



A Task-and-Technique Centered Survey on Visual Analytics for Deep Learning Model Engineering

ARTICLE INFO

Article history:

Received July 13, 2018

Keywords: Visual Analytics, Survey, Deep Learning, Machine Learning, Neural Networks, Visualization, High-Dimensional Visualization

ABSTRACT

Although deep neural networks have achieved state-of-the-art performance in several artificial intelligence applications in the past decade, they are still hard to understand. In particular, the features learned by deep networks when determining whether a given input belongs to a specific class are only *implicitly* described concerning a considerable number of internal model parameters. This makes it harder to construct interpretable hypotheses of *what* the network is learning and *how* it is learning—both of which are essential when designing and improving a deep model to tackle a particular learning task. This challenge can be addressed by the use of visualization tools that allow machine learning experts to explore which components of a network are learning useful features for a pattern recognition task, and also to identify characteristics of the network that can be changed to improve its performance. We present a review of modern approaches aiming to use visual analytics and information visualization techniques to understand, interpret, and fine-tune deep learning models. For this, we propose a taxonomy of such approaches based on whether they provide tools for visualizing a network's architecture, to facilitate the interpretation and analysis of the training process, or to allow for feature understanding. Next, we detail how these approaches tackle the tasks above for three common deep architectures: deep feedforward networks, convolutional neural networks, and recurrent neural networks. Additionally, we discuss the challenges faced by each network architecture and outline promising topics for future research in visualization techniques for deep learning models.

© 2018 Elsevier B.V. All rights reserved.

1. Introduction

One of the main goals of Artificial Intelligence (AI) is to build systems that achieve human-level efficiency in recognition tasks, such as image classification, speech recognition, and sentiment analysis. Although most of these tasks seem trivial to human beings, they are extremely challenging for computer algorithms due to the lack of a formal description of how to solve such problems. For example, humans can, in general, easily recognize if there is a dog in a given image, but it is hard to tell how we got to this conclusion. In other words, it is not clear how to formalize which features in the image makes humans recognize the presence of dogs [1]. Is it the shape of the objects? Is it the color contrast between different regions

in the image? Moreover, even harder, how do humans learn to recognize dogs in the first place? How do we learn to use such features to identify dogs? How can we teach such learning abilities to machines? One of the areas of AI that focus on finding solutions to approach these problems is called Machine Learning (ML). ML algorithms use statistical techniques to optimize functions to progressively achieve better performances in a particular task [2]. To train an ML model, the designer must provide the model with training inputs with known answers—e.g., images of dogs and images without dogs—the model thus automatically learns to model a function that minimizes the chances of wrongly predicting such inputs.

In the area of machine learning, a particular class of techniques has proven increasingly efficient and effective in pattern

recognition applications in the past years: deep learning (DL) techniques. In contrast to what may be called ‘classical’ ML techniques, DL techniques do not rely on the designer to provide a set of hand-engineered features to be used in the learning and decision processes. Rather, they rely on a (large) set of labeled samples to automatically extract and store such features in the so-called architecture of a deep neural network (DNN) [3]. DNNs contain multiple (up to hundreds of) layers that perform simple filtering, thresholding, and aggregation operations on subsets of the input data samples. The power of such architectures relies (1) on their ability to model very complex non-linear decision functions and decision boundaries in the input data space by combining many such simple operations; and (2) on the fact that the set of learned parameters involved in implementing such operations (also known as the *model* learned by the network) can be automatically inferred from a (typically large) set of labeled data samples, by minimizing the error between the predicted and the ground-truth labels.

DNNs have recently shown excellent performance in pattern recognition tasks in part due to the increase in computing power available for training ever larger architectures, *e.g.* by using GPU computing. Early on, Krizhevsky *et al.* proposed the AlexNet network [4], a convolutional neural network (CNN) architecture with six hidden layers that won the 2012 ILSVRC competition on the well-known ImageNet dataset [5], with a top 5 test error rate of 15.4%; i.e. the percentage of images on the test set whose true label was not among the 5 classes considered to be more likely by the model. This architecture was later refined by Zeiler and Fergus [6], which decreased the training set size by one order of magnitude while providing an improved 11.2% error rate. In that publication, the authors have argued that often the “development of better models is reduced to trial and error”, and proposed a pioneering visualization for depicting feature maps to aid the understanding of the network training process. By increasing the network size (or *depth*, as measured in terms of number of layers, and the training set size, subsequent works managed to provide further increases in task accuracy; *e.g.*, in the work of Simonyan and Zisserman [7], where a network with over 100 layers was used; in GoogLeNet [8] (22 layers, achieving top 5 error 6.7% for ILSVRC); and in Microsoft’s ResNet [9] (152 layers, where a 3.6% error was achieved for ILSVRC). Several other similar examples have been published in the past few years.

Even though higher accuracy may be achieved by increasing the architectural complexity of DNNs, such a design choice also brings about several important challenges in deploying such networks. First, the computational training effort required by such systems becomes quite high—DNNs can often require days or even weeks of training time even when using a system composed of many GPUs. Suboptimal configurations of the network architecture, such as its training set size and the network hyperparameters, can also require the entire learning process to be restarted or run from scratch, which is very expensive. As a consequence of these challenges, it becomes increasingly important for one to be able to *understand* the complex interaction between all of these aspects (*i.e.*, the characteristics and size of the training and testing datasets; the net-

work architecture; the network hyperparameters; and the test results) in order to effectively fine-tune a DNN for a specific task. This is by itself a very complex challenge as raised by Liu et al [10], grounded chiefly in two reasons: (1) the size of the space spanned by all these design dimensions is *huge*, so an exhaustive exploration thereof is impossible; moreover, this space is highly non-linear, since small changes to the hyperparameters may cause significant changes to the performance of the corresponding trained model; (2) the *abstract* nature of this space—*i.e.*, the fact that it is hard to intuitively understand the effect of particular changes to hyperparameters in the corresponding performance—makes it hard to understand how and why a DNN behaves in a certain way.

The above constraints determine that, in practice, most designers construct and train their DNN models virtually as ‘black box’ algorithms. More specifically, while designers do have, at least at a high conceptual level, a relatively good idea of what pooling, drop-out, and convolutional layers implement, the joint effect of combining several (tens to hundreds) of such layers and varying their parameters can only, in practice, be assessed via empirical end-to-end measurements of the performance of the resulting DNN. This introduces several important and, so far, only partially answered questions:

- *What has a model learned?* Without a good understanding of *what* decision hypothesis a given trained network has acquired, users and designers of a DNN typically have no idea about what regions of the original data space (sampled by the training process) were ‘internalized’ by the resulting model. As such, it is not clear how to allocate further training efforts to improve the network. Also, a model may have poor generalization performance. Currently, this can only be verified by testing the network on novel data (validation) or by deploying it in field operations—when it may be too late to detect generalization issues. Understanding what a model has inherently learned (or not) is therefore of key importance.
- *Why has a model learned a particular decision hypothesis?* Related to the above point is the challenge of understanding why a model has learned (or not) to generalize certain aspects of the training data. Knowing this may directly provide feedback to the designer as to what needs to be changed in the network architecture, hyperparameters, or in the training data to further strengthen desired properties of the network or to address particular issues. Moreover, this will also give a better understanding of properties of the input data space, *e.g.* in terms of sub-spaces that are particularly challenging for the learning process. Without knowing why a particular model resulted from the training process, DNN optimization can very much be a blind search process.
- *How has a model learned a particular decision hypothesis?* Even if a DNN performs well, one may want to understand how it has internally stored knowledge about the training data; *e.g.*, in which specific layers (or parts thereof) or weight value ranges were particular types of knowledge encoded. This may help in understanding the

1 reason why specific DNN architectures are appropriate for
 2 specific tasks and thus help to extrapolate such knowledge
 3 when tackling new problems. Without this, addressing a
 4 new problem in a different data space might require starting
 5 the entire network engineering process from scratch.

6 Recently, these concerns have been discussed primarily by
 7 the DL community. For example, Marcus [11] argues that the
 8 DL field may be ‘approaching a wall’ and outlines ten chal-
 9 lenges: (1) requirements for huge amounts of labeled data; (2)
 10 limited capacity for transfer between problems; (3) difficulty of
 11 dealing with hierarchical structure; (4) difficulty to deal with
 12 open-ended inference; (5) *DL is not sufficiently transparent*;
 13 (6) it is hard to integrate prior knowledge; (7) it is hard to dis-
 14 tinguish correlation from causation; (8) assumption of a stable
 15 world; (9) *DL’s answers cannot be always trusted*; (10) *DL is*
 16 *hard to engineer with*. Among these, we focus here on chal-
 17 lenges (5), (9), and (10), which directly relate to our previously-
 18 made points regarding the ‘black box,’ unpredictable, and hard
 19 to fine-tune nature of DNNs, respectively. Similar concerns re-
 20 garding these issues have been expressed in the works of Samek
 21 et al. [12] and Ribeiro et al [13].

22 The need to address the above challenges has been rec-
 23 ognized by scientists at the confluence of several domains
 24 (data science, machine learning, and data visualization). One
 25 particular approach to achieving this goal, which we survey in
 26 this paper, is to use *visualization* techniques and tools. This
 27 approach is rooted in earlier works related to understanding
 28 high-dimensional data spaces [14, 15, 16, 17, 18] and visualiz-
 29 ing the operation of classical ML algorithms—in particular,
 30 visualizing feature spaces and the impact of using different
 31 feature selection processes [19, 20, 21, 22, 23, 24, 25, 26, 27].
 32 However, such earlier techniques do not directly address
 33 the challenges that are particular to DNNs, such as the ones
 34 mentioned earlier in this section, and instead focus on un-
 35 derstanding more general correlations between the network
 36 structure, the high-dimensional input-feature spaces used
 37 during training, and the resulting network performance. More
 38 recently, novel techniques in information visualization and
 39 visual analytics (VA) have aimed at supporting and improving
 40 DNN engineering by taking into account their particular
 41 challenges; however, the relatively fast growth of this field has
 42 not yet been fully covered by existing surveys.

43
 44 **Contributions:** We aim to alleviate the problems discussed
 45 above with the following contributions:

- 46 1. we propose a task-and-architecture based taxonomy of
 47 visualization techniques that help engineering DNNs by
 48 whether they tackle one of *three possible tasks*: (1) facil-
 49 itating the visualization of the network structure; (2) facil-
 50 itating the interpretation and analysis of the training pro-
 51 cess; or (3) allowing for feature understanding. We ex-
 52 plain the particularities of these different tasks, discussing
 53 their goals, their challenges and how they are applied in
 54 the context of three types of network architecture: feed-
 55 forward networks, convolutional neural networks, and re-
 56 current neural networks;

- 57 2. we survey a comprehensive body of over 40 papers re-
 58 lated to visualization in DNN engineering, and explain
 59 how these fit within the proposed task taxonomy;
 60 3. we outline important limitations of current work in the area
 61 and suggest directions for future exploration.

62 The structure of this survey is as follows. In section 2, we
 63 introduce import concepts and notation regarding classical
 64 machine learning and deep learning engineering. In Section 3,
 65 we detail the above-mentioned tasks and how visual analytics
 66 techniques relate to the main deep learning architectures in
 67 order to complete such tasks. In Section 5 we introduce
 68 our taxonomy and discuss in details different visualization
 69 techniques used for tackling these tasks in the context of deep
 70 feedforward networks, convolutional neural networks, and
 71 recurrent neural networks. Section 6 discusses how these
 72 techniques cover the needs of different types of DNN designers
 73 and also outlines important technical limitations which can
 74 spawn future research directions. Section 7 concludes the paper.

75 **Relation with other surveys:** Several other surveys partially
 76 touch our focus of interest. The most important existing surveys
 77 include tutorials on deep learning visualization [28, 29, 30, 31];
 78 visualization of convolutional networks [32, 33]; visualiza-
 79 tion of machine learning models [34]; predictive visual analyt-
 80 ics [35]; interactive machine learning [36, 37, 38]; interpretable
 81 machine learning [39]; and surveys of multidimensional visual-
 82 ization techniques [14, 16, 17]. Closer to our focus, Hohman *et*
 83 *al.* [40] present a survey on visual analytics for deep learning.
 84 Their survey follows a human-centered approach to answer the
 85 following questions: (1) why to visualize different aspects of a
 86 deep model or its corresponding training process; (2) who uses
 87 deep learning visualization; (3) what to visualize in deep learn-
 88 ing; (4) how to visualize deep learning; and (5) when in the pro-
 89 cess of designing and training a network the visualization pro-
 90 cess will take place—e.g., during the network engineering step;
 91 throughout training; or after the model parameters and corre-
 92 sponding features are learned. Our survey also addresses these
 93 questions but offers another angle of attack that looks at visual-
 94 ization techniques classified by the task and subtask it addresses
 95 and which model architecture (DFNs, CNNs, RNNs) they apply
 96 it. Furthermore, we focus specifically on how such tools
 97 support a visual analytics approach for solving the above tasks,
 98 rather than on the more general perspective of how to visualize
 99 deep learning data. As such, our survey is complementary to,
 100 and also extends, the work of Hohman *et al.*

2. Classical Machine Learning and Deep Learning

102 In this section, we introduce basic notation relevant both for
 103 classical machine learning algorithms and also for deep learn-
 104 ing models. Later, we will clarify what is the main difference in
 105 how these techniques work.

106 Let $\mathbf{x} \in X$ be an input (e.g., an image or a sound) drawn
 107 from a set or distribution of possible inputs, X (the set of all
 108 possible images). Let $y \in Y$ be a label or output associated
 109 with a particular input \mathbf{x} . If the model is tackling a *regression*

task, y is a (possibly high-dimensional) continuous value. Otherwise, in *classification* tasks, y is a discrete label associated to one or more members of a finite set of classes which elements of X can belong to. Both regression and classification are considered supervised learning techniques, as they are trained with inputs which the actual labels are known. For instance, if \mathbf{x} is an image, Y could be the set {dog, cat}, used to denote the possible animals that may appear on the image. In many practical applications, a particular label y in a set of d possible labels is represented by a vector $\mathbf{y} \in \{0, 1\}^d$, where the i -th element of \mathbf{y} is 1 only if \mathbf{y} corresponds to the i -th possible label in Y . Let $D = \{\mathbf{x}_i, y_i\}$, for $i \in \{1, \dots, N\}$, be a set of N training examples associating particular inputs \mathbf{x} with their corresponding labels y . The objective of a supervised learning algorithm is to analyze a training set D and construct a function $f : X \rightarrow Y$ so that when f is presented with novel inputs (e.g. unseen images) it can correctly predict their corresponding label.

Machine learning algorithms usually optimize their performances by incrementally improving the function f in order to minimize a given cost function C that measures how well f performs; i.e., how well that function correctly predicts the labels of novel inputs. In the regression setting, when f predicts labels that are continuous numbers, C can be a Mean Squared Error such as $C(f) = E_{\mathbf{x} \sim X}[(f(\mathbf{x}) - y)^2]$. In many practical applications, the function f might be easier to learn if the input information is presented to it in a more pre-processed way; for instance, when training an algorithm for detecting cats in pictures, it might be easier to learn an f function that takes as inputs edge and color information instead of raw pixel values. This can be achieved by transforming the inputs $\mathbf{x} \in X$ via a so-called feature function $\Psi : X \rightarrow \mathbb{R}^m$ mapping any element of X to a point in some m -dimensional space. In this case, $f : \mathbb{R}^m \rightarrow \mathbb{R}^d$, i.e., it takes as input some \mathbf{x} and feeds $\Psi(\mathbf{x})$ to f , which produces a prediction in \mathbb{R}^d associated with a given label.

