

Learning Interpretation with Explainable Knowledge Distillation

Raed Alharbi Minh N. Vu My T. Thai

Computer and Information Science and Engineering Department

University of Florida, Gainesville, FL, USA

{r.alharbi, minhvu, mythai}@ufl.edu

Abstract—Knowledge Distillation (KD) has been considered as a key solution in model compression and acceleration in recent years. In KD, a small student model is generally trained from a large teacher model by minimizing the divergence between the probabilistic outputs of the two. However, as demonstrated in our experiments, existing KD methods might not transfer critical explainable knowledge of the teacher to the student, i.e. the explanations of predictions made by the two models are not consistent. In this paper, we propose a novel explainable knowledge distillation model, called XDistillation, through which both the performance the explanations’ information are transferred from the teacher model to the student model. The XDistillation model leverages the idea of convolutional autoencoders to approximate the teacher explanations. Our experiments shows that models trained by XDistillation outperform those trained by conventional KD methods not only in term of predictive accuracy but also faithfulness to the teacher models.

Index Terms—Explainable Machine Learning, Knowledge Distillation, Knowledge Transfer, Neural Network Distillation

I. INTRODUCTION

With the extensive deployment of deep neural networks models (DNNs) on lightweight and low computing resources, such as mobile devices, or Internet-of-Thing (IoT) devices, Knowledge Distillation (KD) has been shown as one of the most promising approaches to transfer knowledge from a large model, called teacher, to a smaller one, called student, without loss of predictive power and validity [1], [2]. The core concept is that the teacher model is utilized during a KD process to guide the student model by passing on substantial information. With similar performance, the student models have much less parameters and can be deployed on less powerful hardware. KD has been successfully used in several application domains, such as computer vision and natural language processing [3].

With an increasing alarm by the public and researchers on the lack of interpretability of current DNNs, that is, they have been used as black-boxes with a little explanation, transferring the explainability is as important as preserving the predictive accuracy in knowledge distillation. Being able to adopt the teacher explanations is significantly desirable because: 1) Explanations provide transparency to the models’ prediction, thereby increasing trustworthiness in using models [4]. 2) Trustworthiness and faithful explanations can identify models’ failures and bias when not all possible scenarios are

testable [5], thereby avoiding several shortcuts learning that existing DNN models have been exhibiting [6].



Fig. 1: The inconsistency between the explanations of the predictions of the two models.

Unfortunately, while existing KD approaches have been focused on preserving the performance accuracy, the explanation on why the model makes its prediction is not transferred from the teacher to the student model. As most of the existing teacher models are not interpretable by themselves, we use a post-hoc explainer SHAP [7] to explain the important features to the model’s prediction. Figure 1 illustrates an inability of transferring the explanation knowledge of existing KD models. As can be seen in Figure 1, the explanations of the teacher and student models of the same image are inconsistent, using SHAP.

Along this direction, this paper proposes an effective KD model, called XDistillation, which not only maintains the teacher’s performance but also approximates the teacher’s explanation. Sitting in the heart of XDistillation is a novel explainable features fusion technique that significantly reduces the total number of student parameters while ensuring consistency between teachers’ and students’ explanations. Extensive experiments demonstrate that XDistillation provides consistent explanations between the teacher and student models while being on par with existing KD techniques in terms of performance accuracy.

The remainder of the paper is structured as follows. Section II presents the related works. Our proposed XDistillation model is introduced in Section III while the experimental analysis are discussed in Section IV. Finally, Section V concludes our paper.

