datapaths, and thus more effective the datapath extraction method.

We computed the top-1, top-3, and top-5 scores of our datapath extraction method. For comparison, we also computed the scores for the DGR method [10]. Table 1 shows the computed scores on the 10 randomly selected classes in the first dataset. It can be seen that our method performs better than the DGR method on all the classes. We further computed the scores on the 100 classes in the second dataset, and the average result shown in the last column of Table 1 further verified the effectiveness of our method.

## 6.2 Case Study

We invited expert  $E_1$ , to evaluate the usefulness of AEVis. As  $E_1$  participated in the aforementioned NIPS 2017 adversarial attack competition, he was interested in using the same **DEV** dataset from the competition [38]. He would like to see whether AEVis could help him gain a better understanding of the misclassification of adversarial examples. The **DEV** dataset contains 1000 images of different classes, and for each image, we generated an adversarial image using the non-targeted attacking method developed by the winning team [37], [48].

To facilitate the analysis, we calculated an adversarial score for each adversarial image using the method in [13]. A higher score means a more obvious adversarial example. These scores, together with the classification results, were presented to  $E_1$  for him to select suitable adversarial examples for analysis.  $E_1$  was interested in misclassified adversarial examples with medium scores, as he commented, ‘Less obvious examples often contain subtle changes

with big influence’. After examining these uncertain adversarial examples,  $E_1$  selected two images for further investigation: an image of a panda head that had been misclassified as a guenon monkey ( $l_2$  in Fig. 1), and an image of a cannon misclassified as a racket (Fig. 12).

### 6.2.1 Panda image

To find and understand the root cause of this misclassification,  $E_1$  selected the adversarial panda image ( $l_2$  in Fig. 1), the normal panda image ( $l_1$  in Fig. 1), and 10 normal guenon monkey images ( $l_3$  in Fig. 1) in the AEVis system. The datapaths of normal and adversarial panda images and a representative monkey image were then automatically extracted for further analysis. In particular, the representative monkey image was selected from 10 randomly sampled monkey images, among which it had the highest datapath similarity with the adversarial panda image (Eq. 5).

**Overview.** The system first displayed the datapaths at the network-level (Fig. 1 (b)). The distances among the three datapaths disclose the diverging-merging patterns through the layers. Following the dataflow in the overview,  $E_1$  found that the datapath of the adversarial panda image began to deviate from the datapath of the normal panda image at layer  $L_A$  (Fig. 6), gradually got closer to the datapath of the monkey image, and finally merged into it at  $L_E$  (Fig. 6). From  $L_A$  to  $L_E$  are the layers where the predictive behaviors of neurons were misled by the adversarial noise. To better understand the working mechanism of adversarial noise, he then analyzed the misclassification from two aspects: the **diverging process** of the adversarial image from the normal source image (panda), and the **merging process** of the adversarial image into the normal target images (monkey).

**Diverging analysis.** To analyze which feature maps in the datapath of the adversarial example were critical for the divergence,  $E_1$  first expanded layer  $L_C$  where the divergence became noticeably large. The encoded value of each feature map was set as the activation difference between the normal source and adversarial images. A large difference indicates that the feature map has detected its learned features in the normal source image but not in the adversarial example. With this understanding,  $E_1$  directly checked the feature map  $F_{C1}$ , which had the largest activation difference in  $L_C$ . By examining the learned features (Fig. 1A) of

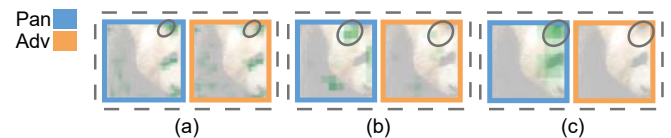


Figure 9: Activation of feature maps in (a)  $L_A$ , (b)  $L_D$  and (c)  $L_E$ .