When deploying an algorithm to learn a function f that minimizes the given cost function C , the designer needs to implicitly specify what is the space of possible functions that will be searched over. This is typically done by representing f via a set of so-called *model parameters* $\Theta \in \mathbb{R}^M$. By assigning different numerical values for each of the M model parameters, we obtain a different function f ; for instance, f could be in the form $f(\mathbf{x}) = \theta_0 + \theta_1 x$ in case of a simple linear prediction model. On the other hand, by using a Θ with a much higher number of weights performing a non-linear function (typical case in neural networks), f may become arbitrarily complex and difficult to interpret function. Prior to using a learning algorithm, it is common to group the feature vectors of all training inputs into a single matrix $\mathbf{X} \in \mathbb{R}^{N \times m}$, which the i -th row is associated with input x_i and where columns store the m components of the feature vector $\Psi(x_i)$. In what follows we denote the i -th row of \mathbf{X} by \mathbf{X}^i .

Classical ML algorithms: Classical ML algorithms include techniques such as *k-nearest neighbors* [41], *support vector machines* [42], *decision trees* [43], *logistic regression* [41], and *random forest classifiers* [44]. These methods, when computing an f function that minimizes the cost C , rely on *manual* specification of a feature function Ψ for transforming arbitrary

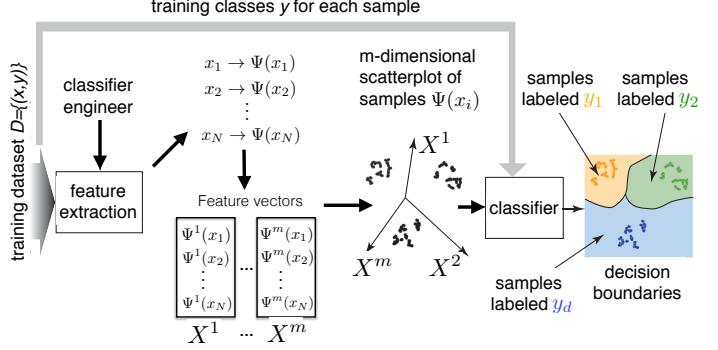


Fig. 1: Design of classical machine learning algorithms. See Sec. 2.

inputs \mathbf{x} into a new representation $\Psi(\mathbf{x}) \in \mathbb{R}^m$ (see Fig. 1). Classical ML algorithms require manual work by the designer in order to construct or identify meaningful features that allow the algorithm to efficiently discriminate between the classes. For instance, Ψ could take as input an image and return another image where the magnitude of each pixel indicates whether an edge exists in that region. Alternatively, Ψ could be the identity function, in which case the learning algorithm operates not over a new feature-based representation of the input \mathbf{x} , but on the values of \mathbf{x} directly. When this is done, the step of computing features is replaced by a step of *selecting* features—essentially, selecting a subset of the columns of \mathbf{X} that the designer identifies as sufficient for the learning algorithm (see Fig. 1). For instance, when classifying whether a person is a male or female based on features corresponding to age, height, weight and eye color, the column corresponding to eye color may be removed by the feature selection process since it is not correlated with gender labels. Manually constructing Ψ and/or manually performing feature selection has two advantages: (1) it may be possible to manually design features that are informative enough so that it becomes easy to identify a clear discrimination rule to determine whether an input belongs to a given label; and (2) features that have an intuitive, application-domain related meaning allow for an effective way to *understand* the entire process of engineering a learning algorithm and analyzing its training process.

Deep learning: Deep learning (DL) techniques target the cases when assumption (1) above does not hold. In such cases, both the extraction of the features \mathbf{X}^j and the construction of f are automatically performed by a Deep Neural Network (DNN). A neural network is one particular form of representing a prediction function f . It is a graph of connected *neurons* typically organized in $L > 2$ layers. The purpose of each layer is to further transform the input data given to the network, automatically building new and more abstract feature representations of it which allow for more efficient classification of that input. Each layer l is composed by a set of neurons or units u_i^l (Fig. 2). Neurons in the first layer take as input all attributes $\Psi(x)$ of a given input data \mathbf{x} ; neurons in subsequent layers ($l > 1$) take as input the output of all neurons of the previous layer. Based on these inputs, each neuron computes its output by linearly combining them and applying some non-linear function, thereby producing an activation or output a_i^l that is passed as input to

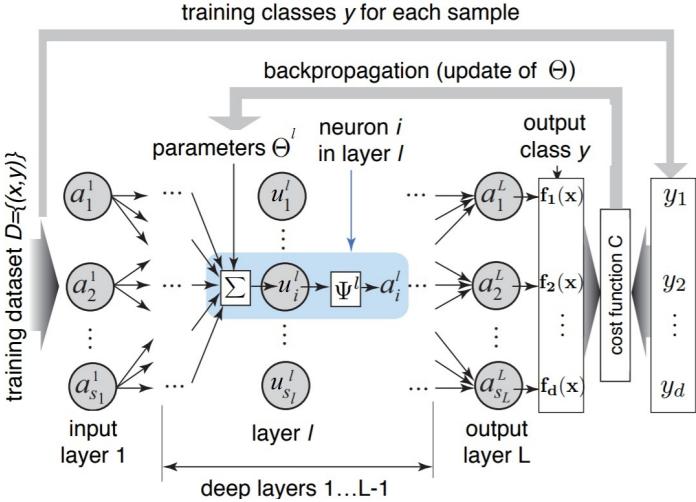


Fig. 2: Architecture of typical DNN (training phase). See Sec. 2.

1 all neurons in the next layer. The output (or *activation* value)
 2 of a neuron in layer l is given as input to neurons in layer $l + 1$
 3 and is associated with a weight parameter θ_{ji} ; in particular, θ_{ji}
 4 is a weight indicating the importance of the output or activation
 5 of the j -th neuron in the previous layer l to the activation of the
 6 i -th neuron in layer $l + 1$. More formally, each neuron j in layer
 7 l computes its output/activation a_i^l as $\phi(\sum_i a_i^{l-1} \theta_{ji})$, where θ_{ji} is
 8 a weight indicating the importance of the activation of the j -th neuron
 9 in the previous layer $l - 1$ to the activation of the i -th neuron
 10 in layer l ; and ϕ is a so-called *activation function*, usually a
 11 non-linear transformation such as tanh, a sigmoid function, or a
 12 rectified linear function [1]. In case the function f modeled as
 13 a DNN outputs $y \in \mathbb{R}^d$, the network's final layer L is composed
 14 of d neurons; when presented with some input x , the network's
 15 label prediction is given by a vector $a^L \in \mathbb{R}^d$ of activations of
 16 each of the d neurons in layer L .

In a DNN, besides layers of neurons as described above, it is also possible to create so-called convolutional layers, which are composed of neurons that essentially implement filters that compute features based on their inputs—e.g., the set of neurons of a particular convolutional layer could implement filters for performing edge detection when given an input image x . The advantage of this in comparison to the approach taken by classical machine learning algorithm is that the filters/features computed by each layer are automatically learned by the algorithm, thereby eliminating the need for manually designing a feature function Ψ . When deploying a DNN, one needs to make a set of design choices: how many layers the network will contain; how many of those will be convolutional layers; which particular activation function ϕ to use in each layer; and how many neurons u_i will compose each layer l . These choices are often referred to as the *architecture* of the DNN and are encoded as a set of hyperparameter P . The choice of hyperparameters is typically manually decided by the DNN engineer based on earlier experience with similar problems. The parameters or weights Θ of the network, on the other hand, are learned by a training algorithm in order to minimize some cost function C . Training algorithms of neural networks are usually based on performing

gradient descent over the cost C with respect to weights Θ on the final layer—*i.e.* by updating $\Theta \leftarrow \Theta - \alpha \nabla_\Theta C$ —and using the *backpropagation algorithm* to propagate this gradient update to previous layer, recursively modifying the weights of each neuron according to the results computed in subsequent layers [1]. However, computing the gradient $\nabla_\Theta C$ requires a linear pass over the entire training set D , which in typical DNN applications may be very large. For this reason, one can alternatively use a so-called mini-batch training process, which splits the dataset D into B subsets of D —each one a smaller data batch compared to the entire set D —and then computes an estimate of the true gradient based on that reduced number of training examples. Updating the network by processing all B batches once is referred to as an *epoch*. This process is repeated for E epochs. Note that this training methodology introduces additional hyperparameters to the problem: the batch size B , a learning rate α , and number of training epochs E . These are determined by the DNN designer before training based on their earlier experiences or heuristics.

Deep neural networks may have several kinds of architectures, each with their own type of neurons performing different set of operations. In this paper, we focus on three architectures that have been very popular in deep learning applications in recent years: *Deep Feedforward Networks* (DFNs); *Convolutional Neural Networks* (CNNs); and *Recurrent Neural Networks* (RNNs). For the purpose of this paper, it is important to differ between these three architecture because the distinctiveness of their neurons brings different challenges when performing the tasks we propose in our taxonomy. As it follows, we introduce in more details the definition of these three architectures.

Deep Feedforward Networks: The most traditional deep learning models are deep feedforward networks (DFNs), also called *multilayer perceptrons* (MLPs). As in other supervised learning approaches, the objective of a DFN is to approximate an unknown function f that can efficiently reproduce the relationship between inputs and outputs of a training set. What distinguishes DFNs from other machine learning techniques is that they are composed of multiple layers, each with multiple neurons. This hugely increases the number of learnable weights the network has (parameters Θ , Fig. 2) and allows it to model functions that are much more complex than those encoded with only a few parameters. In a DFN, layers are fully connected, which means that the i -th neuron in layer l , u_i^l , receives as input a vector a^{l-1} containing the activations of all neurons in the previous layer $l - 1$. The final layer L can have a single neuron—thus, the network produces a single output, used *e.g.* for one-class classification—or multiple neurons, used for discriminative classification goals or multidimensional outputs.

Convolutional Neural Networks: In recent years, many applications in image classification and pattern recognition achieved state-of-the-art performance through the usage of Convolutional Neural Networks (CNNs) [45]. CNNs specialize standard DFNs as they use convolution operations in at least one of their layers. The objective of these operations is to find (small-scale) patterns in the input and send this information to subsequent layers codified via their output activations. The fol-

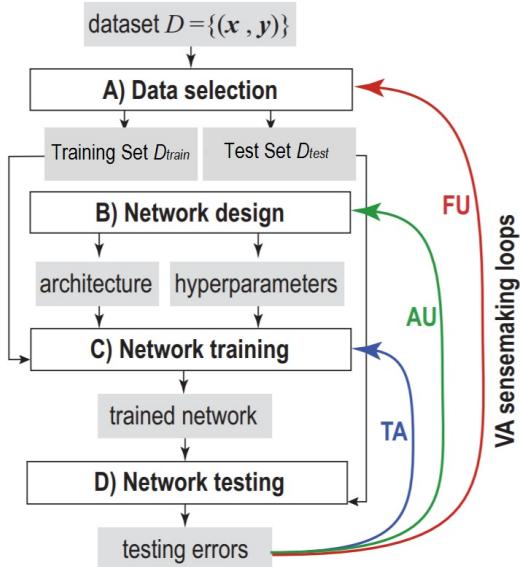


Fig. 3: Visual analytics workflow for DNN workflow engineering support showing the tasks of training analysis (TA), architecture understanding (AU), and feature understanding (FU). See Sec. 3.

lowing layers, in turn, look for more complex patterns, thereby creating a chain of pattern detections that, when trained well, can achieve close-to-human performance in image-related applications. Convolutional layers are usually composed of three stages: multiple convolutional operations; a nonlinear function; and a pooling function that changes the current output value of a layer by aggregating some statistics computed on neighboring outputs [1]. However, the convolutional operation performed by CNNs also bring some challenges to their analysis. A single weight in a convolutional unit acts over every region of the input domain, meaning that small modifications of the weights of a convolutional unit may affect all the domain of the output activation. Additionally, convolutional units usually produce multi-dimensional activation outputs, differently from DFN units, that usually produce a single scalar as activations [10].

Recurrent Neural Networks: Although DFNs and CNNs have achieved impressive results in classification and recognition tasks, they are not suitable for applications where inputs have a temporal or sequential relationship, such as word prediction or machine translation. Recurrent neural networks (RNNs), a different deep learning model, were proposed to handle this type of problem [46]. RNNs are intrinsically different from DFNs and CNNs because their neurons store a different kind of information called *hidden states*. Hidden states have internal values that are combined with the traditional learnable weights Θ of neural networks when computing output activations. In contrast to the parameters Θ , which are frozen after training, the hidden states modify their values each time a new input is processed. This way, the same input instance can, and likely will, generate a different output if the previous inputs in the temporal sequence were different.

3. Deep Learning Engineering: Workflow and Tasks

The process of developing a deep learning model comprises multiple phases. A typical workflow proceeds as follows (Fig. 3). First, given a dataset D consisting of labeled samples (\mathbf{x}, \mathbf{y}) , the designer splits D into two disjoint subsets: the training set D_{train} and the test set D_{test} (Fig. 3A). Such approach is necessary because ML algorithms—particularly those containing a very high number of parameters Θ , as is the case of DNNs—can easily overfit the training data—*i.e.*, they learn to make good predictions for the training set but fail to generalize well to novel inputs, such as the ones in D_{test} . By testing the performance of the network on a set of completely new inputs, one can check whether the model can generalize its predictions to novel data samples that were not in the training set. This selection should be done carefully. Otherwise, even powerful learning algorithms may not perform as desired [1]. Selecting a good training set is not trivial, even if one takes care of issues such as class-label balancing. For instance, it is not clear how well the samples in D_{train} capture the variability of the entire data domain, *i.e.*, how well D_{train} helps to learn all information needed to generalize the task at hand to unseen inputs. The second step in the development workflow is to design a DNN (Fig. 3B). As outlined in Sec. 2, this means choosing an architecture and suitable hyperparameters. Both operations are done largely based on similar designs from the past or heuristics. Yet, it is far from clear how suitable such design choices will be in practice for a given problem. The third step (Fig. 3C) consists of training the designed DNN, following the process detailed in Fig. 2. Finally, the performance of the trained model is measured on the test set D_{test} (Fig. 3D). This is done using aggregated error metrics such as accuracy, precision, recall, or area under the received operator characteristic curve (AUROC) [47, 48]. Alternative schemes involve computing the *confusion matrix* [49] (for classification tasks), which calculates the number of correct predictions for each class and how many times a particular label was mistaken by each one of the alternative labels. This approach is useful for simple models with only a few class choices, but confusion matrices become hard to inspect and visualize for tens or thousands of classes.