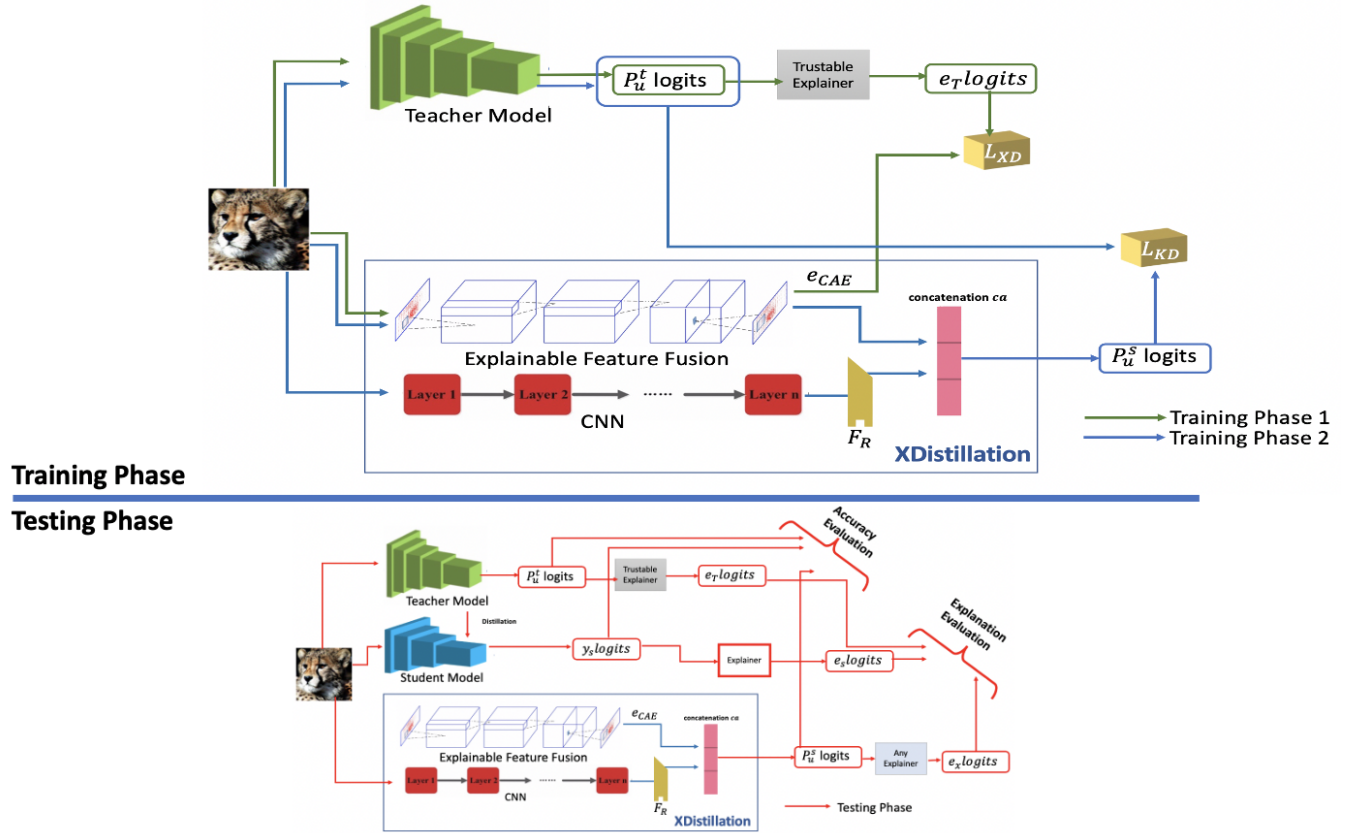


Fig. 2: The overall architecture of XDistillation. The top network illustrates the training phases using two losses functions as described with more details in section III. The first loss minimizes the inconsistency between the explanations of a teacher model and our approximated explanations by a novel explainable feature fusion technique. The second loss minimizes the predictions among the two models. The bottom network depicts the evaluation measures in term of the accuracy performance and explanations.

II. RELATED WORK

Knowledge Distillation. The concept of knowledge distilling relates to the idea of model compression, where a teacher model guides the student model while retaining a high degree of generality [1]. Since then, there have been studies on various methods of knowledge distillation to improve the student model. Zagoruyko *et al.* in [8] proposed the use of attention maps where Heo *et al.* [9] utilize the activated neurons to distill only the valuable information to a student model, and Tung *et al.* [10] follow activation patterns of neurons between a teacher and student model. Yet, no prior attempts have been made to distill the explanations from a teacher a student to improve both a student's performance and explanations.

Explainable Artificial Intelligence (XAI). One of the major deficiencies to the deployment of DNN models is the lack of transparency [11]. The black-box structure of these models permits strong predictions. However, they cannot be explained explicitly. This problem has initiated a new discussion about explainable AI [12]. The purpose of employing interpretability in AI models is to verify the model predictions and the important features that contribute to these predictions can be

interpretable to end-users [11], [12] whereas preserving a high degree of accuracy performance in the same time.

Two popular types of methods have been presented in the literature to tackle the interpretability needs: (1) backpropagation-based such as CAM [13] and GradCAM [14]. (2) perturbation-based such as occlusion analysis [15], LIME [16], and SHAP [7]. On the one hand, predictions and explanations in backpropagation-based methods are both generated by the same fundamental technique. Those methods use the properties of Convolutional Neural Networks (CNNs) by utilizing the activations of the models convolution layers to describe the relation of the input to the output [14]. On the other hand, the perturbation-based methods address the impact of the input perturbations on the output to identify important features for explanations.

Autoencoders (AEs) Using an autoencoder is a common technique for ensuring that the input and output are as comparable as possible [17]. The autoencoder is a form of neural network that is capable of learning a compressed version of the input data in unsupervised manner. Autoencoders have been made significant contributions to the field's application

and research including dimensionality reduction [17], image improvement [18] and detection of outliers [19]. Recently, different variant of AEs has been proposed such as convolutional autoencoders (CAE) [20] and denoising autoencoders (DAE) [21]. One of the primary benefits of AEs is their ability to approximate any function. However, no attempt has been made previously to approximate the explanations of model predictions

In this paper, we select well-known explainers from backpropagation-based and perturbation-based categories to use them during the training and testing phases, as it describes later in Sections III and IV. Particularly, for the first category, GradCAM is chosen as it is computationally efficient (requiring only a single forward and backward pass via a network). For a particular decision, GradCAM gives a value to each neuron utilizing gradient information that flows into the CNN’s last convolutional layer [14]. LIME and SHAP are selected from the second category as providing users with explanations in the form of feature importance score has been seen as a successful technique [7], [16].

III. THE PROPOSED MODEL

We design a simple but effective model, called XDistillation, which provides explanation-distillation from the teacher to increase the performance and enhance the explainability of the student. Figure 2 shows an overview of XDistillation where the key functionalities are explainable features fusion component, how the interpretable features output of this fusion component is fed into Convolutional neural networks (CNNs) via concatenation layers, and how to constrain the signal outputs from CNN (feature reduction- F_R in Figure 2) to allow more explainability freedom from the explainable fusion part.

Due to an importance of the explainable features fusion component, we first describe it in section III-A. Section III-B presents the overall picture of the distillation technique, including the feature reduction in a CNN structure, concatenation layers, and the details of the training phases.

A. Explainable Features Fusion

The goal of the explainable features fusion component is to first extract the explanation features from a teacher model and then approximate those valuable features. The explainable fusion component includes a reliable explainer and CAE as illustrated in Figure 4.

Since existing teacher models work in a black-box manner without explicitly explaining why they make their decisions, we will use post-hoc explainers, such as SHAP, LIME, GradCAM, to extract important features that contribute to the model’s output prediction, called explanations. Next, we leverage the idea of convolutional autoencoder (CAE) [20] to learn and transfer the explanations of the teacher to the student.

Unfortunately, the current representation of explanations from existing explainers such as SHAP, LIME, GradCAM is not ready for us to feed into CAEs. Furthermore, CAEs were designed to imitate its images input as nearly as possible to its output, not to approximate and extract important features

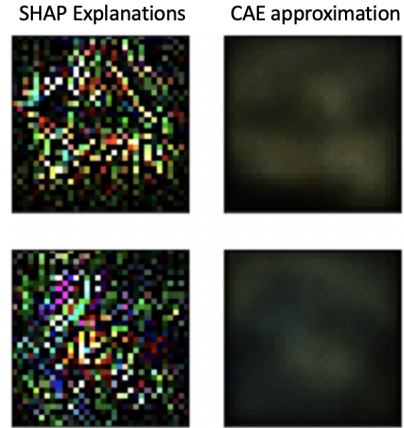


Fig. 3: The output of convolutional autoencoder after feeding the SHAP explanation. The CAE trained on the original SHAP representation. We can observe the sparsity issue on the left images. Hence, the CAE fails to mimic the desired explanations.

to predictions with the lowest possible number of parameters. We will address these two challenges in sections III-A1 and III-A2 accordingly.