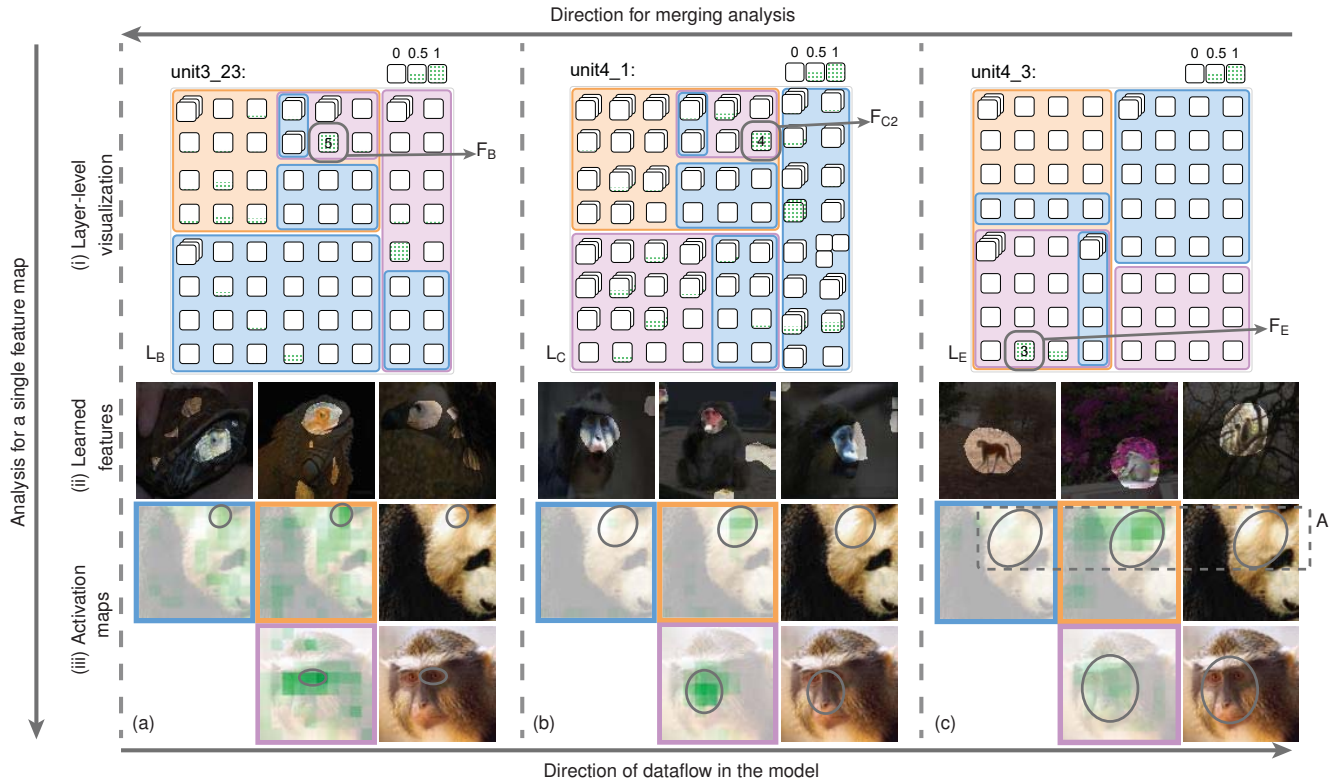


Figure 10:  $E_1$ 's analysis process from the deep layer  $L_E$  to the shallow layer  $L_B$  in order to find the major cause of merging.

this feature map, he discovered that the neurons in this feature map were trained to detect a black circle pattern that resembled an ear or an eye of a panda (Fig. 1B and C). Such a pattern is one of the unique characteristics of a panda and is thus critical for its classification. Then looking at the activation maps (Fig. 1C),  $E_1$  noticed that the neurons covering the ear area were correctly activated for the normal panda image, indicating a successful detection of this critical pattern. However, the same neurons were not activated for the adversarial example.  $E_1$  considered this was an important reason for the misclassification.

To analyze the root cause for this failed detection,  $E_1$  selected the neurons covering the ear area to form an area of interest (Fig. 1D) for a closer examination. He suspected the failed detection was influenced by the feature maps from previous layers. He then set the value encoded in each feature map as the 'contribution' to the selected one, i.e. the area of interest in  $F_{C1}$ , and expanded  $L_A$ , the layer at the beginning of the diverging process. In the treemap-based visualization at  $L_A$ ,  $E_1$  found that feature map  $F_A$  (Fig. 1 (b)) which had the largest contribution to the activation difference in  $F_{C1}$ . By examining its learned features (Fig. 1E),  $E_1$  confirmed that it was trained for low-level detection of black-white boundaries. The activation maps (Fig. 9 (a)) showed that this feature was detected in the normal panda image but not in the adversarial example.  $E_1$  speculated that this failed detection led to the failed detection of the panda's ear in  $F_{C1}$ , and finally led to the failed classification of the adversarial example as a panda.