If testing delivers satisfactory performance, the workflow in Fig. 3 can be finalized. When this is not the case, the key question is: what can the designer do to improve performance? This involves feedback loops at several workflow levels, each one involving a specific *task* as defined in the Contributions section. Additionally, even models performing successfully bring important questions. Often experts want to understand what kind of features the model learned to recognize or understand how the model operates over an input in order to predict its labels. These are additional tasks that we address in this paper. Note that the tasks that we propose to tackle follow the terminology proposed by Munzner [50], in which they define an arbitrary task as a set of both high-level and low-level, domain-specific (but also data-specific) activities that may be involved in answering the aforementioned question at a specific workflow level. The particular tasks that we consider in this paper are the following:

- 1 **1. Architecture Understanding:** it is important to be able
2 to analyse how the network architecture affects its performance
3 to determine *how* a given model (which may be performing poorly) might be updated. To do so, one needs to
4 understand how the network works so they can determine
5 which aspects of the network to modify and when;
- 6 **2. Training Analysis (TA):** one needs to understand *why*
7 training did not perform as expected; otherwise one does
8 not know what to change next when trying to improve net-
9 work performance. Based on insights from TA, one can
10 modify the design of the network, *e.g.*, change its num-
11 ber of layers, neurons per layer, activation functions, inter-
12 layer connections, or hyperparameters;
- 13 **3. Feature Understanding (FU):** at the highest level, one
14 needs to understand *which* aspects of the input data (sam-
15 ples and/or features) affect the quality of the learning pro-
16 cess. By doing this, a designer may choose to, *e.g.*, in-
17 crease the number of convolutional layers in a DNN so
18 that more powerful feature sets can be discovered. Ad-
19 ditionally, before applying a deep learning solution in a
20 practical application, it is desirable to understand exactly
21 what the model is doing, *i.e.*, which features in the input
22 the model into account when deciding the output label and
23 how the model operates on these features in order to calcu-
24 late such result. Without this understanding, it is difficult
25 to ensure the model is working as desired and users may
26 hesitate in apply it in practice.

28 Note that unless suitable support is provided for the TA, AU,
29 and FU tasks, designing an effective DNN is very much a ‘blind
30 search’ process, which requires many costly iterations, either in
31 an automated form (*e.g.* hyperparameter grid search) or done
32 manually [10], by empirically choosing the model architecture
33 and hyperparameter based on the developer expertise. These
34 tasks can profit from visualization and visual analytics methods
35 in several aspects [34, 32]. For instance, visual tools prove to
36 be an efficient way to understand which features a deep model
37 has learned [51] and which particular neurons are responsible
38 for computing those features [52]. In Section 4, we introduce
39 visual analytics and how it can help deep learning engineering,
40 while in Section 5, we further explain the visual analytics role
41 in terms of the taxonomy we propose and how recent techniques
42 have been tackling variations of the tasks above.

54 multiple perspectives, posing increasingly more targeted ques-
55 tions [56]. Hence, the ability to interactively change visual-
56 izations and pose complex on-the-fly defined queries is key to
57 VA. Important VA techniques target problems related to under-
58 standing high-dimensional datasets represented by feature vec-
59 tors like \mathbf{X}^j introduced in Sec. 2. VA has proven effective in
60 many fields such as software maintenance, health science, e-
61 government, homeland security, and social sciences [57, 58].

62 In the last years, VA has increasingly focused on supporting
63 machine learning applications [59]. VA integrates with ML and
64 DL in terms of proposing specific types of sensemaking loops
65 for specific DL tasks. Note that the deep learning tasks (AU,
66 TA, and FU) are typically executed several times during the it-
67 eration of the sensemaking loop (Fig. 3), as typical in VA work-
68 flows. At each iteration, one obtains additional insights and
69 either change the DNN settings to improve it, or digs deeper
70 into querying the available data to make a decision. Separately,
71 visual tools and techniques that support the three tasks are not
72 disjoint. For instance, to understand training results (TA), one
73 can use a network visualization (AU) that shows the roles of
74 neurons in different layers in computing specific class labels.

75 Visualization of DNNs has caught the attention of the re-
76 search community for many years. Pioneering work in this field
77 dates back to the beginning of the century. Streeter *et al.* pro-
78 posed a technique called NVIS [60], where an artificial neu-
79 ral network is represented as a matrix heatmap that encodes
80 the weights of all neurons n_i^l over all layers $1 \leq l \leq L$. An-
81 other early work, Tzeng *et al.* [61] proposes a node-link graph
82 visualization to show a network’s architecture, coloring each
83 neuron according to the strength of its activation a_i^l for a given
84 selected input. While effective for depicting the structure and
85 operation of small networks, such methods do not effectively
86 scale to treat current-day DNNs that have millions of weights
87 and connections. Following these early developments, several
88 new visualization techniques have been proposed to tackle open
89 challenges in DL, as discussed in Section 5. Visualization ap-
90 proaches have already been helping to answer several questions
91 of deep learning expert, such as, for instance, if the features
92 learned by a model can be interpreted by humans [62] and if
93 deep network learn to recognize different classes in a hierarchi-
94 cal way [63].

5. A Taxonomy on Visual Analytics for Deep Learning

95 In this section we present a taxonomy on how VA tech-
96 niques have been applied to support the proposed tasks (AU,
97 TA, FU) for the three main DL architectures. We first discuss
98 the methodology used to create this taxonomy, followed by sec-
99 tions where we details the goals and sub-tasks on each task and
100 describe how recent publications have been using VA to tackle
101 them and how such techniques are applied to different network
102 types.

5.1. Methodology

104 Deep learning visualization is a relatively new topic that has
105 caught the attention of researchers from both machine learn-
106 ing and visual analytics communities over the past few years.
107

4. Visual Analytics of Deep Learning Networks

44 Visual analytics (VA) has emerged as an extension of informa-
45 tion visualization (infovis) having as aim the analytical rea-
46 soning about problems described by large, complex, and ab-
47 stract datasets using interactive visual interfaces [53, 54, 55].
48 While infovis (usually) aims at visually depicting a dataset with
49 the aims of potentially gaining some insights, VA covers the
50 more involved tasks of formulating, refining, and (in)validating
51 hypotheses about the phenomena that lay behind the data at
52 hand. As such, VA techniques and tools propose a so-called
53 ‘sensemaking loop’ in which designers explore the data from

Therefore, publications in this area are widely spread over proceedings and journals of different domains, imposing a need for a strong methodology when selecting papers for a review like this one. To ensure that we were able to find all of the most relevant publications in the area, we searched for contributions in proceedings and journals of several areas, such as machine learning, visual analytics and computer vision. Particularly, we focused on well regarded proceedings in the mentioned areas such as IEEE VAST, IEEE InfoViS, EuroVis, IEEE Transactions on Visualization and Computer Graphics, ICML, NIPS, ACM SIGKDD, ICCV, and CVPR. In particular, about a quarter of the papers presented at IEEE VAST 2017 focused on visual analytics for deep learning. As the interest in deep learning visualization is recent, we decided to focus on publications released from 2010 onwards, although we mention some earlier works that have historical importance in the field [51, 60]. From all the set of publications retrieved in the mentioned proceedings, we filtered the ones mentioning, in their title or abstract, the deployment of visualization methods to understand or analyse neural network models or features. We also searched for papers with keywords such as *model* and *neural network visualization* on online platforms like *arXiv* and *google scholar*. Finally, we searched through the references of all the filtered publications to find other relevant papers.

From the set of collected papers, we identified how VA techniques have been applied to support the AU, TA, and FU tasks for three of the most popular deep learning architectures in the literature: *deep feedforward networks* (DFNs), *convolutional neural networks* (CNNs), and *recurrent neural networks* (RNNs). Table 1 provides an overview of the papers we surveyed and their respectively addressed tasks and network architectures. For completeness, we mention that other DL architectures exist, e.g., *Autoencoders* [64], *Generative Adversarial Networks* (GANs) [65], *Deep Belief Networks* (DBNs) [66], and *Deep Q-Networks* (DQNs) [67] architectures. For a recent overview of such architectures, we refer to Gibson and Patterson [68]. We did not include such architectures in our survey as we did not find enough papers proposing visualization techniques for addressing particular problems of such architectures [69, 70]. That said, it is worth note that such architectures face many challenges similar to the ones addressed in this survey, making the techniques reviewed here also relevant for the analysis of such deep models. It is also important to note that complex applications may require the use of networks combining elements of two or more type of models—e.g., a network containing both convolutional and recurrent layers. In such cases, VA tools must be adapted to tackle the needs of all the network’s components.

Furthermore, we note that these tasks do not exhaustively cover the applications of VA in deep learning engineering. Visual analytics can (or could) be employed in other tasks, such as (1) training data analysis; (2) performance analysis, and (3) model comparison. For (1), visualizations can help identifying the lack of necessary features or bias in a training set. Unfortunately, the literature still shows a lack of approaches addressing this problem, which makes it an interesting topic for future research. For (2), VA can provide powerful tools to analyze the

performance of deep models by providing ways to compare the confusion between multiple classes [71] or to depict the distribution of an output value over the set of input features [72]. For (3), VA can answer why a model performs better than another [73]. This also brings difficult challenges, mainly because it is hard to compare models that do not share the same structure (for instance, number of layers and neurons per layer) since such models end up having completely distinct weights, which make them recognize features in different ways. However, comparing networks with similar structure but different hyperparameters [74, 73, 75] can be an effective way to understand how the parameters affect the final performance of the model.

Training Analysis (TA): The training of a deep model is a hidden process that gives little to no insight to the designer about what is happening. If a model is not performing well, it is very hard for experts to identify what should be changed in the training parameters or even whether there is a problem in the training process at all. As understanding the training process is still an open challenge in deep learning, visual analytics can play a powerful role in addressing it. By visualizing the *evolution of model metrics* during the training process, VA techniques can help to understand how the model achieved some performance and to identify undesirable behaviors. For instance, visualizing metrics about gradient values can help to understand how the updating process is changing the parameters of neurons and hidden units [76, 77], allowing a designer to find out network parts that are not stabilizing or that are not being sufficiently changed by backpropagation. Additionally, analyzing how the neurons’ weights and activations change through the training epochs are key to comprehend how the DNN evolved and how it learned to recognize the relevant features for the relevant task at hand [52]. Visual analytics also uncovers new possibilities for training DNNs, as it can help designers to *analyze the training process in real-time*, allowing the designer to make assumptions about the model, and take corrective decisions, without having to wait for the entire training to finish [78, 10].

Feature Understanding (FU): Although a neural network may be performing well on unseen data, it is not clear when and why this happens. Hence, the machine learning community has put much effort in approaches to figure out which features the input data must have in order to produce the desired output [6, 79] and how the model use these features to compute its label prediction [10, 52]. However, learned features in DNNs are only *implicitly* described in terms of the huge number of parameters Θ of the model, in contrast to the *explicit*, hand-engineered, features used by classical ML techniques (see Sec. 2). As such, VA aims to explain in an interpretable (intuitive) way how the information spread over all neurons in all layers of a DNN captures these features [73]. One sub-task here is to *interpret the model*, i.e. to show how features learned in earlier (closer to input) layers get merged in subsequent layers to identify more complex patterns [10] and how intermediate layers transform the input they receive [52]. Another goal of the FU task is to *explain learned features*, i.e. to identify in a particular input or set of inputs, which of their features were taken into account in order to decide the output label. By understanding what the network is recognizing, experts can give more reliability to their mod-

els, as they can know what they are recognizing or predicting with more certainty [52]. This also allows identifying patterns and features the model has not learned to recognize, but which a human may consider important [76]. Techniques for this task can be classified into *instance-based* and *feature-based* visualizations. The former ones depict the behavior of the model for specific input instances, be it a single one or a subset of many. The main goal of such visualizations is to find which features of the input produce high activations in the network and in which neurons or layers that occur [80]. Feature-based techniques, on the other hand, aim to explain which features an input must have to produce a particular output in the final layer or, more generally, in any layer [52]. Such techniques are well suited when instances are not easily interpretable or in applications with many possible outputs, where analyzing individual inputs can become tedious [80]. As well as for the TA and AU tasks, scalability is also a problem for FU for more complex models, especially when visualizing the activations of many neurons at once [10]. Additionally, some applications may have inputs composed of different data formats. For instance, to classify a social media post, the model could have as input an image, a text, and information on the user who posted it [80]; in such cases, visualizing important features in a meaningful way is harder than when X consists of a single data type.

5.2. Architecture Understanding

As explained in Section 2, modern DNNs may have hundreds of layers and hundreds of thousands of neurons. When using such wide and deep models, it is easy to lose track of all the aspects of their architectures or which computations they do at each unit of the model. Thus, visual analytics can play a key role in helping designers have a better insight about the characteristics and behaviors of their models during the development pipeline.

The main goal of visualizing the architecture of a DNN is to give a good understanding of both high and low-level aspects of the model [81]. At a very high level, showing the network topology (as a graph) quickly tells designers the overall structure of the network, *e.g.*, how many layers L it has, how their sizes s_i vary, and which kind of operation they perform. We call this sub-task *archicteture visualization*. This helps to understand a DNN much in the same way that architectural diagrams reverse-engineered from source code help software maintenance [104]. At finer levels, visualizing the connections between neurons on consecutive layers—encoded as the weights θ_l of a given layer l —may help understanding how simple features get merged into more abstract ones in the classification process [10]. Additionally, one can visualize the combination of DNN structure-and-data, by annotating the DNN graph with weight vectors, activations, and training statistics. This helps understanding how structure correlates with behavior [10, 73]. Showing this helps to find possible inefficiencies of the model, such as inert neurons, layers that are not relevantly contributing to the final prediction [76], or redundant components that recognize the same features [52]. In other words, such techniques allow experts to *validate the architecture* of the model. Given the high number of layers and neurons that current DNNs might achieve,

scalability is an important concern for AU visualizations, as VA techniques that work well for simple models may not be able to handle larger and/or more complex ones [80]. To handle this issue, several approaches were proposed, such as clustering neurons with similar behavior—computed from their activation responses— [10] or hiding parts of the architecture that are not relevant to the analysis [81].

Deep learning architectures are essentially directed acyclic graphs (DAGs) where nodes represent neurons and edges represent connections between subsequent layers [10]. For this reason, *graph visualization* has been a straightforward way to visualize the architecture of deep models [80, 10, 81, 83, 82]. Such visualizations help deep learning engineers in multiple ways. First, they provide an overview of the operations that the model is performing when an input flows through the network [80]. As a refinement to this, showing what the network does in different layers and parts thereof helps the engineer understand whether the chosen architecture is appropriate and, if not, where it should be adapted or modified [81]. We call the graph visualization of a model architecture *topology visualization*.

For such a graph visualization to be effective, it is not enough to depict only the connections between neurons, as this information is typically already known by the architect and is the same in all layers (*i.e.*, drawing the individual neuron connections with no associated value does not bring additional insights). Graph visualizations become effective when they show additional data on the *activity* of the neurons. One way to do this is to show the neuron weights, *e.g.* by color coding. However, weights are hard to interpret, particularly in deeper layers. A more insightful design is to show neuron activations for specific inputs [10]. For this reason, many existing works show the activation vectors \mathbf{a}^l produced for one or more inputs via color-coded matrices or vectors [10, 70].