1) New Representation for Explanations

For the sake of simple demonstration, we will use SHAP as an explainer for our model as SHAP is one of the state-of-the-arts among existing explainers [7]. In our experiments, we will demonstrate the use of our new explaining representation for other explainers.

The current SHAP representation cannot be approximated by CAEs due to a sparsity issue, as shown in Figure 3. The sparsity issue occurs when the available explanations among the channels are not enough for the CAE to mimic the explanations. Therefore, the CAE fails because decisions of filling those empty areas with SHAP explanations (positively or negatively) cannot be taken randomly. In particular, SHAP explainer does not work smoothly with RGB (Red, Green, and Blue) images. The issue is far different for RGB images in comparison to grayscale images. Figure 5 shows SHAP explanation on top of trained VGG16 model [22] for a random image from CIFAR-10 dataset [23]. As can be seen, the SHAP values are quite sparse and spread between the channels. For each channel, most values are zeros, making the issue like a high-dimensional classification problem. Additionally, the non-zero values are quite high (e.g., -5 , $+3$, etc). In other words, a sparse representation of data is a representation in which few parameters or coefficients are not zero, and many are (strictly) zero.

To overcome the sparsity issue while preserving this XAI goal, we divide an image into smaller superpixels, given the SHAP values. For each superpixel patch, we sum all SHAP values within that patch. The superpixel in this way will describe an actual effect of SHAP values that we will have

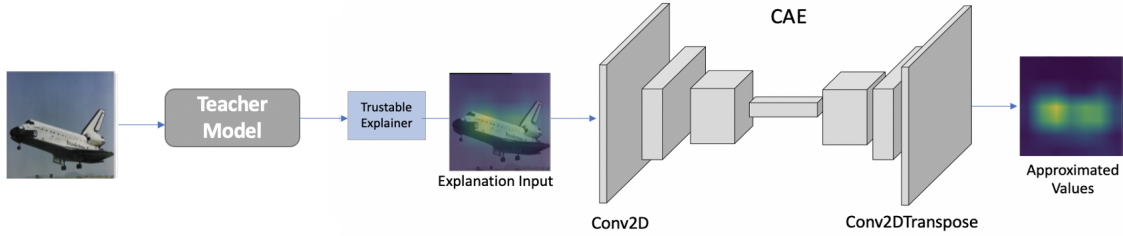


Fig. 4: The overall architecture of the explainable features fusion component. Given an explanation image as input, the fusion component generates the approximated explanation, consisting two main functionalities: (1) a new representation for explanations and (2) an approximation of an image explanations using convolutional autoencoder.

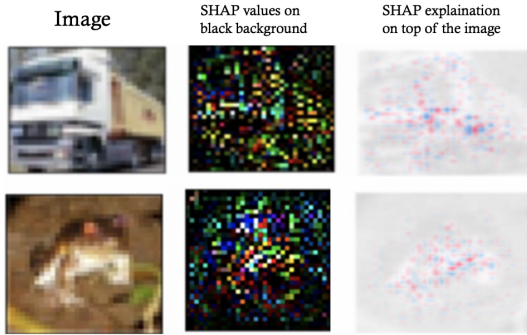


Fig. 5: The sparsity issue of SHAP values. We can observe the sparsity issue between RGB channels for SHAP explanations of the two images. Also, the SHAP explanation of the truck image on the third column does not reveal clear information whether the used model classified the image correctly since the positive and negative values are mixed.

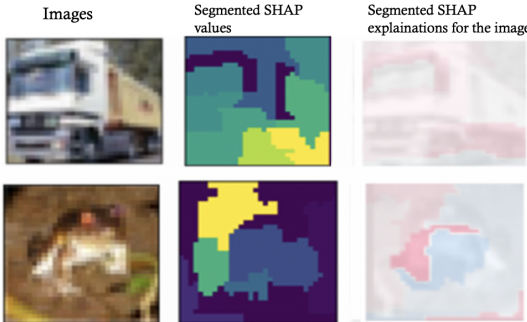


Fig. 6: The new representation of SHAP values. We can observe the improvement of the SHAP explanation of the two images using the new representation. Moreover, we can visually say that the truck image on the third column is classified correctly as the SHAP values react positively with the truck in the image.

on the model's output (either negative or positive). Figure 6 shows the new representation which solves the sparsity issue. The two images in Figure 6 were classified correctly by training the VGG16 model. Then, the SHAP explanations of the two images segmented with different patches (images in

the second column). The new representation of SHAP in the third column clearly demonstrates different dominant areas in the image either positively or negatively, reflecting the model's confidence in classifying those images. For instance, the model was certain that the object in the first image is a truck, as we can see from the positive SHAP values for most patches. However, the model doubts the object in the second image to be a fork.



Fig. 7: Non adequate SHAP segmentation.

In order to do a proper superpixel segmentation to handle the sparsity issue, we adopt a simple linear iterative clustering (SLIC) algorithm [24], which is a local k -mean cluster of 5D pixels. SLIC takes the desired number of approximately equally-sized superpixels k as input and clusters them using a new distance measure that takes superpixel size into account.

Selecting a value k is important. By experiment, k should be in a range of [3, 20]. If $k < 3$, we will have large patches for an image as it is shown in the left image in Fig. 7 which merge the SHAP positive and negative values, and therefore we cannot know what parts contribute to which areas. On the other hand, if $k > 20$, we will have the same sparsity issue for some parts of the image as shown in the right image of Fig. 7. In our experiments shown in Section IV, we set $k = 19$.