To confirm his speculation,  $E_1$  repeated the analysis on  $L_D$  and  $L_E$ . In these two layers, he selected several feature maps with a big activation difference between the normal and adversarial images. From the activation maps, he found more significant misses in the detection of critical patterns in the adversarial example (Fig. 9 (b),

and (c)). Selecting the area of interest and tracing the contribution back to  $L_A$ ,  $E_1$  identified the same feature map  $F_A$  as the biggest contributor to the failed detection in  $L_D$  and  $L_E$ . At this point,  $E_1$  was convinced that the missing detection of the black-white boundary in  $F_A$  was the root cause for the failed detection of critical patterns in higher levels that finally led to the failed classification of the adversarial example as a panda.

**Merging analysis.** After analyzing the reason why the adversarial example was not classified as a panda,  $E_1$  turned his attention to why it was classified as a monkey. He suspected that the same region that led to the failed classification of panda actually contributed to its classification as a monkey. Therefore, he retained the same area of interest and expanded layer  $L_E$  (Fig. 6), the merging point for the datapaths of the adversarial and the monkey images. To find the feature maps that had the main contribution to the misclassification, he set the encoded value of each feature map as the 'contribution' to the prediction of the adversarial panda image and identified the feature map ( $F_E$  in Fig. 10 (c)) with the largest contribution. After examining its learned features (Fig. 10 (c)(ii)),  $E_1$  discovered that it was trained to detect monkeys in various situations. Comparing the activation maps of the adversarial example and normal monkey image (Fig. 10 (c)(iii)), he found that the monkey face was activated in the monkey image as expected. However, it was hard to explain the activation at the top part of the adversarial example. Intuitively, there were no indications of a monkey in that part. Thus,  $E_1$  decided to trace back to the lower levels to seek more clues.

Guided by the larger activation on the activation map of the adversarial example,  $E_1$  first adjusted the area of interest to include the most activated neurons at this layer (Fig. 1D), and then analyzed the contributions to  $F_E$  from the feature maps in previous layers.



In layer  $L_C$ , he found feature map  $F_{C2}$  (Fig. 10 (b)), which had the highest contribution to  $F_E$ . After a closer look, the neurons in  $F_{C2}$  seemed to detect the face of a monkey (Fig. 10 (b)(ii)). Activation on the monkey image was correctly located in the middle of the monkey face, but in the adversarial panda image, it was again located on the top right part as was the case in  $F_E$  (Fig. 10(c)). ‘Are there patterns of a monkey face?’ With this question in mind,  $E_1$  carefully compared the adversarial example and the monkey image (Fig. 10 (b)(iii)). He found that for the adversarial example, in the activated region, the dark strip with two lighter patches on either side did resemble the look of a monkey’s nose with lighter cheeks next to it.

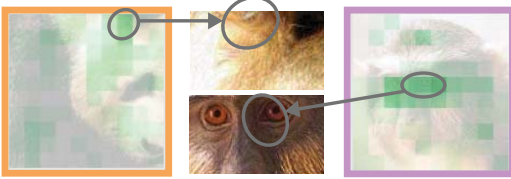


Figure 11: A small region with the appearance of an eye was detected at the top right corner of the adversarial panda image.

However, the same pattern was present in the normal panda image, but it was not detected by the neurons in  $F_{C2}$ .  $E_1$  thus wanted to investigate the more fundamental cause for the detection. Again, he adjusted the area of interest according to the activation map and expanded the previous layer  $L_B$  (Fig. 6). He identified that feature map  $F_B$  (Fig. 10 (a)) was in the datapath intersection for both the adversarial and monkey images and had the largest contribution to  $F_{C2}$ . Examining the learned features of this feature map (Fig. 10 (a)(ii)),  $E_1$  found that the neurons inside were trained to identify the eyes of different animals. Further inspecting the activation of the three images, he found that a small region with the appearance of an eye was detected at the top right corner of the adversarial panda image. (Fig. 11). And compared with the normal image, this seemed attributed to the added adversarial noise. In addition, the position of the ‘eye’ combined with the position of the ‘nose’ detected in  $F_E$  resembled the layout of a real monkey face. At this point,  $E_1$  figured out the effect of the adversarial noise. The top right corner of the adversarial panda image had some similarities with a monkey’s face. In particular, the imperceptible adversarial noise misled the model to detect a monkey’s eye first, then the subtle changes in image layout misled the model to detect a monkey’s face. It was like a domino effect and finally led to the misclassification of the adversarial example as a monkey.