As DNNs become deeper and wider to tackle more complex applications, scalability becomes an important issue for visualizing the network’s graph structure. One solution for this is to visually cluster neurons having similar activations [10] and next use *edge bundling* to visually group connections linking neurons in the same clusters [10, 81]. By clustering groups of neurons with similar activations, architecture visualization can be made clearer. Also, this approach can highlight large groups of neurons that may be involved in correctly predicting a particular label, as well as to highlight classes that are not being sufficiently learned by the model (*e.g.*, few neurons respond to inputs of that class) [10]. Edge bundling has proven very effective to trade clutter for overdraw when creating simplified views of graphs of millions of edges [105, 106] and hence it has the required scalability for handling very large DFNs. Another approach proposed to improve the scalability of graph visualizations is the omission of non-critical operations (*e.g.*, *pooling* layers, which implement a type of pre-processing on the outputs of a given layer using fixed operations that are not updated by the training process, and thus sometimes may be omitted in a visual analysis). Finally, to improve scalability, it is also possible to highlight particular network regions with similar properties—*e.g.*, parts of a DNN with similar weights

Technique	Taxonomy				Networks		
	Architecture Understanding	Training Analysis	Feature Understanding	DFN	CNN	RNN	
CNNVis [10]	AVis and AVal		FE	●	●		
Activis [80]	AVis		FE and PV	●	●		
TensorFlow Graph Visualizer [81]	AVis			●	●		
TensorFlow Playground [82]	AVis and AVal	RTA	MI	●			
DGMTracker [70]	AVis and AVal	MME	FE	●	●	●	
ReVACNN [83]	AVis	RTA	FE	●	●		
Harley [84]	AVis		FE	●	●		
BIDViz [78]		RTA and MME		●	●	●	
DeepEyes [76]	AVal		MI and FE	●	●		
Rauber et al. [52]		MME	MI	●			
Deep View [85]		MME	FE	●	●		
Grad-CAM [86]			FE		●		
RNNbow [77]		MME					●
Zahavy et al. [69]			MI	●	●		
Yosinski et al. [62]	AVal		MI and FE		●		
CNNComparator [74]	AVal and MC			●	●	●	
Alsallakh et al. [63]			MI	●	●	●	
Zeiler and Fergus [6]			FE		●		
Nguyen et al. 2016-1 [87]			MI		●		
Nguyen et al. 2016-2 [88]			MI		●		
Aubry and Russell [89]			MI		●		
Simonyan et al. [90]			MI and FE		●		
Wei et al. [91]			MI and FE		●		
Mahendran and Vedaldi 2015 [92]			FE		●		
Mahendran and Vedaldi 2016 [93]			FE		●		
Zintgraf et al. 2016 [94]			FE		●		
Dosovitskiy and Brox [95]			FE		●		
Erhan et al. [51]			MI		●		
Samek et al. [12]			FE		●	●	
Montavon et al. [28]			FE		●	●	
FeatureVis [33]			FE		●		
Zintgraf et al. 2017 [96]			FE		●		
Heyi Li et al [97]			FE		●		
VisualBackProp [98]			FE		●		
RNNVis [73]			MI and FE				●
LSTMVis [99]			MI and FE				●
Jiwei Li et al. [100]			MI and FE				●
LAMVI [101]			FE				●
Ding et al. [102]			FE				●
Karpathy et al. [103]			FE				●

Table 1: Taxonomy of VA-related publications related to different DNN architectures and engineering tasks. Tasks were further refined as follows: Architecture Understanding: architecture visualization (AVis), architecture validation (AVal) and model comparison (MC); Training Analysis: real-time analysis (RTA) and model metrics evolution (MME); and finally Feature Understanding: model interpretability (MI), feature explainability (FE) and performance validation (PV).

and activations [81]. Graph visualizations are strongly aided by interactivity which can help designers focus on and explore in more detail particular parts of the model they deem more interesting [81].

Visualizing the graph structure is not the only way to analyse the architecture of a DNN though. Other kinds of visualizations can also be helpful in give insights if the chosen architecture is the right one, without explicitly visualizing the model’s graph workflow. For instance, in DeepEyes [76], the authors propose a visualization called *perplexity histogram*, a histogram encoding how the response of the layer for particular inputs are changing over time. With this technique, they

are able to identify layers that are either stable, improving or decreasing their recognition capability. Another important goal of architecture understanding is to verify if the chosen structure is efficient, particular related to the number and order of layers and setting of hyperparameters. This is what we call *topology validation*. Several works have used activation heatmaps in order to identify units or layers that are not being helpful for the network, either because they are not producing significant high activations for any kind of input—and thus can be dropped from the architecture—or because they are activating too often for too different inputs—what may indicate the need of more units [76, 10, 70].

1 Additionally, another sub-task worth mentioning is *model
2 comparison*. In many context, it would be useful for experts to
3 compare different models to understand, for instance, why one
4 performs better than other in some application. This is partic-
5 ular useful when comparing two model with same architecture
6 but different hyperparameters. In CNNComparator [74], the au-
7 thors do that by comparing the difference in learnable parame-
8 ters after training using heatmaps and histograms. Additionally
9 in this area, there is a whole field in VA dedicated to under-
10 stand how a model is affected by the tuning of hyperparameters
11 called *Visual Parameter Space Analysis* (VPSA) [75].

12 **Deep Feedforward Networks:** Visualizing the architecture of
13 traditional DFNs is pretty straightforward. In such networks,
14 each neuron in a particular layer i receives as input the out-
15 put activation of all neurons in the previous layer $i - 1$. Also,
16 each individual connections (i.e., between a particular neuron
17 in $i - 1$ and another particular neuron in i) is associated to a
18 single weight or activation value. For this reason, it is easier to
19 represent and visualize such values using color encodings and
20 heatmaps, either on the edges [82, 83] or on the nodes [70, 84].

21 **Convolutional Neural Networks:** Just as DFNs, CNNs can
22 also be seen as directed acyclic graphs. However, CNNs have
23 what it is called *convolutional layers*. In such layers, neu-
24 rons perform convolutional operations in the input data received
25 from the previous layer. This brings some differences for the
26 analysis of the architecture of a CNN if compared to DFNs.
27 First of all, the analysis must take into account that a single
28 weight in a convolutional neuron is applied to not only one but
29 several input values. For instance, if an input image is sent
30 to a convolutional layer, a given weight in the layer will op-
31 erate all over the image values, and not in just a single pixel
32 as it happens in DFNs. Additionally, visualizing the activation
33 of CNNs layers is more challenging as each unit produces a
34 multidimensional output activation, and not a single value as
35 DFNs do — e.g., each unit in the convolutional layer will pro-
36 duce an *activation map* with similar dimensions as the input,
37 while DFN units produce a single unidimensional value. Such
38 particularities must be taken into account when analysing the
39 architecture of CNN models. Nonetheless, graph visualization
40 still is a natural way to visualize CNNs [10, 83, 84, 70]. A typi-
41 cal approach when visualizing convolutional layers is to use the
42 nodes of the graph to display either the convolutional filter of
43 the unit (i.e., the unit weights in the exact order they are applied
44 on the input [74], the activation map produced by a particular
45 input [83, 84, 70], or inputs that produce strong activations on
46 that neuron [10]. Figure 4 shows an example of CNN architec-
47 ture visualization [10]. Here, neurons with similar activations
48 are clustered together, and only the features that produce the
49 strongest activations on them are displayed, to limit visual clutter.

51 **Recurrent Neural Networks:** The visualization of the struc-
52 ture of a RNN is a difficult and, to the best of our knowledge,
53 still unexplored topic. One of the issues to overcome when
54 building architectural visualizations for RNNs is a large num-
55 ber of possible output formats. In some applications, such as
56 sentiment analysis, the RNN processes the whole sequence of

57 inputs and then returns a single output [107]. However, in other
58 applications such as machine translation, the network must out-
59 put a new value at each element of the input stream that is pro-
60 cessed [108]. Other variations, usually used for sequence pre-
61 diction tasks, aim to predict an output equal to the following
62 element of the input stream every time an input unit is pro-
63 cessed [1]. Additionally, RNNs have a recursive structure in
64 which input elements are processed in a sequential manner.
65 When one element is processed, the *hidden state* of the unit
66 is modified and this can, and probably will, modify the behav-
67 ior of the unit for future input elements. When visualizing the
68 architecture of RNN models, analystis should be aware of this
69 particularities, as understanding the how the hidden state is be-
70 ing affected by input and how its affecting the outputs is as im-
71 portant as understanding the structure of input weights and out-
72 put activations. All in all, visualizations supporting architecture
73 understanding for RNNs are weakly developed and, given the
74 complexity of these architectures, we consider this a promising
75 future research topic.

5.3. Training Analysis

76 As most machine learning techniques, neural networks are
77 trained via gradient-based methods, such as gradient descent,
78 that minimize a cost function C (see Sec. 2). Since DFNs have
79 multiple layers, the traditional gradient-based method used (as
80 described in earlier sections) is the backpropagation algorithm,
81 which updates all weights in the network, starting from the ones
82 in final layer and moving towards shallower layers of the net-
83 work, in order to minimize the prediction error for inputs in the
84 training set [1].

85 We consider two sub-tasks in the analysis of training pro-
86 cess: *model metric evolution* and *real-time analysis*. In the for-
87 mer, the expert analyses how a particular metric evolves over
88 the training of the network. Such metrics are user-defined and
89 can contain information helping to understand, for instance, if
90 the gradient updating is being performed as expected [77], if the
91 model is improving its performance [85] or if particular units or
92 layers are indeed converging to some learning [76]. Real-time
93 analysis, on the other hand, is the analysis of anything related to
94 the network—architecture, learned features, metrics—in real-
95 time, allowing the designer to recognize problems and imper-
96 fections as soon as possible, and thus decreasing the time ex-
97 pended in the training of the model. The two sub-tasks are,
98 of course, not exclusive, as evolutionary metrics are often a
99 good option to be displayed in real-time while the training pro-
100 gresses [76].

101 The traditional way to visualize the training process is to plot
102 the accuracy of the model for the training or test set through the
103 training epochs. While this gives a good intuition on how accu-
104 racy changes with training, it does not provide insights on how
105 to improve it apart from more training. For instance, informa-
106 tion about how the weights are being updated on individual lay-
107 ers or neurons along the training is not provided. Qi *et al.* [78]
108 alleviate this by allowing designers to define and plot their own
109 metrics in real-time, allowing a more effective guiding of the
110 training process. As the learning is not uniform in all layers,
111 developers may want to identify layers or neurons that are not

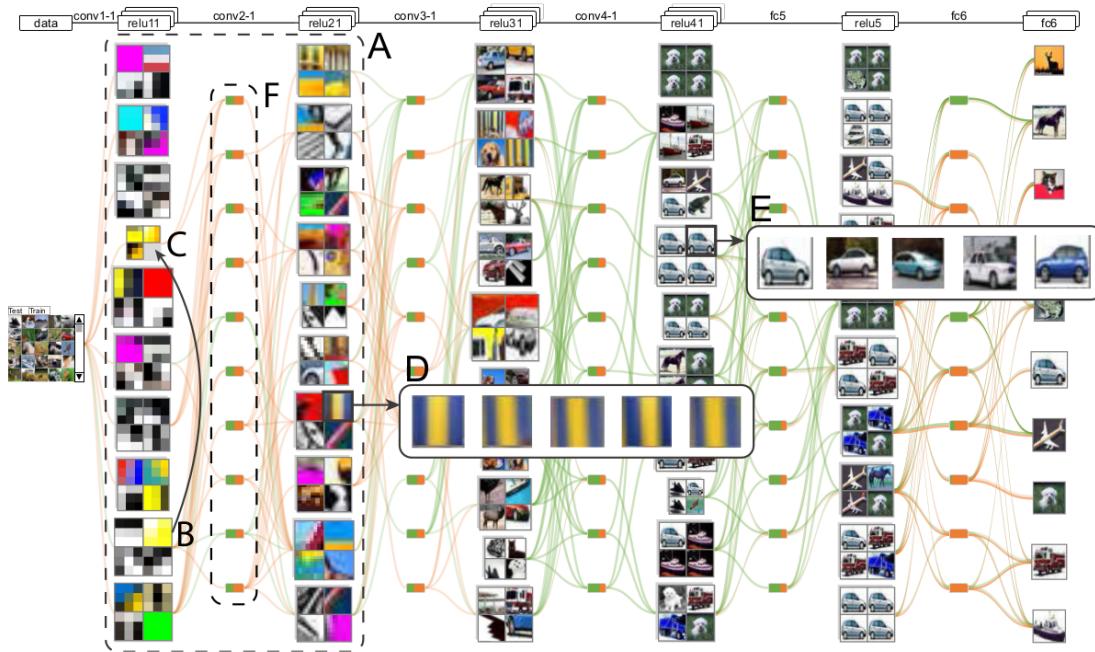


Fig. 4: CNNVis Tool [10]: the structure of the network is shown as a directed acyclic graph where nodes display learned features for clusters of similar neurons.

being well trained. To support this, Pezzotti *et al.* [76] propose so-called perplexity histograms, a visual technique to find stable layers, *i.e.*, the ones that already stopped receiving relevant weight updates. By visualizing the progress of individual layers during training, this helps to identify if a given layer’s weights converge to a good solution or if the layer is not learning to recognize any useful patterns in the input data. An alternative approach, proposed by Zhong *et al.* [85], uses heatmaps to depict how neurons weights and activations change over the training process. For this, they propose two metrics: discriminability—measuring how different is the average activation produced in a layer by elements of a particular class from elements of any other class—; and density—which evaluates the quantity of higher activations a particular neuron produces.

Deep Feedforward Networks: Most methods to visualize the training process of DFNs focus on displaying how a particular metric evolves over time by a given layer or neuron. Focusing on a single element of the model rather than in the whole network is effective because then the designer can get more information about where the model is underperforming and how to modify it. One recent approach aiming at this is proposed by Rauber *et al.* [52]. For each sample $\mathbf{x}_i \in D_{\text{train}}$ in the training set, they consider the activation vector $\mathbf{a}(\mathbf{x}_i)^{L-1}$ of the last hidden layer. Next, these vectors are projected from $\mathbb{R}^{s_{L-1}}$ —where s_{L-1} is the number of neurons in layer $L - 1$ —to \mathbb{R}^2 using standard dimensionality reduction (DR) methods such as t-SNE [109], yielding a 2D scatterplot of points \mathbf{p}_i , one per sample \mathbf{x}_i . Key to DR methods is their ability to place points with similar high-dimensional vectors close to each other in 2D, and points with dissimilar vectors far apart in 2D, respectively—thus showing how similar are the high-dimensional vectors. This process is repeated for all training epochs, yielding a 2D trajectory of points per sample \mathbf{x}_i . These trajectories are then

colored according to the classes \mathbf{y}_i of the corresponding training samples, and bundled to yield a simplified, though suggestive, view of how the neurons of the last hidden layer get increasingly more class-specialized (farther apart in the 2D projection) as training progresses. Figure 5 (middle) shows an example hereof. The same type of visualization can be used to show how the network layers learn to discriminate between the different classes (Fig. 5 (right)). Here, each trail represents the projection of activations of all hidden layers $2, \dots, L - 1$ of a test sample \mathbf{x}_i , after training. As in the previous image, one can see how same-class images yield increasingly more similar activations as the data flows deeper through the network.

Convolutional Neural Networks: The training process of a CNN is typically done via backpropagation, in a very similar way to the training of DFNs. As such, most of the visual approaches used to analyze the training of DFNs can also be used to CNNs. However, the differences between convolutional neurons and fully-connected ones have an impact on the training process. First, since the former applies convolutional operations on their inputs, weights act on a set of *neighboring* inputs—sometimes in two or more dimensions. Secondly, such weights are shared across the input, *i.e.*, the same weight parameters are applied to multiple regions of the input. Because of these impacts, designers must take into account, for instance, that a neuron activation may be affected by changes in multiple parts of the input.