2) Convolutional Autoencoder (CAE)

Preliminaries and Notations. Generally, convolutional autoencoders (CAEs) are a form of neural networks that consist of two parts: (1) encoder, which learns to encode the important information of an input into compressed latent representation. (2) decoder, which reconstructs the input again from the latent representation [20]. Mathematically, CAE maps an input $x \in R^d$ to a latent representation $h \in R^l$, where usually $d > l$. The latent representation of the i -th feature map for a single-channel input x is given by

$$h^i = \sigma(x * W^i + b^i) \quad (1)$$

where the bias b is applied to the entire map, σ is an activation function, W is the trainable weights and $*$ denotes 2D convolution layer.

Current Usages and Limitation. The most popular use of CAE is feature extractor for classification [20]. This can be done in two ways: (1) Eq. (1) is used to preserve the most important information from an input via a convolution layer. Then, the latent representation is used with different CNNs for a classification task. (2) the encoder (Eq. (1)) and decoder (reconstruction in CAEs [20]) parts are used, but with replacing some of the convolution layers in the decoder part with fully connected layers to do the classification task.

Since the primary goal of using CAEs is not to do the classification task, most of current CAEs use the fully connected layers as part of their frameworks which increase the total number of parameters. If naively applying this approach, the number of parameters in a student model significantly exceeds that of a teacher model. Also, using existing modified CAEs to do the classification task requires us to use Softmax on the last layer, which makes us lose critical information of explanations.

Our Proposed CAE. To overcome the limitations, we first adopt the convolutional autoencoder as an image transformation function. The input of the auto-encoder x will be an image, and instead of a reverse mapping function, its purpose is to approximate the SHAP explanations of the teacher model $t(x)$, which will be a convolutional neural network. Our objective here is to reconstruct (approximate) SHAP explanation of the teacher model as it is shown in Figure 4.

We use Eq. (1) to calculate the latent representation of the i -th feature map for a single-channel input x . The estimation of the SHAP values e_T of $t(x)$ (the first part of the training phase in Figure 2) is approximated as follows:

$$e_{CAE} = \sigma \left(\sum_{i \in H} h^i * \tilde{W}^i + b \right) \quad (2)$$

where \tilde{W} denotes the flip operations over the dimension weights and H is a set of latent feature map.

We use different rule of thumbs to preserve the parameters constraint for teacher-student model while allowing CAE to approximate SHAP explanations successfully at the same time. First, the non-linearity data processing is the primary goal of an activation function after each layer [25]. Because convolution/deconvolution operations are multiplications in nature, our visualization revealed that the result of the convolution/deconvolution operations (the values of feature maps) increased sharply from layer to layer, preventing the CAE model from converging during training. Thus, the use of hyperbolic tangent activation functions, which constrain the resulting values of feature maps to the interval $[-1,1]$, sets appropriate limits on the values of feature maps at the end of the decoder part as well as maintaining the SHAP valuable information since SHAP values will be within the range for

a model that outputs a probability [26], and provides good convergence of the whole model.

Second, the basic structure of the autoencoder is extended by changing the completely connected layers to convolution layers in the convolutional autoencoder because (1) leveraging convolutional operations not only for slower training time and chances of reducing overfitting but also maintaining a reasonable total number of parameters. For instance, the total number of parameters for the small CAE version in Table I (Section IV) are reduced from 279,872 to only 1077. (2) The 2D image structure is ignored by fully connected AEs or partially by fully connected layers in CAEs. This not only creates redundancy in the parameters while handling inputs of exact magnitude.

A sample of successful SHAP approximation by the fusion component is shown in Figure 4 where SHAP explainer is first employed on top of a correct prediction of the teacher model for a plane image. The new representation of SHAP values is then prepared as described in Section III-A1 and passed to our new CAE, which produces the explainable features that describe the plane visually, as it can be seen in Figure 4.

B. XDistillation Model

We are now ready to describe the rest of our proposed XDistillation model. Figure 2 demonstrates the overall architecture of XDistillation model, which consists of (1) explainable feature fusion component, and (2) CNN model. We use the first component, our CAE model, to transfer the explanations knowledge from a teacher model to the CNN model. Additionally, we use the feature reduction F_R to allow more explainability freedom associated with the first component.

For a given image x , a teacher model generates a vector of prediction scores

$$s^t(x) = [s_1^t(x), s_2^t(x), \dots, s_u^t(x)],$$

where t refers to the teacher model, x is the input and u is the value of the scores which are then transformed to probabilities:

$$p_u^t(x) = \frac{\exp^{s_u^t(x)}}{\sum_j \exp^{s_j^t(x)}}.$$

Since the trained neural networks generate probability distributions with peaks that might be less instructive, we utilize the method in [1] for temperature scaling to "Smooth" those probabilities:

$$\tilde{p}_u^t(x) = \frac{\exp^{s_u^t(x)/\tau}}{\sum_j \exp^{s_j^t(x)/\tau}}, \quad (3)$$

where $\tau > 1$ is a temperature hyperparameter.

Similarly, the CNN model inside the XDistillation model in Figure 2 also returns a smooth output denoted by $\tilde{p}^s(x)$ where s refer to the CNN model.

Figure 2 illustrates 2 training phases of XDistillation. In the first training phase, SHAP explainer is used to obtain ground truth explanations of $t(x)$. Following that, We train CAE using

the Mean Absolute Error (MAE) as follows (L_{XD} loss in Figure 2):

$$L_{XD} = \frac{1}{2n} \sum_{i=1}^n (e_{CAE_i} - e_{T_i}), \quad (4)$$

where e_{CAE_i} is described in Eq. (2) and e_T is the ground truth explanations of $t(x)$. The weights are then updated with classic backpropagation with stochastic gradient descent.

With the trained CAE model, we feed its output to a CNN; the resulting representations are flattened and concatenated before they are used as an input of the classifier part of the model. An illustration of our architecture can be found in Figure 2.