**Summary.** From the above analysis,  $E_1$  summarized two effects of the adversarial noise. The first one was that the outline of the panda’s ear was affected by the noise, which led to the failed detection of the ear and resulted in the large decrease of the predicted probability of the panda class. The second one is that the adversarial noise misled the model to detect a monkey’s eye in the same region, which further led to the detection of a monkey’s face. Then the probability of monkey class largely increased and resulted in the final misclassification.

### 6.2.2 Cannon image

$E_1$  carried out a similar examination on the adversarial cannon image. Examining the datapaths of the normal and adversarial cannon images and a representative racket image (Fig. 12),  $E_1$

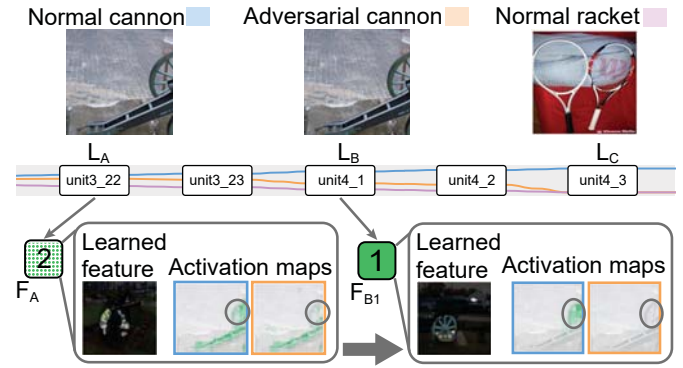


Figure 12: Diverging analysis for an adversarial cannon image.

first identified the diverging and merging points ( $L_B$  and  $L_C$  in Fig. 12). To analyze the root cause for divergence,  $E_1$  followed the same process as above. He first discovered the activation difference between the normal and adversarial images on feature map  $F_{B1}$  that was trained to detect the wheel of a cannon (Fig. 12). He then traced the failed detection to feature map  $F_A$  in  $L_A$ , where the wheel shafts were not detected in the adversarial example (Fig. 12).  $E_1$  speculated the added noise blurred the edges of the shafts which resulted in the failed detection, and further led to the failed detection of the wheel, and finally the misclassification.

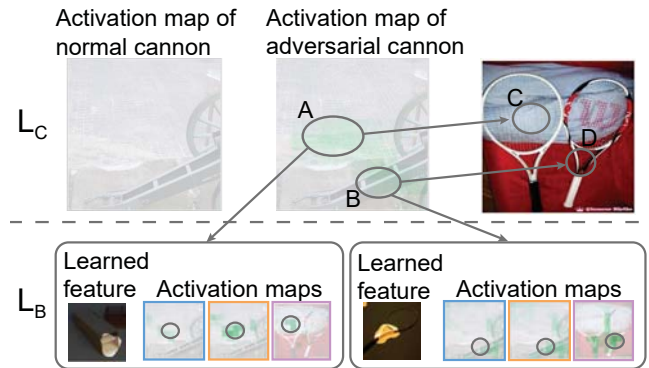


Figure 13: Merging analysis for adversarial cannon image.

To understand why the adversarial example was misclassified as a racket,  $E_1$  turned his attention to the merging point  $L_C$ , and identified the feature map  $F_C$  which had the largest contribution to the misclassification (Fig. 13). Comparing the activation maps (Fig. 13  $L_C$ ), he found two regions (Fig. 13A and B) that were wrongly activated in the adversarial cannon image. Selecting each region and tracing back to layer  $L_B$ ,  $E_1$  noticed the feature maps that contributed the most to each of the regions were trained to detect ‘net’ and ‘racket throat’ respectively (Fig. 13C and D). There are similarities between the streaks on the ground and a racket net, and between the gun mount and a racket throat. However, the added noise in the adversarial image creates a stronger activation in the two feature maps (Fig. 13  $L_B$ ).  $E_1$  thus speculated that the stronger activation together with the failed detection of the wheel misled the model to detect the net and the throat of a racket in adjacent regions, which finally led to the misclassification of the adversarial image.

## 7 DISCUSSION

AEVis can effectively illustrate the prediction mechanism of adversarial examples and help discover the root cause that leads to incorrect predictions. However, it still has several limitations, which may shed light on future research directions.