Liu *et al.* [70] propose a specific visualization for CNN training—particularly for generative models—that displays how features are learned through the training process by plotting line charts showing various designer-selected statistics of interest (activations, gradients or weight updates) over time. If the designer spots a layer with an interesting or abnormal behavior—such as an abrupt change of many activations in a

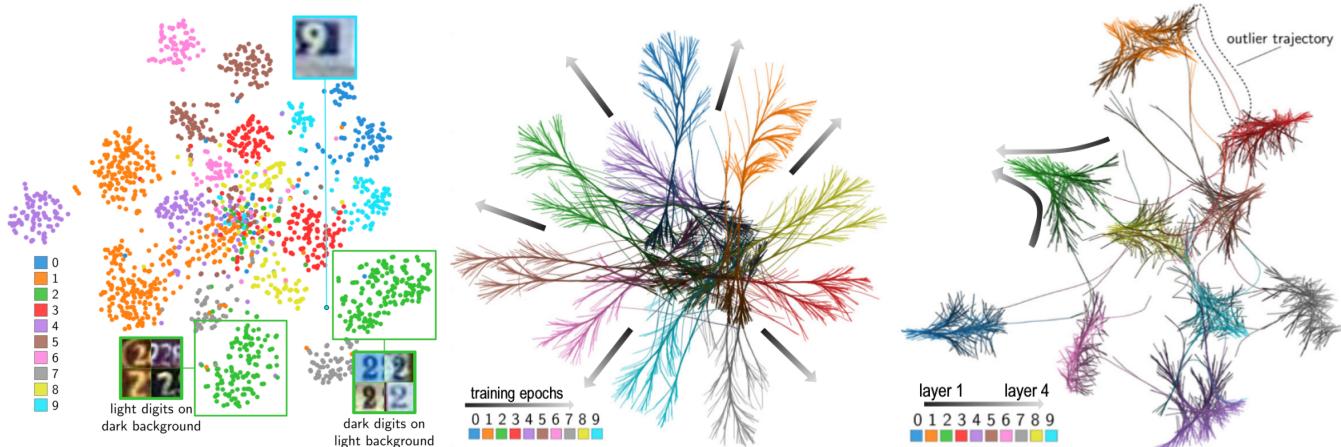


Fig. 5: By projecting activation vectors onto a bidimensional space, Rauber *et al.* [52] are able to explore the learned features (left), how they evolve over the training epochs (middle), and how they are identified by different layers (right).

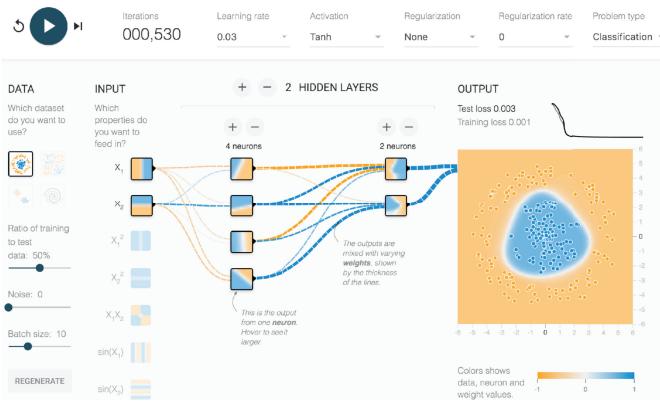


Fig. 6: TensorFlow Playground tool [10] showing jointly the structure and activations of a simple DFN, and allowing interactive control and monitoring of the training.

single epoch—, they can explore the activations of that layer’s neurons in more details and visualize the subsets of input data that lead to such activations.

Recurrent Neural Networks: So far, the visual analysis of the training process of RNNs is a topic that still has not caught large attention from the research community. The only contribution we found in this direction was proposed by Cashman *et al.* [77] aiming at understanding the so-called gradient flow during the training process. Specifically, they address the problem of how gradient values change during the backpropagation phase of training. Understanding how these weights change throughout the process helps by providing insights to the designer regarding how new training examples modify the learned model; this, in turn, may help in determining how to improve the model hyperparameters, if necessary. Two common problems that can occur during backpropagation process are the vanishing gradient and the exploding gradient [110] problem. In the former, the gradient becomes insignificantly small very quickly, making the training process unable to relevantly change the weights of layers far away from the output layer. In the latter case, the gradient keeps an exaggeratedly large value for many layers,

changing the parameters Θ so drastically that the model never stabilizes. To understand this type of phenomena related to the gradient flow, Cashman *et al.* propose a stacked bar chart visualization (Fig. 11). Each stacked bar represents the magnitude of the gradient that produced the model’s weight update at a single time step. Each partition of a bar represents how far in time each part of the gradient came from. Using this visualization, the authors were able to effectively identify cases of vanishing gradients, where large gradient partitions in one stack bar quickly become very narrow partitions in the next time steps.

5.4. Feature Understanding

In contrast to feature engineering methods, DL methods do not have an explicit representation of the data features they use (Sec. 2). The representation of such features is actually ‘scattered’ in the parameter weights \mathbf{P} of all neurons over all layers. Unfortunately, because DNNs are basically a composition of nonlinear functions, it is very hard to retrieve any interpretable information from the weight values, particularly the ones in deeper layers. An alternative way to understand the learned features is to analyze the *activations* produced by these weights when a given input is fed into the network. This can be done at several levels, as follows. By visualizing the activations vectors produced by a *single* input over the *entire* model, using *e.g.* heatmaps or matrix plots, one can understand how the sample information flows through the network until an output is obtained [80]. Alternatively, by visualizing the activation vectors of multiple input samples for a single layer or even single neuron, one can find which patterns the layer, or neuron, is learning to recognize [52]. Yet another approach is to visualize how the network divides the input data space X at each neuron [82] (see also Fig. 6). This approach is very intuitive to use, as it effectively shows how each neuron u_i classifies every possible sample $\mathbf{x} \in X$, by color-coding a 2D plot of X with the respective activations $a_i(\mathbf{x})$, and also allows real-time changes of the architecture and hyperparameters and training monitoring. However, it only works for data spaces $X \subset \mathbb{R}^2$ and relatively small networks (tens of neurons).

The aforementioned heatmap technique is probably one of the most widespread approaches for feature understanding. Here, activations are displayed as a matrix where rows are input samples or sample classes and columns are the neurons of the DNN,[80, 76]. The colors of the matrix cells encode the activation produced by each input or class (row) to each neuron of the model. When visualizing activations for particular inputs, heatmap matrices work as *instance-based* techniques that show which parts of the neural model learn features for that input[76]. Conversely, when displaying activations for a single class or subset of inputs—done usually by computing the average activation of all considered input samples—heatmap matrices behave as *feature-based* techniques, showing which network parts specialize in recognizing that class or subset[80]. However, one should note that such visualizations can be misleading when inputs belonging to the same class present distinct input features—for instance, in an image classification task, images of cats may have very different background colors but need to be assigned the same label. Even in a network with high performance on this task, neuron activations in hidden layers may be significantly different for samples of this type (qualitatively different input values but same associate label)[88], and therefore different activation patterns should not necessarily be interpreted as evidence that the network did not successfully learn appropriate intermediate features. An example of such a heatmap matrix produced by Kahng *et al.* [80] is shown in Fig. 7 (B).

As outlined in Sec. 2, the activations $\mathbf{a}^l(\mathbf{x})$ produced by a hidden layer l of a DNN for an input sample \mathbf{x} form a high-dimensional vector that captures the different features of \mathbf{x} . Comparing such vectors for all samples \mathbf{x} in a training or test set allows one to visualize how the network succeeds (or not) in discriminating between these. To do this, dimensionality reduction methods can be used to create 2D scatterplots of these samples. Close points in these scatterplots indicate samples which are found to be similar to the network [80, 76, 52, 69]. In particular, by projecting activations from multiple inputs, designers can identify instances wrongly predicted and build hypotheses of why this happened by comparing them with inputs with similar activations [52]. Figure 5(left) shows such a projection for the well-known SVHN image dataset [111], with points colored by the class label, based on the last hidden layer activations after training [52]. This image allows a designer to quickly observe that data points belonging to each class are divided into two compact clusters—one which represents light digits on a dark background, and one which represents dark digits on a light background. Hence, this visualization lets one see that the network has learned *irrelevant* information of digit-vs-background contrast, which is not useful when classifying the digit images—and this is a finding that can help fine-tuning the network to learn more effectively [52]. The main downside of dimensionality reduction techniques is the intrinsic information loss related to the compression of a high dimensional data to a bidimensional representation. Also, different dimensionality reduction methods—or the same one but with different parameter tuning—can give different and equally useful projections, as datasets may have multiple features that cannot be uncov-

ered in a single projection. To overcome this limitation, several authors have proposed interactive dimensionality reduction approaches [112]. These methods can be worth when visualizing very high dimensional data such as layer activations and weight vectors.

Deep Feedforward Networks: In DFN models, we can interpret the fully-connected layers as operations performing non-linear transformations in the multidimensional space the training data lies. For this reason, techniques aiming to visualize high-dimensional spaces such as dimensionality reduction are often good alternatives to analyse DFNs [52]. However, even though such approach helps to interpret the function of the layer as a whole, it brings little insight about individual units. Techniques that show how neurons respond to different kind of inputs, such as heatmaps [10] are more desirable in such cases, as they give a clear understanding of the role that neuron is playing in the recognition task.

Convolutional Neural Networks: Visualizing activations is as important to understand CNNs as it is to understand traditional DFNs. Particularly, heatmap matrices have stood out as an effective tool for this task [84, 10, 62, 63], as they can compactly display the activations of neurons for multiple input samples [10] or in multiple convolutional filters [62]. By using a heatmap matrix to compare activation vectors for different inputs, Alsallakh *et al.* [63] showed that this type of analysis can lead to non-trivial conclusions—such as the identification of a hierarchical structure in the classes present in a training set, given that classes that share similar features usually have similar activations in shallower layers and more distinct activations in deeper ones.

DR projections have also been used in CNNs to compare activation vectors of different inputs [83]. For instance, Aubry *et al.* [89] compare the activation vectors of slightly modified images to understand the manifold created by these modifications in the input (image) space. Nguyen *et al.* [88] use DR projections of images belonging to the same class in a training set to identify image clusters that are characterized by distinct types of features. This allows identifying different styles, or kinds, of images that must be assigned to the same label by the model, which can be a challenging task for improperly trained models.

Since CNNs are often used in applications that take images as inputs, heatmaps have also been used to show which regions of the input image produce the strongest activations in a given convolutional filter. For instance, Yosinski *et al.* [62] depict the learned filters of convolutional layers as images with the same resolution as the input image, with pixels colored based on to how they contribute to the activations produced by that filter. Zeng *et al.* [74] also use heatmaps to visualize the differences in the weights of CNNs with the same architecture but trained with distinct hyperparameters, aiming at understanding how they affect the model’s performance. Many variations of such image-based techniques have been proposed for the tasks of CNN training analysis and feature understanding. As such, we discuss them in greater detail next.

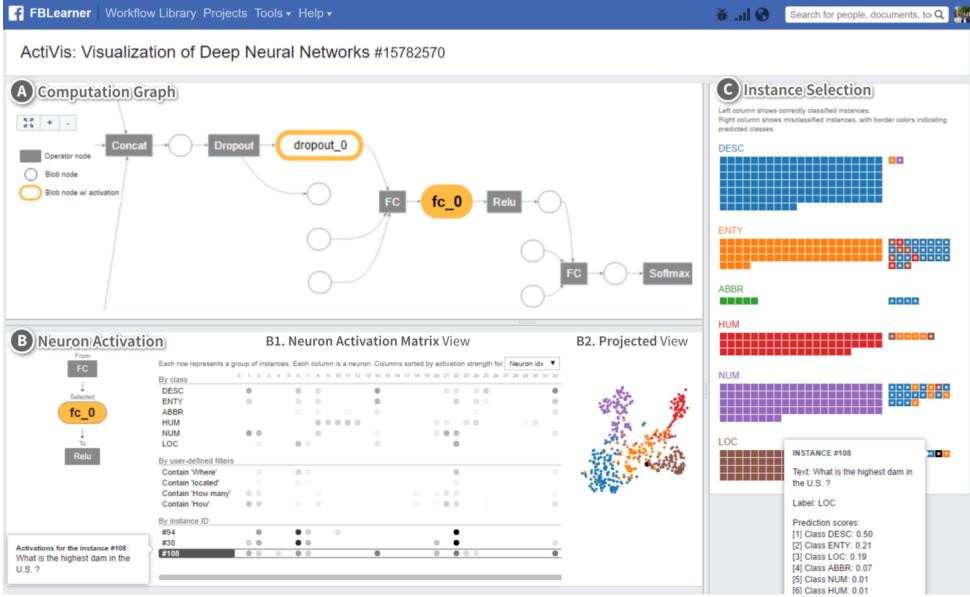


Fig. 7: The Activis tool [80] allows the designer to: (a) visualize the structure of the model; (b) explore the learned features; and c) analyze its performance.

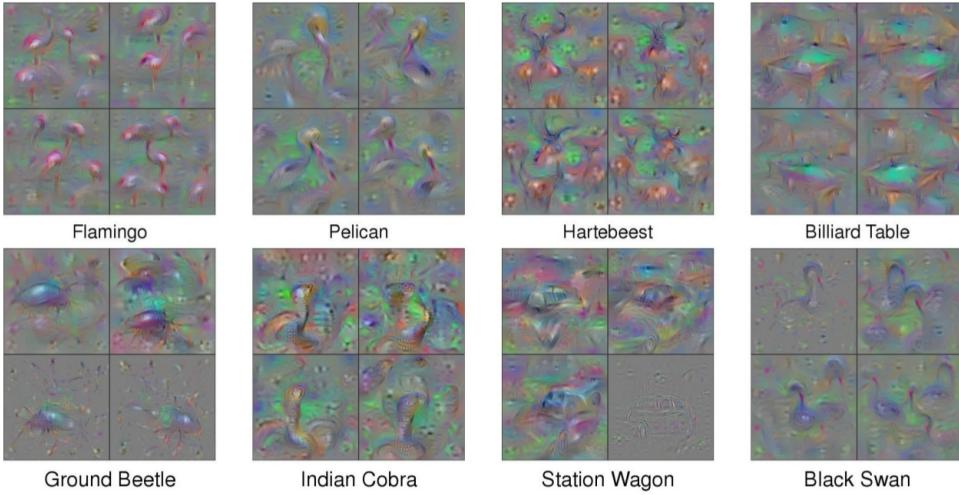


Fig. 8: Activation maximization algorithm applied to eight different classes of the ImageNet dataset [5]. Image adapted from Yosinski *et al.* [62].

5.4.1. Image-based techniques

Given that CNNs are particularly useful for applications that use images as input, many methods have proposed to visually understand them in the image domain. These methods can be divided into two main types: *instance-based* and *feature-based* methods. Methods from both classes are reviewed and discussed next.

5.4.2. Instance-based techniques

Although deep CNN models can achieve high accuracy in image classification tasks, it is not trivial to figure out which features an image should have to be assigned to some particular class. Understanding this may help in identifying mistakes made by the model, such as considering a recurrent background object or feature as intrinsically part of the class. Figure 9 shows two examples hereof. Image (a) shows a sample used as input for Google's Inception neural network [8]. Image (b)

shows which parts of this input image have been relevant for detecting the class 'Labrador' in the input. We see that such parts contain both relevant information (the dog's face) but also spurious pixels (the bottom part of the character's shirt). Image (c) shows the image of a dog that was incorrectly classified as 'Wolf' using the same neural network [13]. Image (d) shows that this incorrect classification relied solely on the presence of a particular background. Using this insight, the designer of the network discovered that wolf images in the training set all had snow as a background. We call such techniques *instance-based* visualizations, as they emphasize parts of input instances responsible for a high activation vector \mathbf{a}^l in some layer l (see also Fig. 10a).