In the second training phase, we use the knowledge distillation loss (L_{KD} loss in Figure 2) which can be defined as follows:

$$\mathcal{L}_{KD} = -\tau^2 \sum_b p_u^t(x) \log p_u^s(x) \quad (5)$$

We extended the objective function with feature reduction F_R for the last layer after the concatenation ca . In combination with the knowledge distillation loss defined in Equation 5, this yields the total loss function:

$$\mathcal{L} = \mathcal{L}_{cls} + (1 - \alpha)\mathcal{L}_{KD} + F_R(\lambda \|c\|_2^2) \quad (6)$$

where λ is signal controller hyperparameter, and c is the weight of last layer of the concatenation part. The feature reduction F_R adds coefficient squared magnitude to the loss function as penalty term in last layer (concatenation layer as it is shown in Figure 2). The cross entropy function is defined as:

$$L_{cls} = - \sum_{a=1}^M y \log(p(y_a))$$

where M is the number of classes, y is the truth labels and $p(y_a)$ is the softmax probability for a^{th} class.

The λ in F_R acts as a signal gate controller between the incoming signals from both CNN and the feature fusion component. In other words, if $\lambda = 0$, then the signals pass the gate without constraints. However, if λ is large, then the incoming signals from the fusion component have more freedom to pass the gate and feed into the CNN model with implicit interpretable features.

Therefore, we increase the feature reduction value to constrain the coming signals from the CNN model, which computes various kinds of features in an image such as edges, curves. This gives more freedom to approximated explanations from the explainable fusion component. Those features carry important explanations features of the teacher model. Then we merge the approximated SHAP explanation features of the teacher model with CNN features as follow:

The operation of concatenation layer in the context of tensors, is the operation of joining tensors along one dimension using Linear layer. In particular, given $x \in \mathbb{R}^{n \times g}$ and

$y \in \mathbb{R}^{m \times g}$, where x is CAE features, and y is CNN features. The resulting tensor after concatenation ca will then be:

$$\begin{aligned} cat(x, y) &= [x_1, \dots, x_n, y_1, \dots, y_m] \\ &\text{where } x_i \in \mathbb{R}^{1 \times g} \\ &\text{and } y_i \in \mathbb{R}^{1 \times g} \end{aligned}$$

Figure 2 illustrates the details of XDistillation. It is crucial to note that in Xdistillation at the testing phase, we use the frozen weights of the CAE because the CAE trains separately and only once. The purpose behind using the pre-trained CAE and not train it with CNN model is that: (1) reducing the training time of using our model with different types of student models. (2) providing flexibility to our method to be used with various explainers without retraining the CAE. Training the explanations of only one explainer, SHAP values, for teacher model in Figure 4 and then use it in XDistillation improves the explanation output of the proposed model, which is similar to the actual explanations of the teacher model.

IV. EXPERIMENTAL ANALYSIS

In the following experiments, we analyze XDistillation in terms of predictive accuracy and explanation similarity (testing phase in Figure 2). In this section, we first describe our experimental setup and description of testing datasets. Next, we show the benefits of the explanation infusion step on the predictive performance of different models and compare XDistillation with other state-of-the-art distillation methods in term of test accuracy. Finally, we evaluate the consistency of the student model’s explanations generated by different knowledge distillation methods.

TABLE I: The hyper-parameters Details of the convolutional autoencoders

CAE_Model	Type	Kernel size	Stride	Padding	#parameters
Small_CAE	Conv2d	3		1	1077
	MaxPool2d	2	2		
	Conv2d	3		1	
	MaxPool2d	2	2		
	ConvT	2	2		
	ConvT	2	2		
Large_CAE	Conv2d	4	2	1	198,796
	Conv2d	4	2	1	
	Conv2d	4	2	1	
	Conv2d	4	2	1	
	ConvT	4	2	1	
	ConvT	4	2	1	
	ConvT	4	2	1	
	ConvT	4	2	1	

A. Experiment Setting

Dataset. Our experiments are conducted on MNIST [27] and CIFAR-10 [23] datasets. The MNIST dataset contains images of 28×28 grayscale hand-written digits. In our experiments, the training-testing split is 60,000 : 10,000. On the other hand, CIFAR-10 consists of 60,000 32×32 colour images in 10 classes, with 6000 images per class. There are 50,000 training images and 10,000 test images.

TABLE II: Model Generalization

Dataset	autoencoder	Model	Accuracy %	#parameters
MNIST	Small_CAE	Teacher(LeNet5)	99	61,701
		Baseline student(Net)	97	8,297
		KD (Net)	98	8,297
		Xdistillation (Net)	98	17,386
CIFAR-10	Large_CAE	Teacher (VGG16)	93.78	14,728, 266
		Baseline student(VGG7)	89.2	2,781,386
		KD (VGG7)	90.2	2,781,386
		Xdistillation (VGG7)	90.9	3,521,276
CIFAR-10	Large_CAE	Teacher(VGG16)	93.78	14,728, 266
		Baseline student (VGG11)	92	9,231,114
		KD (VGG11)	92.41	9,231,114
		Xdistillation (VGG11)	92.59	9,624,598

Teacher Models. In our experiments, the teacher models are implemented using LeNet5 [28] and VGG16 architecture [22] for the MNIST and CIFAR-10 dataset respectively. The LeNet is trained using MSE loss and Adam optimizer with learning rate 0.001 and number of epochs 150. In CIFAR-10, our training included data augmentations with random horizontal flips and random crops of size 32 with a possible padding of 4 pixels. The training utilizes with a batch size of 128 for 500 epochs. The optimizer is stochastic gradient descent with a momentum of 0.9 and a learning rate schedule of 0.1, 0.01 and 0.001 starting from the epochs 0, 150 and 250 respectively. The models’ accuracy and number of parameters are shown in Table III.

Student Models. The student models in our experiments are smaller LeNet5 for MNIST and VGG7 [22] for CIFAR-10. Their testing accuracy and number of parameters can also be found in Table III. The student models are trained using some state-of-the-art distillation knowledge techniques, including knowledge distillation (KD) [1], attention transfer (AT) [8], neural selective transfer (NST) [29], and activation boundary (AB) [9].