**Visual scalability.** We have demonstrated that AEVis is able to analyze a state-of-the-art CNN (ResNet101), which has 101 layers and is much deeper than traditional CNNs (e.g., VGG-Net). More recently, deeper CNNs with thousands of layers [7] have been developed. When handling such deep neural networks, the layers of interest at low levels of the hierarchy are difficult to fit in one screen, even with the help of our segmented DAG. A possible solution to alleviate this issue is to employ a mini-map to help the expert track the current viewpoint, which has proven effective in TensorFlow [30].

Currently, we utilize a river-based visual metaphor to illustrate the diverging and merging patterns. The layout of the datapaths is calculated using a rule-based method (Sec. 5.2). Such a design echoes the most common analytical task when three datapaths need to be compared (adversarial examples, the normal source examples, and the normal target examples). If more datapaths are to be analyzed, an optimization-based layout method can be applied. For example, we can minimize the mean-square-error between the vector of real datapath distances and their screen distances with a constraint so that the order of real datapath distances is maintained. The above optimization problem is convex (convex functions over convex sets) and guaranteed to achieve a global minimum. As we have not observed such needs, we leave this method in the discussion here. Apart from the river-based visualization, the treemap-based visualization in the layer level is the other factor that limits the ability to analyze a lot of datapaths. The intuitive treemap-based design is suitable for comparing several datapaths [49] and has been proven effective in the case studies. We can further improve its scalability by adopting a less intuitive but more scalable set of visualization techniques, such as PowerSet [50].

**Generalization.** AEVis aims to analyze the adversarial examples for CNNs because most research on adversarial attacks focuses on generating adversarial images for CNNs.

In addition to attacking CNNs, there are several initial attempts to attack other types of DNNs [8], such as recurrent neural networks (RNNs), autoencoders (AEs), and deep generative models (DGMs). In these types of DNNs, there are also neurons that are critical for predictions. For example, Ming et al. [19] demonstrated that some neurons in an RNN were critical for predicting the sentiment of a sentence, such as the neurons for detecting positive/negative words. Such neurons and their connections form a datapath for an RNN. Thus, AEVis can be extended to help understand the root cause of adversarial examples for these DNNs. The main extension required is the development of suitable datapath extraction and visualization methods for different types of DNNs. For example, to visualize the datapath of RNNs, we can first unfold the architecture of an RNN to a DAG [51], and then employ a DAG layout algorithm to calculate the position of each unfolded layer.

In addition to images, there are adversarial attacks on other types of data [8], such as adversarial documents and adversarial videos. To generalize AEVis to different types of data, we need to change the visual hint for neurons (learned features and activation maps) according to the target data type. For example, when analyzing adversarial documents, we can use a word cloud to

represent the ‘learned feature’ of a neuron [19], and select the keywords that strongly activate the neuron.

## 8 CONCLUSION

We have presented a robustness-motivated visual analysis tool, AEVis, to help machine learning experts investigate the prediction process and understand the root cause of incorrect predictions of adversarial examples. The visualization at multiple levels, together with the constrained datapath extraction, allows efficient identification of critical layers from datapaths’ diverging-merging patterns and critical neurons from the activation maps. The contribution analysis and the rich interactions further enable users to trace the root cause of the misclassification of adversarial examples. We conducted a quantitative experiment to evaluate the datapath extraction method and a representative case study with an expert to demonstrate the usefulness of AEVis in explaining the reasons behind the misclassification of adversarial examples.

There are several directions we could follow in our future research. First, based on the discovered root cause for misclassification, an interesting and important step forward is to develop targeted defense solutions. We will continue working with machine learning experts to explore an effective route from discovered cause to targeted solutions for developing more adversarial robust DNN models. Second, complementary to developing defense solutions, another avenue is to detect potential adversarial examples online and remove them from further processing. A set of streaming visualizations that can incrementally integrate the incoming log data with existing data is the key to online monitoring. Third, as discussed in Sec. 7, an interesting direction is to generalize AEVis to analyze the noise robustness of other types of DNNs, and to tackle other types of data. Improving the scalability to deeper DNNs and the visualization of more datapaths is also an area of future interest. Different datapath extraction algorithms and suitable visualization designs would be interesting research topics.