Several examples of such instance-based visualizations exist. Simonyan *et al.* [90] rank pixels in the input images by how much they contribute to a particular class assignment. Montavon *et al.* [113] backpropagate the activation from deeper lay-

17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33

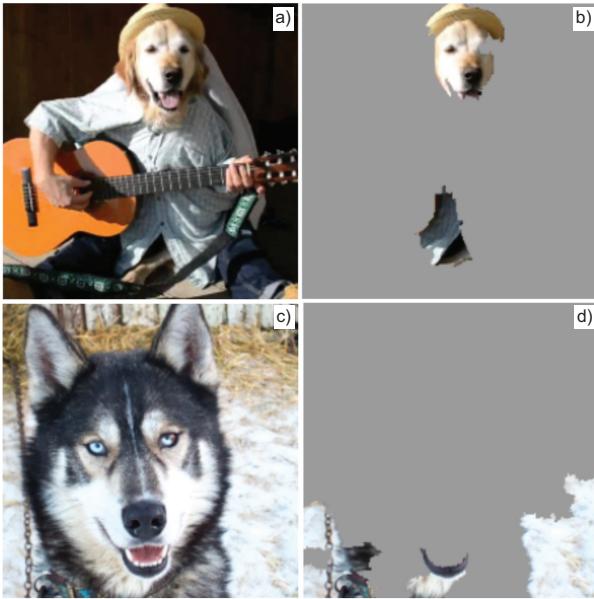


Fig. 9: Image-based visualizations explaining the input elements that lead to a classification outcome. See Sec. 5.4.1.

ers to identify relevant pixels in the input. Zintgraf *et al.* [94, 96] iteratively remove different patches of the input image and check whether the model is still capable of recognizing the correct class. More recent works [98, 86] analyze the weighted gradient in the last convolutional layer to understand how information is flowing through the model. Li *et al.* use a two-step algorithm based on *Layer-wise Relevance Propagation* (LRP) [114] to recognize the more relevant pixels to the activation. The results of such approaches are usually displayed as so-called saliency maps [90], *i.e.*, heatmaps where pixels are colored based on their relevance to the classifier’s output.

A separate challenge of training deep models is to discover how much of the input information has been lost over the layers. This can be done by measuring whether (and how much) it is possible to reconstruct an image from its activation vector produced in a given layer of the network. As a neural network is a complex composition of nonlinear functions, it is likely that a perfect inverse transformation (from activations to the input) is not possible, especially from the deeper layers. Yet, some of the image’s most relevant features may be reconstructed in this way. The main technique to approach this problem is called *code inversion* [92]. To recover a synthetic image $\tilde{\mathbf{x}}$ from the activation vector \mathbf{a}' in a deep network layer l caused by an actual input image \mathbf{x} , gradient descent is used to synthesize an image $\tilde{\mathbf{x}}$ that generates an activation vector as close as possible to the one generated by \mathbf{x} . Mahendran *et al.* [92] found that activations of lower layers can reconstruct the input image \mathbf{x} more faithfully than those from deeper layers. This supports the hypothesis that deeper layers learn more abstract (less detailed) representations of the input. In some networks, even deep layers can preserve image-specific features like object position and colors [95]. Mahendran *et al.* [93] present a comprehensive study measuring the quality of images generated by code inversion according to criteria such as reconstruction similarity, natural-

ness, interpretability, and classification consistency. Other approaches use deconvolutional or up-convolutional networks to synthesize such approximate images [95, 6]. Deconvolutional networks try to produce the inverse operations performed by a CNN and recover an image from an activation vector. By contrast, up-convolutional networks are conventionally trained CNNs that learn to predict images from the activation vectors produced by the CNN under analysis.

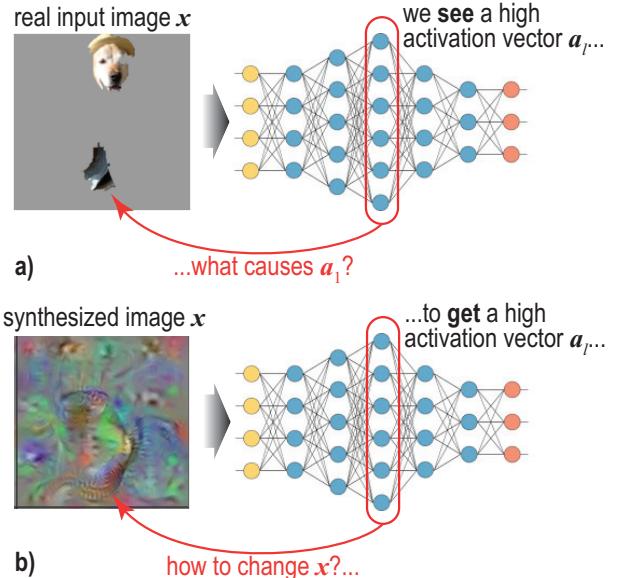


Fig. 10: Instance-based *vs* feature-based visualizations for feature understanding. See Sec. 5.4.1.

Another instance-based technique related to image synthesis is *caricaturization* [93]. In this approach, an input image is modified to sharpen features that cause high activations in some given neurons of a CNN model. This is usually done by an optimization algorithm similar to the one used in activation maximization techniques (see Sec. 5.4.3). However, the aim here is to exaggerate the most relevant features of real input, thereby allowing the designer to find out what the model is learning from that input. Such features were found to predominate in the final (deepest) layers of the CNN model.

An important drawback of instance-based techniques is that they produce insights that are valid only for a single input \mathbf{x} . This leaves some important questions unanswered. It is usually not clear if all the instances $\{\mathbf{x}_i\}$ belonging to some class \mathbf{y} must contain the features highlighted by the explanation of a single such image \mathbf{x}_i —for instance, it is not evident from the examples shown in Fig. 9 that all images that will be labeled as *wolf* will be so due to the presence of background snow. To reach such a conclusion, one needs to manually browse through the explanations of multiple inputs \mathbf{x}_i belonging to the class *wolf*.

5.4.3. Feature-based techniques

Whereas instance-based techniques reflect activation information associated with different inputs, feature-based techniques work in the opposite direction (Fig. 10b): they create synthetic input images \mathbf{x} that cause high activation vectors \mathbf{a}' in some given layer l , or in a neuron u'_i thereof. When applied to

1 a neuron on the network’s last layer L , this technique produces
 2 images that capture the features that the network deems relevant
 3 for the class that that neuron is responsible for [90]. This way,
 4 one can find which generic features the network searches over
 5 the input image space to predict specific labels. This approach
 6 can also be applied to hidden layers to give insights on which
 7 features these layers are looking for when predicting a class [51, 62].

9 As outlined in Sec. 5.4.2, instance-based techniques are, by
 10 construction, limited to showing learning explanations for a single input (e.g. which particular pixels of an input image where
 11 deemed relevant for determining its class). By contrast, feature-
 12 based techniques do not have this problem as they aim to show
 13 features that are typically relevant for a whole class. Some tech-
 14 niques do this by displaying images from the training set or test
 15 set that produce high activations for the neuron or layer of interest.
 16 Yet, it can be hard for the designer to figure out which specific
 17 features in the displayed image set are responsible for the high
 18 activations [87]. To cope with this, many techniques use
 19 optimization methods to synthesize a ‘summary’ image that
 20 maximizes the activation vector of interest. This type of analysis
 21 allows designers to uncover many properties of deep models.
 22 For example, Nguyen *et al.* [88] present a case where a neuron in a
 23 hidden fully-connected layer activates for different underwater
 24 objects, such as sharks, turtles, and scuba divers, indicating
 25 that these layers often learn high-level concepts that are
 26 present in multiple classes.

28 One of the first techniques to synthesize images that cause
 29 high activations to a particular class is *activation maximization*
 30 [51]. As the activation of a neuron can be seen as a non-linear
 31 function $\phi(\mathbf{x})$ where \mathbf{x} is a given input (see also Fig. 2), the
 32 input $\hat{\mathbf{x}}$ that maximizes the above function can be found by
 33 applying an optimization method—for instance, traditional
 34 gradient ascent—on the input \mathbf{x} with a fixed Θ equal to the
 35 network weights learned after training. This is a non-convex
 36 optimization problem where convergence to an optimal global
 37 value cannot be guaranteed. However, different and meaningful
 38 local optima can be found. These often represent different
 39 facets, or aspects, of the class of interest [88], such as different
 40 properties of an image that may all be equally relevant
 41 for determining its class. This technique was first applied to
 42 deep belief networks [51] and later generalized to deep convolu-
 43 tional networks [90]. Figure 8 shows the results of the
 44 activation maximization technique produced by eight different
 45 classes dataset [5]. The produced synthetic images highlight
 46 typical features that actual images of the respective classes tend
 47 to have.

48 A drawback of activation maximization is that the synthesized
 49 images may be too abstract to interpret, as can be seen in
 50 Fig. 8. This is not surprising, given that the space of all
 51 images of a given class is too large to be captured by a single
 52 ‘average’ image [90, 79]. To address this, several regulariza-
 53 tion methods have been proposed. These methods constrain
 54 the generation of synthetic images by enforcing various criteria
 55 that are typically present in real-world images [92, 90, 91]. For
 56 example, Yosinski *et al.* [62] propose four regularization tech-
 57 niques that aim to synthesize images with more realistic fea-

tures: L_2 decay; Gaussian blur; small norm pixel clipping; and
 58 small contribution pixel clipping. The two first regularizations
 59 aim to remove high brightness amplitudes and high frequencies
 60 that rarely appear in natural images; the last two regulariza-
 61 tions remove pixels with negligible influences on the activa-
 62 tions, letting the designer focus on the important features of
 63 the synthesized image. Activation maximization often produces
 64 images with repeated features, such as multiple objects or ex-
 65aggerated objects of the analyzed class, as such exaggerations
 66 increase the class-specific activation values. For instance, op-
 67 timizing an image for a flamingo-recognizing activations leads
 68 to multiple pelicans scattered over the synthetic image (Fig. 8
 69 top-left). This is one of the causes of the artificial look of such
 70 synthetic images. Nguyen *et al.* [88] alleviate this problem by
 71 using a center-biased regularization penalizes changes close to
 72 the borders of the image, thereby forcing the optimization to
 73 synthesize features closer to the image center.

74 While the above improvements create more interpretable
 75 synthetic images, they still exhibit non-natural colors and bor-
 76 ders. Generating realistic images has, however, been success-
 77 fully achieved by generative neural networks (GNNs) [65, 115].
 78 Nguyen *et al.* [87] proposed to use a GNN model as a prior to
 79 generating realistic images that maximize the activation of a
 80 CNN neuron. GNNs are deep models that learn to generate
 81 novel samples from the distribution in which the training set
 82 lies. In this case, the GNN is trained to generate realistic im-
 83 ages from a numeric vector input. After that, this numeric input
 84 is optimized to generate an image that maximizes the activation
 85 of a given neuron of the CNN.

86 Another issue faced when synthesizing images with activa-
 87 tion maximization is that classes may have very distinct in-
 88 stances. For example, a ‘store’ class could be represented both
 89 by images of the outdoor facade of the store or by images of the
 90 inside of the building. Hence, neurons—especially the ones in
 91 deeper layers—that recognize features of such classes must be
 92 able to activate for very different sets of input features. Such
 93 neurons are then said to be *multifaceted* [88]. In such cases, tra-
 94 ditional optimization methods like gradient ascent end up mix-
 95 ing features of different facets into the resulting synthetic im-
 96 age, which renders it ambiguous and confusing. To overcome
 97 this problem, Nguyen *et al.* introduce a multifaceted feature vi-
 98 sualization [88] that initializes the activation maximization al-
 99 gorithm with an image obtained by averaging instances of the
 100 training set that belong to the same class facet. This creates a
 101 bias aiming to make the algorithm synthesize an image with the
 102 features of that facet.

103 Recurrent Neural Networks:

104 One of the main challenges for RNNs, just like for CNNs
 105 and DFNs, is to understand the activation patterns produced by
 106 specific inputs and which features these inputs are capturing,
 107 since being able to do so is key to understanding the behavior
 108 of the trained model. Like for CNNs and DFNs, heatmaps have
 109 also been used to visualize activations of recurrent models [101,
 110 73]. While it is possible to get good insights on which hidden
 111 states produce higher activations for a single input by looking
 112 at an activation heatmap, it is not easy to understand if the same
 113 hidden states share similar behavior for a group of inputs, *e.g.*,

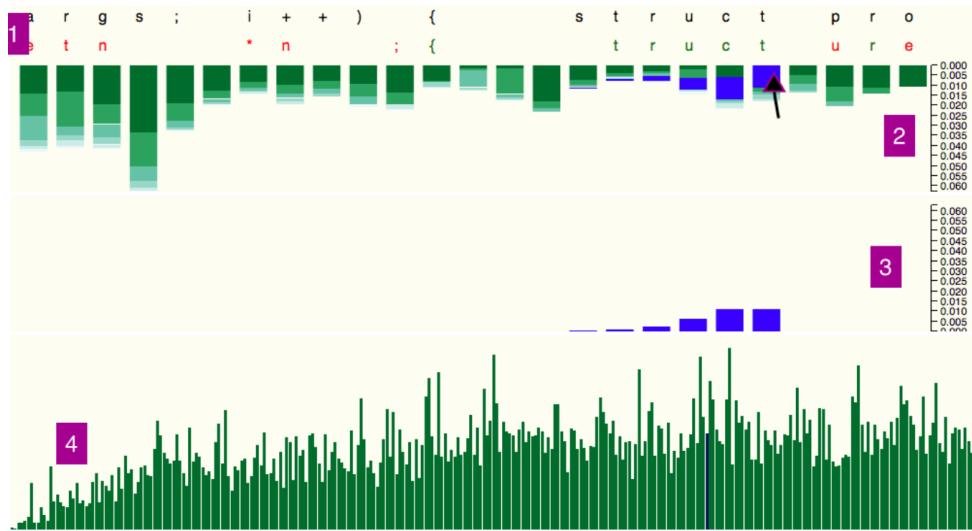


Fig. 11: RNNbow Tool [77] uses stacked barcharts to visualize how the gradient loss progresses through the hidden states of a recurrent neural network.

words with similar meanings. To address this problem, Ming *et al.* [73] propose a technique that co-clusters hidden states and input elements—in particular, it creates two cluster models: one of the activation values of a selected hidden layer and one of the input values given to the network; then it tries to identify correlations between types of inputs that consistently generate activations associated with one particular activation cluster.

Another application of heatmaps for RNN models is to measure the importance of each input unit. Li *et al.* [100] proposed a saliency heatmap that shows the saliency score of each word in the input of a word classification task. This score is calculated by checking how much each input unit contributed to the final label assignment. Ding *et al.* [102] also use heatmaps to display the relevance of each input word for each output word in a machine translation model. To calculate this relevance, they use a layer-wise relevance propagation (LRP) algorithm. Unlike gradient-based approaches, LRP does not require the activations to be differentiable, which confers additional robustness to the approach.

Heatmaps have also been used to display characteristics of the input text data. For instance, Karpathy *et al.* [103]—one of the first works proposing the utilization of visualization to understand RNNs—build a heatmap measuring the relevance of each input character or word to the activation of a given hidden unit. This approach was able to identify interpretable semantic units—for instance, one specialized in identifying new lines—in a character prediction application. Strobelt *et al.* [99] used heatmaps to display designer-defined metrics of interest on the input text, such as to which part-of-speech (POS) a given word belongs shown to (Fig. 12). Conceptually, this visualization is of the same kind of explanatory type as the instance-based visualizations for CNNs (Figs. 9 and 10a), as they highlight parts of an input sample that cause the network to choose a given output class.

RNNs are intrinsically designed to handle sequential input data such as text or time series. When a new input item is processed, the values of the model’s hidden units change. This affects the result not only for the current input item but also for

a (potentially long) range of subsequent input items. Understanding how each input item changes these hidden units and affects latter computations helps identifying critical aspects of the model that may lead to undesirable performance. Given the sequential nature of RNNs, time series charts emerge as a straightforward way to visualize changes in the hidden units’ values. Yet, we have found only a few works in this direction. For instance, Strobelt *et al.* [99] proposed a parallel coordinates plot (PCP) visualization where each dimension is an input unit, and each polyline displays how a single hidden state varies with the input. Hence, the horizontal PCP axis can be interpreted as temporal order, and the polylines are analogous to time series. By allowing the designer to highlight units with high activation for a given range of the input, Strobelt *et al.* were able to find distinct regions of similar behavior in the input text stream, and concluded that such regions had similar semantics.