XDistillation Model. To train XDistillation, we first use SHAP [7], LIME [16] and GradCAM [14] explainers to generate the teacher’s explanations. Then, the explanations are transformed as described in section III-A1 and the convolutional autoencoder as described in section III-A2 is trained based on these transformed explanations. The autoencoder is trained using the L1 loss and the Adam optimizer with a learning rate of 0.001. The detail of the large CAE is shown in Table I. The CNNs used in XDistillation have the same architecture and parameters as the student models. The CNN is trained using the loss Eq. (6). We typically use $\alpha = 0.9$ and temperature $\tau = 1$.

One key parameter that determines XDistillation’s performance is the choice of the regularized parameter λ controlling the model’s signal from the CNN and the CAE as described in equation (6). Figure 8 shows the accuracy of XDistillation during training with different values of λ . In the following experiments, we use $\lambda = 5e - 4$.

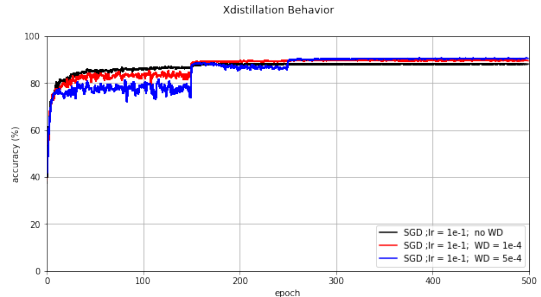


Fig. 8: The testing accuracy of XDistillation during training with different regularized parameter λ (equation 6)

TABLE III: Performance Comparison

Model	Accuracy %	#parameters
Teacher	93,78	14,728, 266
Baseline student	89.2	2,781,386
Knowledge distillation (KD)	90.2	2,781,386
Attention transfer (AT)	90	2,781,386
Neural selective transfer (NST)	89.47	2,781,386
Activation boundary (AB)	89.36	2,781,386
Xdistillation	90.9	3,521,276

B. Results

Predictive performance. We first demonstrate the advantage of the explanation infusion step via CAE by comparing the test accuracy of XDistillation models with the baseline student models and the KD models. As shown in Table II, the introduction of the CAE generally increases the overall predictive accuracy of the student models.

Table III shows the predictive performance of the XDistillation model and compares it with other distillation knowledge methods. We can see that XDistillation outperform the other models in term of performance with an accuracy of 90.9%. The second close model to ours is the KD with 90.2% accuracy, whereas the performance of the remaining models was below 90%. Generally, the accuracy improved in comparison to the baseline student model and other state-of-the-art methods.

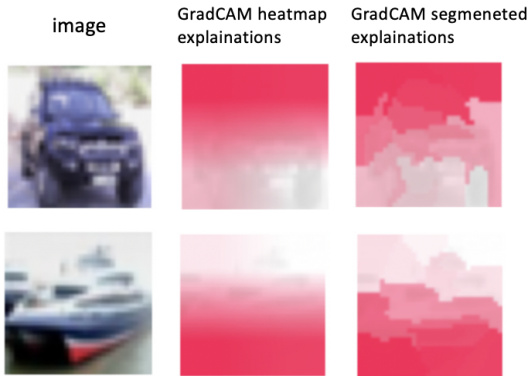


Fig. 9: GradCAM representations where the red color indicates the important features

TABLE IV: Scoring Distance

MSE	Explainer	SHAP	GradCAM	LIME
	KD	0.027	0.025	0.006
AT	0.040	0.239	0.014	
NST	0.037	0.246	0.015	
AB	0.033	0.199	0.014	
Xdistillation (our)	0.026	0.024	0.006	

Explanation consistency. We now examine the ability of XDistillation model in providing consistent explainable features to the teacher model. Here, we use SHAP, GradCAM and, LIME explainers for our evaluation. Since there are some differences in the representation of the explainers’ output. In particular, GradCAM proliferates as heatmap over an image, shown in Figure 9. On the other hand, LIME explainer utilizes quickshift segmentation approach [30]. Thus, for a fair comparison, we apply the explanation’s transformation step described in subsection III-A1 onto GradCAM and the SLIC segmentation procedure on LIME.

Given explanations of the same input of the same predictions from two models, we use the mean square error (MSE) [31] to measure the similarity between the two explanations. Table IV reports the average MSE on 10,000 pairs of explanations in CIFAR-10 experiments. In general, we can observe that the explanations of XDistillation are more similar

TABLE V: Scoring explanations

Model	Explainer	SHAP	GradCAM	LIME
	KD	68.3%	56.1%	50.1%
AT	67.9%	52.1%	44.3%	
NST	67.6%	58.6%	44.4%	
AB	67.3%	55.6%	44.7%	
Xdistillation (our)	70.5%	62%	54.4%	

to the teacher than other methods.

Since different explainers return different representations, beside MSE, we also measure the explanation consistency by measuring the overlapping area of explained super-pixels returned by each explainer. Specifically, for all correctly classified inputs by both teacher and student models, we collect the k most important super-pixels of the corresponding explanations and compute the $sign(x)$ function as follows:

$$h(x) = \begin{cases} 1 & \text{where } sign(M(x)) = sign(T(x)) \\ 0 & \text{otherwise} \end{cases}$$

where $M(x)$ is the explanation value of overlapping SLIC super-pixels patch that classified correctly by the student model, $T(x)$ is the explanation number of overlapping SLIC super-pixels patch that classified correctly by the teacher model, and $sign(\cdot)$ is a function where it returns 1 if the explanation segments of the teacher and the student have same sign. We then compare the counts of the correct explanation patches $h(x)$ for both models, our XDistillation and other KD, based on the baseline model (VGG7), with normalization.

Table V shows the overlapping-score of XDistillation and other methods. As can be seen, our proposed model has more similar explanations to the teacher model in comparison to knowledge distillation methods and baseline model. For instance, the overlap area of GradCAM for XDistillation method improved by about 6% compare to the second close model to our method in term of the explanation performance which means our model is more similar to the teacher model. For reference, Figure 11 provides examples of explanations of some predictions generated by SHAP, LIME and GradCAM along with the corresponding overlapping-scores.