## REFERENCES

- [1] Y. LeCun, Y. Bengio, and G. Hinton, “Deep learning,” *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [2] S. Liu, X. Wang, M. Liu, and J. Zhu, “Towards better analysis of machine learning models: A visual analytics perspective,” *Visual Informatics*, vol. 1, no. 1, pp. 48–56, 2017.
- [3] S. M. Moosavi-Dezfooli, A. Fawzi, O. Fawzi, and P. Frossard, “Universal adversarial perturbations,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 86–94.
- [4] A. Nguyen, J. Yosinski, and J. Clune, “Deep neural networks are easily fooled: High confidence predictions for unrecognizable images,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 427–436.
- [5] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, “Intriguing properties of neural networks,” in *Proceedings of the International Conference on Learning Representations*, 2014.
- [6] S. Zheng, Y. Song, T. Leung, and I. Goodfellow, “Improving the robustness of deep neural networks via stability training,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 4480–4488.
- [7] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [8] N. Akhtar and A. Mian, “Threat of adversarial attacks on deep learning in computer vision: A survey,” *IEEE Access*, vol. 6, no. 1, pp. 14 410–14 430, 2018.
- [9] S. Komkov and A. Petiushko, “Advhat: Real-world adversarial attack on arcface face id system,” *arXiv preprint arXiv:1908.08705*, 2019.
- [10] Y. Wang, H. Su, B. Zhang, and X. Hu, “Interpret neural networks by identifying critical data routing paths,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 8906–8914.