6. Discussion and Open Challenges

Although a substantial amount of work has been done over the past few years regarding the use of visual analytics in deep learning techniques, there are still many challenges that need to be successfully addressed by future research. We discuss these challenges from the point of view of the three types of *tasks* that we structured our survey along (i.e., architecture understanding, training analysis, and feature understanding) and, additionally, from the point of view of end-to-end *requirements* that engineers developing deep models face in practice.

6.1. Architecture understanding

Hyperparameter exploration: One problem that is still open even for machine learning experts is how hyperparameters affect the training results. Developers of neural networks usually make architectural choices such as the number of neurons per layer, number of layers, activation function types, training batch size, learning rate, and number of training epochs, by empirical and trial-and-error approaches. This considerably adds to the



Fig. 12: LSTMVis tool [99] showing the activations of the hidden states in a RNN model for a sequence of input units chosen by the designer.

cost of fine-tuning the training of deep models. In a more broad sense, analysing the parameter space of a simulation model is a topic that has caught a strong interest from the VA community in the past years [75]. However, that is still a lot of space for novel contributions that can help designers to more precisely tune their deep networks. Interactive ‘drawing board’ solutions where all these aspects can be directly controlled by a designer, while their effects are visualized in real time, exist [82] but cover only very small networks (tens of neurons) with simple two-dimensional inputs.

It would be interesting for machine learning experts to visually analyze the hyperparameter space of a neural model and how values in this space affect the model performance [34]. To do this, visually more scalable techniques (capable of depicting large network architectures) are needed, as well as methods that annotate the hyperparameter space with measured network performance. There is work on the use Bayesian optimization for actively selecting which hyperparameter values to try next, based on annotating particular settings of such values with the corresponding model performance and trying to infer which ones might work better [116]. Although there is some work about applying such techniques to deep models [117, 118], more research is needed in this direction.

Additionally, faster training pipelines, possibly based on multiscale techniques, are needed to close the sensemaking loop at interactive rates, thereby allowing designers to effectively ‘steer’ the architecture design as they observe its behavior.

6.2. Training analysis

Training data exploration: To be properly trained, deep models often require large and high-dimensional training sets. However, to date, there are only a few solutions to understanding such training sets, and in particular, which aspects of the input data affect training in specific ways. Training instances may contain hidden biases or mistakes, such as mislabeling, irrelevant correlations of input features with classes [13, 52], and, at

a higher level, unbalanced coverage of the entire input space X by training samples. All of these aspects can severely harm the effectiveness of the trained model. Visualization techniques, notably the ones focused on the analysis of high-dimensional datasets [14, 52, 76], are a promising alternative to address this issue. Additionally, a better understanding of the training set characteristics can lead to more efficient choices of architectures and hyperparameters.

Training guiding and interaction: Due to the difficulty in understanding machine learning techniques and the long time required to train deep models, interactive solutions have received significant interest from the machine learning community [36, 37, 38]. Visualization methods are key to achieve such interaction, as they can quickly show the designer what is happening during the training process. However, only a few approaches have focused on using visualization to provide real-time feedback to the designer [78, 83, 82]. The challenge here is directly related to hyperparameter exploration, *i.e.* providing both visually scalable and computationally scalable (fast) metaphors for the training process.

6.3. Feature understanding

Explainable models: Visual analytics has proven to be an effective tool in explaining the features learned by neural models. Current visualization methods can show which features (from the input data) have been learned by a given model. Many types of features can be considered by these techniques, such as pixels in an image, words in a text document, and value ranges of input data attributes (columns in a data table), each of which computed either per input sample or per set of related instances, *e.g.* class or class facet. Solutions produced by this type of visualization highlight the most discriminative features for determining a class [13] and follow the intuitions used by earlier methods that aimed at achieving a similar goal but for classifiers that used hand-engineered features [20]. The need for more explainable models, however, is still noticeable. In particular, if

one could say *which* input features are responsible for a model’s decision, the next step would be to show *how* that decision was made. This involves explaining the responsibilities of groups of layers or neurons of a DNN and how these work together to calculate the final output [34]. In the long run, this will lift the current feature understanding goal to cover the more important goal of understanding how a model as a whole took a given decision.

6.4. Non-Functional Requirements

Apart from the functional requirements for the visual understanding process of deep learning mentioned above, some non-functional requirements exist, as discussed next.

Fidelity: Modern DNNs can have hundreds of thousands of neurons spread over hundreds of layers and millions of parameters [4, 8, 9]. The sheer amount of data embedded in, or produced by, such models demands novel VA techniques that scale effectively. The key issue here is that of fidelity or trust: when data is aggregated or simplified, how can we be sure that we trust what the visualization shows? For instance, many approaches use dimensionality reduction (DR) techniques to visualize the high-dimensional data produced by DNNs. Small changes on the hyperparameters of DR methods can lead to massively different visualizations that may easily convey different or wrong insights [119, 120, 121]. The goal of optimizing dimensionality reduction so that it accurately conveys the high-dimensional data structure of large datasets (millions of observations, hundreds of dimensions) is an ongoing research endeavor [122, 123].

Scalability: The large size of modern DNNs poses two scalability problems. First, we need to develop *visually* scalable in-fovis metaphors to depict the large amount of high-dimensional, temporal, and relational data spanned by such networks. This is in itself a key challenge in information visualization, for which answers are yet to be found. Separately, we need access to *computationally* scalable ways of performing DNN training, so that insights found in this process can be communicated to the designer at interactive rates, for the VA sensemaking loop to be effectively closed. In cases where interactive-rate training of DNNs is simply not possible due to the size of the problem, approximation methods could be developed that deliver a less accurate, but still insightful, view on the training process at interactive rates.

Different applications: Currently, most research works are focusing mainly on models handling image classification or natural language processing tasks that use well-behaved training sets. However, more complex applications usually have to handle training sets that may be composed of different types of data and that may come from distinct sources, thus requiring more complicated architectures [80]. Visual analytics tools directing towards these models could be an interesting perspective for future research.

7. Conclusion

In this article, we reviewed works aiming at using visual analytics techniques to understand deep neural networks—a topic

that has been widely discussed by the research community in the past few years. We classified these publications into three categories, depending on the particular visualization goal that they tried to achieve: network architecture understanding, visualization to support training analysis, and feature understanding. In particular, we were able to identify that most of the reviewed publications have been mainly focused on understanding which features a given learned model can recognize, how they do it, and how the learning of such features occurs. These visualization approaches have proven to be effective in validating the performance of deep models and providing more intuition of their inner workings to the designer of the respective neural network model. However, there is still a lack of contributions aiming at developing techniques that can interactively guide the development of a deep neural network, with only a few approaches addressing this issue. Given that deep networks are usually difficult and slow to train, we consider this as a promising topic of future research, with visual analytics playing a key role in providing such visual interactivity to machine learning experts.

References

- [1] Goodfellow, I, Bengio, Y, Courville, A. Deep Learning. The MIT Press; 2016.
- [2] Samuel, AL. Some studies in machine learning using the game of checkers. IBM Journal of Research and Development 1959;3(3):210–229.
- [3] LeCun, Y, Bengio, Y, Hinton, G. Deep learning. Nature 2015;521(7553):436–444.
- [4] Krizhevsky, A, Sutskever, I, Hinton, GE. ImageNet classification with deep convolutional neural networks. In: Proc. International Conference on Neural Information Processing Systems; vol. 1. 2012, p. 1097–1105.
- [5] Deng, J, Dong, W, Socher, R, Li, LJ, Li, K, Fei-Fei, L. ImageNet: A large-scale hierarchical image database. In: Proc. IEEE Conference on Computer Vision and Pattern Recognition. 2009, p. 248–255.
- [6] Zeiler, MD, Fergus, R. Visualizing and understanding convolutional networks. In: Proc. European Conference on Computer Vision. Springer; 2014, p. 818–833.
- [7] Simonyan, K, Zisserman, A. Very deep convolutional networks for large-scale image recognition. 2015. [arXiv:1409.1556](https://arxiv.org/abs/1409.1556).
- [8] Szegedy, C, Liu, W, Jia, Y, Sermanet, P, Reed, S, Anguelov, D, et al. Going deeper with convolutions. In: Proc. IEEE Conference on Computer Vision and Pattern Recognition. 2015, p. 1–9.
- [9] He, K, Zhang, X, Ren, S, Sun, J. Deep residual learning for image recognition. In: Proc. IEEE Conference on Computer Vision and Pattern Recognition. 2016, p. 165–173.
- [10] Liu, M, Shi, J, Li, Z, Li, C, Zhu, J, Liu, S. Towards better analysis of deep convolutional neural networks. IEEE Transactions on Visualization and Computer Graphics 2017;23(1):91–100.
- [11] Marcus, G. Deep learning: A critical appraisal. 2018. [arXiv:1801.00631](https://arxiv.org/abs/1801.00631).
- [12] Samek, W, Wiegand, T, Müller, KR. Explainable artificial intelligence: Understanding, visualizing and interpreting deep learning models. 2017. [arXiv:1708.08296v1](https://arxiv.org/abs/1708.08296v1).
- [13] Ribeiro, MT, Singh, S, Guestrin, C. Why should i trust you? explaining the predictions of any classifier. 2016. [arXiv:1602.04938](https://arxiv.org/abs/1602.04938).
- [14] Liu, S, Majovic, D, Wang, B, Bremer, PT, Pascucci, V. Visualizing high-dimensional data: Advances in the past decade. Computer Graphics Forum 2016;23(3):1249–1268.
- [15] van der Maaten, L, Postma, E, van den Herik, H. Dimensionality reduction: A comparative review. 2009. Technical Report 2009-005, Univ. Tilburg, Netherlands.
- [16] Cunningham, JP, Ghahramani, Z. Linear dimensionality reduction: Survey, insights, and generalizations. Journal of Machine Learning Research 2015;16:2859–2900.
- [17] Sorzano, C, Vargas, J, Montano, AP. A survey of dimensionality reduction techniques. 2014. [arXiv:1403.2877](https://arxiv.org/abs/1403.2877).

- [18] Blum, MGB, Nunes, MA, Prangle, D, Sisson, SA. A comparative review of dimension reduction methods in approximate bayesian computation. *Statistical Science* 2013;28(2):189–208.
- [19] Rauber, P, Falcão, A, Telea, A. Projections as visual aids for classification system design. *Information Visualization* 2017;.
- [20] Rauber, P, da Silva, R, Feringa, S, Celebi, M, Falcão, A, Telea, A. Interactive image feature selection aided by dimensionality reduction. In: Proc. EuroVA. 2015, p. 67–74.
- [21] Krause, J, Perer, A, Bertini, E. INFUSE: Interactive feature selection for predictive modeling of high dimensional data. *IEEE Transactions on Visualization and Computer Graphics* 2014;20(12):1614–1623.
- [22] Yuan, X, Ren, D, Wang, Z, Guo, C. Dimension projection matrix/tree: Interactive subspace visual exploration and analysis of high dimensional data. *IEEE Transactions on Visualization and Computer Graphics* 2013;19(12):2625–2633.
- [23] Tatu, A, Maas, F, Farber, I, Bertini, E, Schreck, T, Seidl, T, et al. Subspace search and visualization to make sense of alternative clusterings in high-dimensional data. In: Proc. IEEE VAST. 2012, p. 63–72.
- [24] Turkay, C, Filzmoser, P, Hauser, H. Brushing dimensions: A dual visual analysis model for high-dimensional data. *IEEE Transactions on Visualization and Computer Graphics* 2011;17(12):2591–2599.
- [25] Noris, B. MLdemos: Open source visualization tool for machine learning algorithms. 2017. <http://mldemos.epfl.ch>.
- [26] Gleicher, M. Explainers: Expert explorations with crafted projections. *IEEE Transactions on Visualization and Computer Graphics* 2013;19(12):2042–2051.
- [27] Gleicher, M. A framework for considering comprehensibility in modeling. *Big Data* 2016;4(2):75–88.
- [28] Montavon, G, Samek, W, Müller, KR. Methods for interpreting and understanding deep neural networks. *Digital Signal Processing* 2018;73:1–15.
- [29] Samek, W, Binder, A, Montavon, G, Lapuschkin, S, Müller, KR. Evaluating the visualization of what a deep neural network has learned. *IEEE Transactions on Neural Networks and Learning Systems* 2017;28(11):2660–2673.
- [30] Yeager, L, Heinrich, G, Mancewicz, J, Houston, M. Effective visualizations for training and evaluating deep models. In: Proc. International Conference on Machine Learning Workshop on Visualization for Deep Learning. 2016.,
- [31] Zeng, H. Towards better understanding of deep learning with visualization. 2016. MSc thesis, Dept. of Computer Science and Engineering, Hong-Kong Univ. of Science and Technology.
- [32] Seifert, C, Aamir, A, Balagopalan, A, Jain, D, Sharma, A, Grottel, S, et al. Visualizations of deep neural networks in computer vision: A survey. In: *Transparent Data Mining for Big and Small Data*. Springer; 2017, p. 123–144.
- [33] Grün, F, Rupprecht, C, Navab, N, Tombari, F. A taxonomy and library for visualizing learned features in convolutional neural networks. 2016. [arXiv:1606.07757](https://arxiv.org/abs/1606.07757).
- [34] Liu, S, Wang, X, Liu, M, Zhu, J. Towards better analysis of machine learning models: A visual analytics perspective. *Visual Informatics* 2017;1(1):48 – 56.
- [35] Lu, Y, Garcia, R, Hansen, B, Gleicher, M, Maciejewski, R. The state-of-the-art in predictive visual analytics. *Computer Graphics Forum* 2017;36(3):539–562.
- [36] Amershi, S, Cakmak, M, Knox, WB, Kulesza, T. Power to the people: The role of humans in interactive machine learning. *AI Magazine* 2014;35(4):105–120.
- [37] Bernardo, F, Zbyszynski, M, Fiebrink, R, Grierson, M, et al. Interactive machine learning for end-user innovation. In: Proc. Designing the User Experience of Machine Learning Systems (AAAI Spring Symposium Series). 2017..
- [38] Sacha, D, Sedlmair, M, Zhang, L, Lee, JA, Weiskopf, D, North, S, et al. Human-centered machine learning through interactive visualization. In: Proc. European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning. 2016.,
- [39] Lipton, Z. The mythos of model interpretability. 2016. [arXiv:1606.03490](https://arxiv.org/abs/1606.03490).
- [40] Hohman, F, Kahng, M, Pienta, R, Chau, DH. Visual analytics in deep learning: An interrogative survey for the next frontiers. 2018. [arXiv:1801.06889](https://arxiv.org/abs/1801.06889).
- [41] Murphy, K. Machine learning: a Probabilistic Perspective. The MIT Press; 2012.
- [42] Bishop, CM. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. 2006.
- [43] Hastie, T, Tibshirani, R, Friedman, J. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer; 2009.
- [44] Breiman, L. Random forests. *Machine Learning* 2001;45(1):5–32.
- [45] LeCun, Y, Boser, B, Denker, JS, Henderson, D, Howard, RE, Hubbard, W, et al. Backpropagation applied to handwritten zip code recognition. *Neural Computation* 1989;1(4):541–551.
- [46] Elman, JL. Finding structure in time. *Cognitive Science* 1990;14(2):179–211.
- [47] Park, SH, Goo, JM, Jo, CH. Receiver operating characteristic (ROC) curve: Practical review for radiologists. *Korean Journal of Radiology* 2004;5(1):1118.
- [48] Powers, D. Evaluation: From Precision, Recall and F-Measure to ROC, Informedness, Markedness and Correlation. *Journal of Machine Learning Technology* 2011;2(1):37–63.
- [49] Fawcett, T. An introduction to ROC analysis. *Pattern Recognition Letters* 2006;27(8):861–874.
- [50] Brehmer, M, Munzner, T. A multi-level typology of abstract visualization tasks. *IEEE Transactions on Visualization and Computer Graphics* 2013;19(12):2376–2385.
- [51] Erhan, D, Bengio, Y, Courville, A, Vincent, P. Visualizing higher-layer features of a deep network. University of Montreal 2009;Technical Report 1341.
- [52] Rauber, P, Fadel, SG, Falcão, A, Telea, A. Visualizing the hidden activity of artificial neural networks. *IEEE Transactions on Visualization and Computer Graphics* 2017;23(1):101–110.
- [53] Wong, PC, Thomas, J. Visual analytics. *IEEE Computer Graphics and Applications* 2004;24(5):20–21.
- [54] Keim, DA, Mansmann, F, Schneidewind, J, Thomas, J, Ziegler, H. Visual analytics: Scope and challenges. In: *Lecture Notes in Computer Science (LNCS 4404)*. Springer; 2008, p. 76–90.
- [55] Lu, J, Chen, W, Ma, Y, Ke, J, Li, Z, Zhang, F, et al. Recent progress and trends in predictive visual analytics. *Frontiers of Computer Science* 2017;11(2):192–207.
- [56] Pirolli, P, Card, S. The sensemaking process and leverage points for analyst technology as identified through cognitive task analysis. In: Proc. International Conference on Intelligence Analysis. 2005.,
- [57] Keim, D, Andrienko, G, Fekete, JD, Görg, C, Kohlhammer, J, Melancon, G. Visual analytics: Definition, process, and challenges. In: *Information Visualization – Human-Centered Issues and Perspectives*. Springer; 2008, p. 154–175.
- [58] Keim, D, Kohlhammer, J, Ellis, G, Mansmann, F. Mastering the information age: Solving problems with visual analytics. *Eurographics Association*; 2010.
- [59] IEEE VAST 2017 Symposium. 2017. <http://ieelevis.org/year/2017/info/papers>.
- [60] Streeter, MJ, Ward, MO, Alvarez, SA. NVIS: an interactive visualization tool for neural networks. In: Proc. SPIE Visual Data Exploration and Analysis; vol. 4302. 2001, p. 1–8.
- [61] Tzeng, FY, Ma, KL. Opening the black box – data driven visualization of neural networks. In: Proc. IEEE Visualization. 2005, p. 383–390.
- [62] Yosinski, J, Clune, J, Fuchs, T, Lipson, H. Understanding neural networks through deep visualization. In: Proc. International Conference on Machine Learning Workshop on Deep Learning. 2015.,
- [63] Alsallakh, B, Jourabloo, A, Ye, M, Liu, X, Ren, L. Do convolutional neural networks learn class hierarchy? *IEEE Transactions on Visualization and Computer Graphics* 2018;24(1):152–162.
- [64] Hinton, GE, Zemel, RS. Autoencoders, minimum description length and helmholtz free energy. In: Proc. International Conference on Neural Information Processing Systems. NIPS'93; San Francisco, CA, USA: Morgan Kaufmann Publishers Inc.; 1993, p. 3–10.
- [65] Goodfellow, I, Pouget-Abadie, J, Mirza, M, Xu, B, Warde-Farley, D, Ozair, S, et al. Generative adversarial nets. In: *Advances in Neural Information Processing Systems 27*. Curran Associates, Inc.; 2014, p. 2672–2680.
- [66] Hinton, GE, Osindero, S, Teh, YW. A fast learning algorithm for deep belief nets. *Neural Computation* 2006;18(7).
- [67] Mnih, V, Kavukcuoglu, K, Silver, D, Rusu, AA, Veness, J, Bellemare, MG, et al. Human-level control through deep reinforcement learning. *Nature* 2015;518(7540).