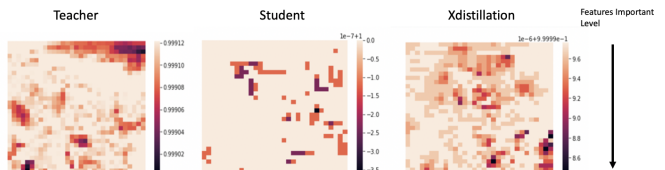


Fig. 10: Occlusion Analysis

Occlusion Analysis. Occlusion analysis is a simple method for understanding which parts of an image are most important for a model. Small perturbations of the data can be used to measure a network’s sensitivity to occlusion in different regions of the data. The process for an occlusion experiment is as follows: we first mask part of an image before feeding it into three different trained models, including the (VGG16), our proposed model, and KD (VGG7) [1]. Then, we draw a heatmap of class scores for each masked image. Lastly, we slide the masked area to a different spot and repeat the step process again until we cover all the images. The reasoning behind all these is that if the class score for a partially occluded image is different than the true class, then the occluded area was likely very important.

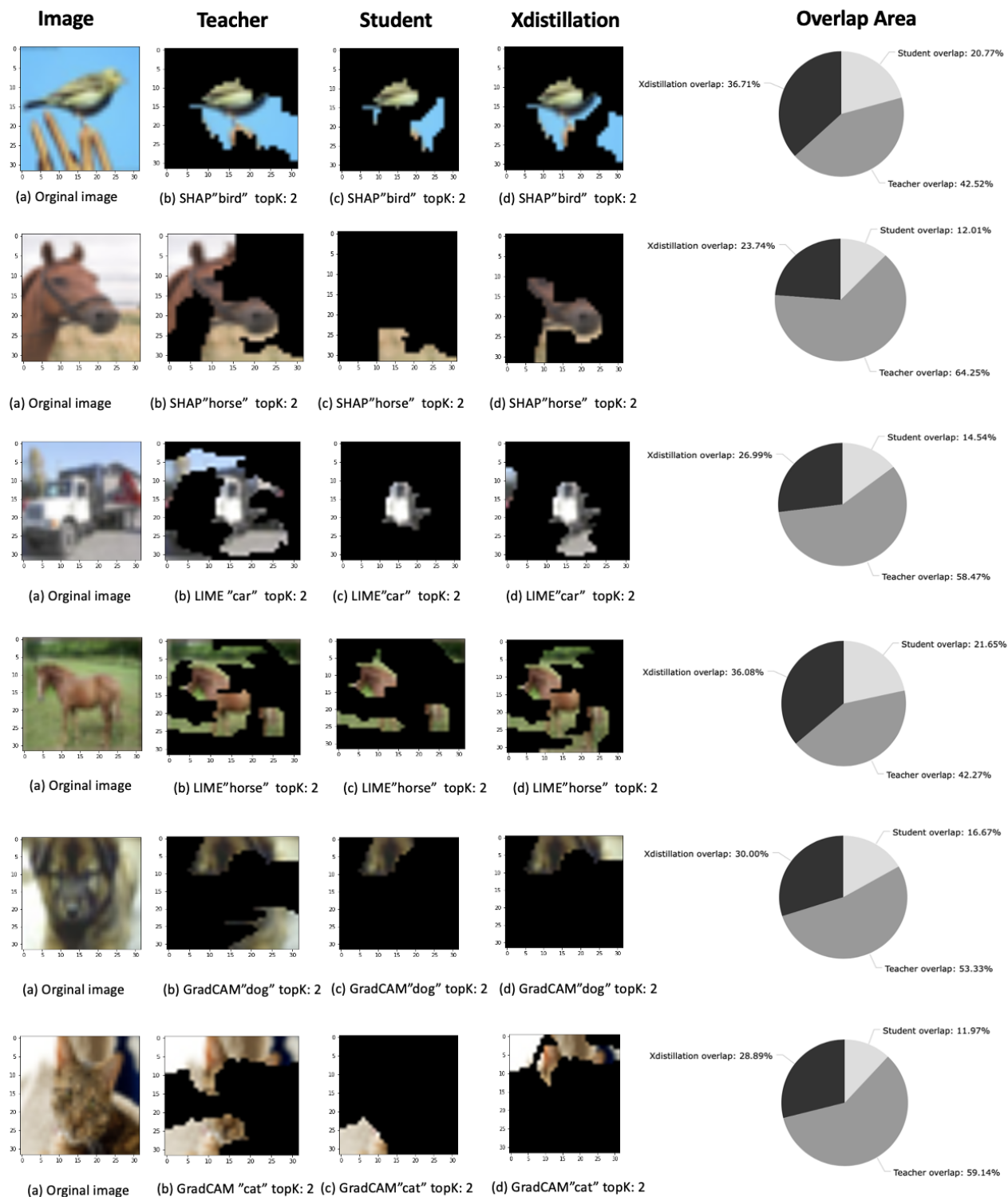


Fig. 11: We use SHAP, LIME and GradCAM to compare the overlap area of teacher, student and our proposed method to a original image for topK: 2

Figure 10 shows samples of the occlusion for the three models. In our human-subjective test, we apply the experiment for randomly 30 samples for each model whose predictions made *correctly*. thirty users were asked to choose which two of the three models are most likely similar to each other based on heatmap features. 18 out of 30 select the most likely models to each other are the teacher and XDistillation model where 5 individuals choose XDistillation and KD model. We can observe that the the most important features captured by the two models, teacher and XDistillation, are most likely similar to each other.

V. CONCLUSION

Knowledge distillation (KD) tackles the issue of transferring knowledge from a vast and complicated neural network to a smaller one. The standard methodology reduces the KD divergence between a teacher and student model’s outputs. However, current KD techniques ignore an important information of an explanatory network of the teachers. In this paper, we present XDistillation, explanatory model that mimics teacher explanations. We show experimentally that our proposed model outperforms the existing KD techniques. The Xdistillation method offers the ability to eliminate incoherence of the explanations between the teacher and the student patterns apart from current KD strategies.