- [11] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," in *Proceedings of the International Conference on Learning Representations*, 2015.
- [12] F. Liao, M. Liang, Y. Dong, T. Pang, X. Hu, and J. Zhu, "Defense against adversarial attacks using high-level representation guided denoiser," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 1778–1787.
- [13] T. Pang, C. Du, Y. Dong, and J. Zhu, "Towards robust detection of adversarial examples," in *Proceedings of the Advances in Neural Information Processing Systems*, 2018, pp. 4579–4589.
- [14] Y. Bengio, A. Courville, and P. Vincent, "Representation learning: A review and new perspectives," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 8, pp. 1798–1828, 2013.
- [15] M. Liu, S. Liu, H. Su, K. Cao, and J. Zhu, "Analyzing the noise robustness of deep neural networks," in *Proceedings of the IEEE Conference on Visual Analytics Science and Technology*, 2018, pp. 60–71.
- [16] N. Pezzotti, T. Hilt, J. Van Gemert, B. P. F. Lelieveldt, E. Eiseemann, and A. Vilanova, "DeepEyes: Progressive visual analytics for designing deep neural networks," *IEEE Transactions on Visualization and Computer Graphics*, vol. 24, no. 1, pp. 98–108, 2018.
- [17] S. Nie, C. Healey, K. Padia, S. Leeman-Munk, J. Benson, D. Cairra, S. Sethi, and R. Devarajan, "Visualizing deep neural networks for text analytics," in *2018 IEEE Pacific Visualization Symposium (PacificVis)*, 2018, pp. 180–189.
- [18] A. Bilal, A. Jourabloo, M. Ye, X. Liu, and L. Ren, "Do convolutional neural networks learn class hierarchy?" *IEEE Transactions on Visualization and Computer Graphics*, vol. 24, no. 1, pp. 152–162, 2018.
- [19] Y. Ming, S. Cao, R. Zhang, Z. Li, Y. Chen, Y. Song, and H. Qu, "Understanding hidden memories of recurrent neural networks," in *2017 IEEE Conference on Visual Analytics Science and Technology*, 2017, pp. 13–24.
- [20] H. Strobelt, S. Gehrmann, H. Pfister, and A. M. Rush, "LSTMVis: A tool for visual analysis of hidden state dynamics in recurrent neural networks," *IEEE Transactions on Visualization and Computer Graphics*, vol. 24, no. 1, pp. 667–676, 2018.
- [21] B. C. Kwon, M. Choi, J. T. Kim, E. Choi, Y. B. Kim, S. Kwon, J. Sun, and J. Choo, "RetainVis: Visual analytics with interpretable and interactive recurrent neural networks on electronic medical records," *IEEE Transactions on Visualization and Computer Graphics*, vol. 25, no. 1, pp. 299–309, 2019.
- [22] H. Strobelt, S. Gehrmann, M. Behrisch, A. Perer, H. Pfister, and A. M. Rush, "Seq2seq-Vis: A visual debugging tool for sequence-to-sequence models," *IEEE Transactions on Visualization and Computer Graphics*, vol. 25, no. 1, pp. 353–363, 2019.
- [23] M. Liu, J. Shi, K. Cao, J. Zhu, and S. Liu, "Analyzing the training processes of deep generative models," *IEEE Transactions on Visualization and Computer Graphics*, vol. 24, no. 1, pp. 77–87, 2018.
- [24] J. Wang, L. Gou, H. Yang, and H. Shen, "GANViz: A visual analytics approach to understand the adversarial game," *IEEE Transactions on Visualization and Computer Graphics*, vol. 24, no. 6, pp. 1905–1917, 2018.
- [25] M. Kahng, N. Thorat, D. H. Chau, F. B. Vigas, and M. Wattenberg, "GAN Lab: Understanding complex deep generative models using interactive visual experimentation," *IEEE Transactions on Visualization and Computer Graphics*, vol. 25, no. 1, pp. 310–320, 2019.
- [26] J. Wang, L. Gou, H. Shen, and H. Yang, "DQNViz: A visual analytics approach to understand deep Q-networks," *IEEE Transactions on Visualization and Computer Graphics*, vol. 25, no. 1, pp. 288–298, 2019.
- [27] F. Hohman, M. Kahng, R. Pienta, and D. H. Chau, "Visual analytics in deep learning: An interrogative survey for the next frontiers," *IEEE Transactions on Visualization and Computer Graphics*, vol. 25, no. 8, pp. 2674–2693, 2019.
- [28] M. Liu, J. Shi, Z. Li, C. Li, J. Zhu, and S. Liu, "Towards better analysis of deep convolutional neural networks," *IEEE Transactions on Visualization and Computer Graphics*, vol. 23, no. 1, pp. 91–100, 2017.
- [29] F. Y. Tzeng and K. L. Ma, "Opening the black box - data driven visualization of neural networks," in *Proceedings of the IEEE Visualization*, 2005, pp. 383–390.
- [30] K. Wongsuphasawat, D. Smilkov, J. Wexler, J. Wilson, D. Man, D. Fritz, D. Krishnan, F. B. Vigas, and M. Wattenberg, "Visualizing dataflow graphs of deep learning models in TensorFlow," *IEEE Transactions on Visualization and Computer Graphics*, vol. 24, no. 1, pp. 1–12, 2018.
- [31] P. E. Rauber, S. G. Fadel, A. X. Falco, and A. C. Telea, "Visualizing the hidden activity of artificial neural networks," *IEEE Transactions on Visualization and Computer Graphics*, vol. 23, no. 1, pp. 101–110, 2017.
- [32] J. Wang, L. Gou, W. Zhang, H. Yang, and H. Shen, "DeepVID: Deep visual interpretation and diagnosis for image classifiers via knowledge distillation," *IEEE Transactions on Visualization and Computer Graphics*, vol. 25, no. 6, pp. 2168–2180, 2019.
- [33] Y. Ma, T. Xie, J. Li, and R. Maciejewski, "Explaining vulnerabilities to adversarial machine learning through visual analytics," *IEEE Transactions on Visualization and Computer Graphics (accepted)*, pp. 1–1, 2019.
- [34] A. W. Harley, "An interactive node-link visualization of convolutional neural networks," in *Proceedings of the International Symposium on Visual Computing*, 2015, pp. 867–877.
- [35] M. Kahng, P. Y. Andrews, A. Kalro, and D. H. Chau, "ActiVis: Visual exploration of industry-scale deep neural network models," *IEEE Transactions on Visualization and Computer Graphics*, vol. 24, no. 1, pp. 88–97, 2018.
- [36] L. v. d. Maaten and G. Hinton, "Visualizing data using t-SNE," *Journal of Machine Learning Research*, vol. 9, no. Nov, pp. 2579–2605, 2008.
- [37] Y. Dong, F. Liao, T. Pang, H. Su, J. Zhu, X. Hu, and J. Li, "Boosting adversarial attacks with momentum," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 9185–9193.
- [38] N. I. P. Systems, "Nips 2017: Adversarial attack," <https://nips.cc/Conferences/2017/CompetitionTrack>, 2017, Last accessed 2019-10-16.
- [39] M. D. Zeiler and R. Fergus, "Visualizing and understanding convolutional networks," in *Proceedings of the European Conference on Computer Vision*, 2014, pp. 818–833.
- [40] E. Alexander and M. Gleicher, "Task-driven comparison of topic models," *IEEE Transactions on Visualization and Computer Graphics*, vol. 22, no. 1, pp. 320–329, 2016.
- [41] M. Gleicher, "Considerations for visualizing comparison," *IEEE Transactions on Visualization and Computer Graphics*, vol. 24, no. 1, pp. 413–423, 2018.
- [42] R. Tibshirani, "Regression shrinkage and selection via the lasso," *Journal of the Royal Statistical Society: Series B (Methodological)*, vol. 58, no. 1, pp. 267–288, 1996.
- [43] A. Krizhevsky, "Learning multiple layers of features from tiny images," University of Montreal, Tech. Rep., 2009.
- [44] T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein, *Introduction to algorithms*. MIT press, 2009.
- [45] W. Cui, S. Liu, L. Tan, C. Shi, Y. Song, Z. Gao, H. Qu, and X. Tong, "TextFlow: Towards better understanding of evolving topics in text," *IEEE Transactions on Visualization and Computer Graphics*, vol. 17, no. 12, pp. 2412–2421, 2011.
- [46] M. Bruls, K. Huizing, and J. J. Van Wijk, "Squarified treemaps," in *Data Visualization*. Springer, 2000, pp. 33–42.
- [47] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, "Imagenet large scale visual recognition challenge," *International Journal of Computer Vision*, vol. 115, no. 3, pp. 211–252, 2015.
- [48] Y. Dong, "Non targeted adversarial attacks," <https://github.com/dongyp13/Non-Targeted-Adversarial-Attacks>, 2017, Last accessed 2019-10-16.
- [49] B. Alsallakh, L. Micallef, W. Aigner, H. Hauser, S. Miksch, and P. Rodgers, "Visualizing sets and set-typed data: State-of-the-art and future challenges," in *Proceedings of the Eurographics Conference on Visualization*, 2014, pp. 1–21.
- [50] B. Alsallakh and L. Ren, "PowerSet: A comprehensive visualization of set intersections," *IEEE Transactions on Visualization and Computer Graphics*, vol. 23, no. 1, pp. 361–370, 2017.
- [51] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016.