- [68] Gibson, A, Patterson, J. Deep Learning. O'Reilly Media Inc.; 2017.
- [69] Zahavy, T, Ben-Zrihem, N, Mannor, S. Graying the black box: Understanding dqn's. In: Proc. International Conference on Machine Learning. 2016, p. 1899–1908.
- [70] Liu, M, Shi, J, Cao, K, Zhu, J, Liu, S. Analyzing the training processes of deep generative models. *IEEE Transactions on Visualization and Computer Graphics* 2018;24(1):77–87.
- [71] Ren, D, Amershi, S, Lee, B, Suh, J, Williams, JD. Squares: Supporting interactive performance analysis for multiclass classifiers. *IEEE Transactions on Visualization and Computer Graphics* 2017;23(1):61–70.
- [72] Mühlbacher, T, Piringer, H. A partition-based framework for building and validating regression models. *IEEE Transactions on Visualization and Computer Graphics* 2013;19(12):1962–1971.
- [73] Ming, Y, Cao, S, Zhang, R, Li, Z, Chen, Y, Song, Y, et al. Understanding hidden memories of recurrent neural networks. In: Proc. IEEE Visual Analytics Science and Technology (VAST). 2017..
- [74] Zeng, H, Haleem, H, Plantaz, X, Cao, N, Qu, H. CNNComparator: Comparative analytics of convolutional neural networks. In: Proc. Workshop on Visual Analytics for Data Learning (VADL). 2017..
- [75] Sedlmair, M, Heinzel, C, Bruckner, S, Piringer, H, Müller, T. Visual parameter space analysis: A conceptual framework. *IEEE Transactions on Visualization and Computer Graphics* 2014;20(12):2161–2170.
- [76] Pezzotti, N, Holtl, T, Gemert, JV, Lelieveldt, BP, Eisemann, E, Vilanova, A. DeepEyes: Progressive visual analytics for designing deep neural networks. *IEEE Transactions on Visualization and Computer Graphics* 2018;24(1):98–108.
- [77] Cashman, D, Patterson, G, Mosca, A, Chang, R. RNNbow: Visualizing learning via backpropagation gradients in recurrent neural networks. In: Proc. Workshop on Visualization for Deep Learning (VADL). 2017..
- [78] Qi, H, Liu, J, Zou, X, Tang, A. BIDViz: Real-time monitoring and debugging of machine learning training processes. Master's thesis; EECS Department, University of California, Berkeley; 2017.
- [79] Nguyen, A, Yosinski, J, Clune, J. Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. In: Proc. IEEE Conference on Computer Vision and Pattern Recognition. 2015..
- [80] Kahng, M, Andrews, PY, Kalro, A, Chau, DH. Activis: Visual exploration of industry-scale deep neural network models. *IEEE Transactions on Visualization and Computer Graphics* 2018;24(1):88–97.
- [81] Wongsuphasawat, K, Smilkov, D, Wexler, J, Wilson, J, Man, D, Fritz, D, et al. Visualizing dataflow graphs of deep learning models in tensorflow. *IEEE Transactions on Visualization and Computer Graphics* 2018;24(1):1–12.
- [82] Smilkov, D, Carter, S, Sculley, D, Vigas, FB, Wattenberg, M. Direct-manipulation visualization of deep networks. 2017. [arXiv:1708.03788](#).
- [83] Chung, S, Suh, S, Park, C, Kang, K, Choo, J, Kwon, BC. RevaCNN: Real-Time visual analytics for convolutional neural network. In: Proc. ACM SIGKDD Workshop on Interactive Data Exploration and Analytics (IDEA). 2016..
- [84] Harley, AW. An interactive node-link visualization of convolutional neural networks. In: Proc. International Symposium on Advances in Visual Computing (ISVC). Springer; 2015, p. 867–877.
- [85] Zhong, W, Xie, C, Zhong, Y, Wang, Y, Xu, W, Cheng, S, et al. Evolutionary visual analysis of deep neural networks. In: Proc. International Conference on Machine Learning Workshop on Visualization for Deep Learning. 2017..
- [86] Selvaraju, RR, Cogswell, M, Das, A, Vedantam, R, Parikh, D, Batra, D. Grad-CAM: Visual explanations from deep networks via gradient-based localization. 2016. [arXiv:1610.02391](#).
- [87] Nguyen, A, Dosovitskiy, A, Yosinski, J, Brox, T, Clune, J. Synthesizing the preferred inputs for neurons in neural networks via deep generator networks. In: Proc. International Conference on Neural Information Processing Systems. 2016, p. 3395–3403.
- [88] Nguyen, A, Yosinski, J, Clune, J. Multifaceted feature visualization: Uncovering the different types of features learned by each neuron in deep neural networks. 2016. [arXiv:1602.03616](#).
- [89] Aubry, M, Russell, BC. Understanding deep features with computer-generated imagery. In: Proc. IEEE International Conference on Computer Vision (ICCV). 2015..
- [90] Simonyan, K, Vedaldi, A, Zisserman, A. Deep inside convolutional networks: Visualising image classification models and saliency maps. 2013. [arXiv:1312.6034](#).
- [91] Wei, D, Zhou, B, Torralba, A, Freeman, W. Understanding intra-class knowledge inside CNN. 2015. [arXiv:1507.02379](#).
- [92] Mahendran, A, Vedaldi, A. Understanding deep image representations by inverting them. In: Proc. IEEE Conference on Computer Vision and Pattern Recognition. 2015..
- [93] Mahendran, A, Vedaldi, A. Visualizing deep convolutional neural networks using natural pre-images. *International Journal of Computer Vision* 2016;120(3):233–255.
- [94] Zintgraf, LM, Cohen, TS, Welling, M. A new method to visualize deep neural networks. 2016. [arXiv:1603.02518](#).
- [95] Dosovitskiy, A, Brox, T. Inverting visual representations with convolutional networks. In: Proc. IEEE Conference on Computer Vision and Pattern Recognition. 2016..
- [96] Zintgraf, LM, Cohen, TS, Adel, T, Welling, M. Visualizing deep neural network decisions: Prediction difference analysis. 2017. [arXiv:1702.04595](#).
- [97] Li, H, Mueller, K, Chen, X. Beyond saliency: understanding convolutional neural networks from saliency prediction on layer-wise relevance propagation. 2017. [arXiv:1712.08268](#).
- [98] Bojarski, M, Choromanska, A, Choromanski, K, Firner, B, Jackel, L, Müller, U, et al. VisualBackProp: Efficient visualization of CNNs. 2016. [arXiv:1611.05418](#).
- [99] Strobelt, H, Gehrmann, S, Pfister, H, Rush, AM. LSTMVis: A tool for visual analysis of hidden state dynamics in recurrent neural networks. *IEEE Transactions on Visualization and Computer Graphics* 2018;24(1):667–676.
- [100] Li, J, Chen, X, Hovy, E, Jurafsky, D. Visualizing and understanding neural models in nlp. 2015. [arXiv:1506.01066](#).
- [101] Rong, X, Adar, E. Visual tools for debugging neural language models. In: Proc. International Conference on Machine Learning Workshop on Visualization for Deep Learning. 2016..
- [102] Ding, Y, Liu, Y, Luan, H, Sun, M. Visualizing and understanding neural machine translation. In: Proc. Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers); vol. 1. 2017, p. 1150–1159.
- [103] Karpathy, A, Johnson, J, Fei-Fei, L. Visualizing and understanding recurrent networks. 2015. [arXiv:1506.02078](#).
- [104] Diehl, S. Software Visualization: Visualizing the Structure, Behaviour, and Evolution of Software. Springer; 2007.
- [105] van der Zwan, M, Codreanu, V, Telea, A. CUBu: Universal real-time bundling for large graphs. *IEEE Transactions on Visualization and Computer Graphics* 2016;22(12):2250–2263.
- [106] Lhuillier, A, Hurter, C, Telea, A. State of the art in edge and trail bundling techniques. *Computer Graphics Forum* 2017;36(3):619–645.
- [107] Zhang, L, Wang, S, Liu, B. Deep learning for sentiment analysis: A survey. 2018. [arXiv:1801.07883](#).
- [108] Zhang, J, Zong, C. Deep neural networks in machine translation: An overview. *IEEE Intelligent Systems* 2015;30(5):16–25.
- [109] Maaten, Lvd, Hinton, G. Visualizing data using t-SNE. *Journal of Machine Learning Research* 2008;9(Nov):2579–2605.
- [110] Hochreiter, S, Bengio, Y, Frasconi, P, Schmidhuber, J. Gradient flow in recurrent nets: the difficulty of learning long-term dependencies. IEEE Press; 2001, p. 464.
- [111] Netzer, Y, Wang, T, Coates, A, Bissacco, A, Wu, B, Ng, AY. Reading digits in natural images with unsupervised feature learning. In: Proc. Neural Information Processing Systems. 2011, p. 5–12.
- [112] Sacha, D, Zhang, L, Sedlmair, M, Lee, JA, Peltonen, J, Weiskopf, D, et al. Visual interaction with dimensionality reduction: A structured literature analysis. *IEEE Transactions on Visualization and Computer Graphics* 2017;23(1):241–250.
- [113] Montavon, G, Lapuschkin, S, Binder, A, Samek, W, Müller, KR. Explaining nonlinear classification decisions with deep Taylor decomposition. *Pattern Recognition* 2017;65:211 – 222.
- [114] Bach, S, Binder, A, Montavon, G, Klauschen, F, Müller, KR, Samek, W. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLOS One* 2015;10.
- [115] Dosovitskiy, A, Brox, T. Generating images with perceptual similarity metrics based on deep networks. In: Advances in Neural Information Processing Systems 29. Curran Associates, Inc.; 2016, p. 658–666.
- [116] Brochu, E, Cora, VM, De Freitas, N. A tutorial on bayesian optimiza-

- 1 tion of expensive cost functions, with application to active user modeling
2 and hierarchical reinforcement learning. 2010. [arXiv:1012.2599](https://arxiv.org/abs/1012.2599).
- 3 [117] Snoek, J, Larochelle, H, Adams, RP. Practical bayesian optimization
4 of machine learning algorithms. In: Advances in Neural Information
5 Processing Systems 25. Curran Associates, Inc.; 2012.,
- 6 [118] Snoek, J, Rippel, O, Swersky, K, Kiros, R, Satish, N, Sundaram, N,
7 et al. Scalable bayesian optimization using deep neural networks. In:
8 Proc. International Conference on Machine Learning. 2015.,
- 9 [119] Wattenberg, M. How to use t-SNE effectively. 2017. <https://distill.pub/2016/misread-tsne>.
- 10 [120] Martins, R, Coimbra, D, Minghim, R, Telea, A. Visual analysis of
11 dimensionality reduction quality for parameterized projections. Computers & Graphics 2014;41:26–42.
- 12 [121] Rauber, P, Falcão, A, Telea, A. Visualizing time-dependent data using
13 dynamic t-SNE. In: Proc. EuroVis – short papers. 2016, p. 137–142.
- 14 [122] McInnes, L, Healy, J. Umap: Uniform manifold approximation and
15 projection for dimension reduction. 2018. [arXiv:1802.03426](https://arxiv.org/abs/1802.03426).
- 16 [123] Pezzotti, N, Höllt, T, Lelieveldt, BP, Eisemann, E, Vilanova, A.
17 Hierarchical stochastic neighbor embedding. Computer Graphics Forum
18 2016;35(3):21–30.
- 19
20