ACKNOWLEDGMENT

This work was supported in part by the National Science Foundation Program on Fairness in AI in collaboration with Amazon under award No. 1939725.

REFERENCES

- [1] G. Hinton, O. Vinyals, and J. Dean, “Distilling the knowledge in a neural network,” in *NIPS Deep Learning and Representation Learning Workshop*, 2015.
- [2] S. Han, J. Pool, J. Tran, and W. J. Dally, “Learning both weights and connections for efficient neural network,” in *NIPS*, 2015.
- [3] J. H. Cho and B. Hariharan, “On the efficacy of knowledge distillation,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 4794–4802.
- [4] R. Alharbi, M. N. Vu, and M. T. Thai, “Evaluating fake news detection models from explainable machine learning perspectives,” in *ICC 2021-IEEE International Conference on Communications*. IEEE, 2021, pp. 1–6.
- [5] A. B. Arrieta, N. Díaz-Rodríguez, J. Del Ser, A. Bennetot, S. Tabik, A. Barbado, S. García, S. Gil-López, D. Molina, R. Benjamins *et al.*, “Explainable artificial intelligence (xai): Concepts, taxonomies, opportunities and challenges toward responsible ai,” *Information Fusion*, vol. 58, pp. 82–115, 2020.
- [6] R. Geirhos, J.-H. Jacobsen, C. Michaelis, R. Zemel, W. Brendel, M. Bethge, and F. A. Wichmann, “Shortcut learning in deep neural networks,” *Nature Machine Intelligence*, vol. 2, no. 11, Nov 2020.
- [7] S. M. Lundberg and S.-I. Lee, “A unified approach to interpreting model predictions,” in *NIPS*, 2017.
- [8] N. Komodakis and S. Zagoruyko, “Paying more attention to attention: improving the performance of convolutional neural networks via attention transfer,” in *ICLR*, 2017.
- [9] B. Heo, M. Lee, S. Yun, and J. Y. Choi, “Knowledge transfer via distillation of activation boundaries formed by hidden neurons,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, no. 01, 2019, pp. 3779–3787.
- [10] F. Tung and G. Mori, “Similarity-preserving knowledge distillation,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 1365–1374.
- [11] M. Vu and M. T. Thai, “Pgm-explainer: Probabilistic graphical model explanations for graph neural networks,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 12 225–12 235, 2020.
- [12] Z. C. Lipton, “The mythos of model interpretability: In machine learning, the concept of interpretability is both important and slippery,” *Queue*, vol. 16, no. 3, pp. 31–57, 2018.
- [13] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, and A. Torralba, “Learning deep features for discriminative localization,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2921–2929.
- [14] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, “Grad-cam: Visual explanations from deep networks via gradient-based localization,” in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.
- [15] R. Fong, M. Patrick, and A. Vedaldi, “Understanding deep networks via extremal perturbations and smooth masks,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 2950–2958.
- [16] M. T. Ribeiro, S. Singh, and C. Guestrin, “‘‘ why should i trust you?’’ explaining the predictions of any classifier,” in *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, 2016, pp. 1135–1144.
- [17] Y. Wang, H. Yao, and S. Zhao, “Auto-encoder based dimensionality reduction,” *Neurocomputing*, vol. 184, pp. 232–242, 2016.
- [18] K. G. Lore, A. Akintayo, and S. Sarkar, “Llnet: A deep autoencoder approach to natural low-light image enhancement,” *Pattern Recognition*, vol. 61, pp. 650–662, 2017.
- [19] J. Chen, S. Sathe, C. Aggarwal, and D. Turaga, “Outlier detection with autoencoder ensembles,” in *Proceedings of the 2017 SIAM international conference on data mining*. SIAM, 2017, pp. 90–98.
- [20] J. Masci, U. Meier, D. Cireşan, and J. Schmidhuber, “Stacked convolutional auto-encoders for hierarchical feature extraction,” in *International conference on artificial neural networks*. Springer, 2011, pp. 52–59.
- [21] A. Ashfahani, M. Pratama, E. Lughofer, and Y.-S. Ong, “Devdan: Deep evolving denoising autoencoder,” *Neurocomputing*, vol. 390, pp. 297–314, 2020.
- [22] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” in *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, Y. Bengio and Y. LeCun, Eds., 2015.
- [23] A. Krizhevsky, G. Hinton *et al.*, “Learning multiple layers of features from tiny images,” 2009.
- [24] W. F. Noh and P. Woodward, “Slic (simple line interface calculation),” in *Proceedings of the fifth international conference on numerical methods in fluid dynamics June 28–July 2, 1976 Twente University, Enschede*. Springer, 1976, pp. 330–340.
- [25] M. Ranzato, F. J. Huang, Y.-L. Boureau, and Y. LeCun, “Unsupervised learning of invariant feature hierarchies with applications to object recognition,” in *2007 IEEE conference on computer vision and pattern recognition*. IEEE, 2007, pp. 1–8.
- [26] S. M. Lundberg and S.-I. Lee, “A unified approach to interpreting model predictions,” in *Advances in Neural Information Processing Systems 30*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, Eds. Curran Associates, Inc., 2017, pp. 4765–4774. [Online]. Available: <http://papers.nips.cc/paper/7062-a-unified-approach-to-interpreting-model-predictions.pdf>
- [27] L. Deng, “The mnist database of handwritten digit images for machine learning research,” *IEEE Signal Processing Magazine*, vol. 29, no. 6, pp. 141–142, 2012.
- [28] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, “Gradient-based learning applied to document recognition,” *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [29] Z. Huang and N. Wang, “Like what you like: Knowledge distill via neuron selectivity transfer,” *arXiv preprint arXiv:1707.01219*, 2017.
- [30] A. Vedaldi and S. Soatto, “Quick shift and kernel methods for mode seeking,” in *European conference on computer vision*. Springer, 2008, pp. 705–718.
- [31] H. Rezatofoghi, N. Tsoi, J. Gwak, A. Sadeghian, I. Reid, and S. Savarese, “Generalized intersection over union: A metric and a loss for bounding box regression,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 658–666.