**Kelei Cao** is now a master student at Tsinghua University. His research interest is visual analytics for explainable deep learning. He received a BS degree in the Department of Computer Science and Technology, Tsinghua University.







**Mengchen Liu** is a Senior Researcher at Microsoft. His research interest includes explainable AI and computer vision. He received a BS in Electronics Engineering and a Ph.D in Computer Science from Tsinghua University. He has served as PC members and reviewers in various conferences and journals.



**Hang Su** is an assistant professor at Tsinghua University. He received his B.S., M.S., and Ph.D. Degrees from Shanghai Jiaotong University. His research interests lie in the development of computer vision and machine learning algorithms for solving scientific and engineering problems arising from artificial learning, reasoning, and decision-making. His current work involves the foundations of interpretable machine learning and the applications of image/video analysis. He has served as senior PC members in the dominant international conferences.



**Jing Wu** is a lecturer in computer science and informatics at Cardiff University, UK. Her research interests are in computer vision and graphics including image-based 3D reconstruction, face recognition, machine learning and visual analytics. She received BSc and MSc from Nanjing University, and PhD from the University of York, UK. She serves as a PC member in CGVC, BMVC, etc., and is an active reviewer for journals including Pattern Recognition, Computer Graphics Forum, etc.



**Jun Zhu** received his BS and PhD degrees from the Department of Computer Science and Technology in Tsinghua University, where he is currently a professor. He was an adjunct faculty and postdoctoral fellow in the Machine Learning Department, Carnegie Mellon University. His research interests are primarily on developing statistical machine learning methods to understand scientific and engineering data arising from various fields. He regularly serves as Area Chairs at prestigious conferences, such as ICML, NeurIPS. He is an associate editor-in-chief of IEEE TPAMI.



**Shixia Liu** is an associate professor at Tsinghua University. Her research interests include visual text analytics, visual social analytics, interactive machine learning, and text mining. She worked as a research staff member at IBM China Research Lab and a lead researcher at Microsoft Research Asia. She received a B.S. and M.S. from Harbin Institute of Technology, a Ph.D. from Tsinghua University. She is an associate editor-in-chief of IEEE Trans. Vis. Comput. Graph.