

TVE: Learning Meta-attribution for Transferable Vision Explainer

Guanchu Wang¹, Yu-Neng Chuang¹, Fan Yang², Mengnan Du³, Chia-Yuan Chang⁴,
Shaochen Zhong¹, Zirui Liu¹, Zhaozhuo Xu⁵, Kaixiong Zhou⁶, Xuanting Cai⁷, Xia Hu¹

¹Rice University, ²Wake Forest University, ³New Jersey Institute of Technology, ⁴Texas A&M University,

⁵Stevens Institute of Technology, ⁶North Carolina State University, ⁷Meta Platforms, Inc.

{gw22, yc146, shaochen.zhong, zl105, xia.hu}@rice.edu; yangfan@wfu.edu; mengnan.du@njit.edu;
cychang@tamu.edu; zxu79@stevens.edu; kzhou22@ncsu.edu; caixuanting@meta.com

Abstract

Explainable machine learning significantly improves the transparency of deep neural networks. However, existing work is constrained to explaining the behavior of individual model predictions, and lacks the ability to transfer the explanation across various models and tasks. This limitation results in explaining various tasks being time- and resource-consuming. To address this problem, we introduce a **T**ransferable **V**ision **E**xplainer (TVE) that can effectively explain various vision models in downstream tasks. Specifically, the transferability of TVE is realized through a pre-training process on large-scale datasets towards learning the meta-attribution. This meta-attribution leverages the versatility of generic backbone encoders to comprehensively encode the attribution knowledge for the input instance, which enables TVE to seamlessly transfer to explain various downstream tasks, without the need for training on task-specific data. Empirical studies involve explaining three different architectures of vision models across three diverse downstream datasets. The experimental results indicate TVE is effective in explaining these tasks without the need for additional training on downstream data.

1 Introduction

Explainable machine learning (ML) contributes to enhancing the transparency of deep neural networks (DNNs) for human comprehension [12]. It significantly facilitates the deployment of DNNs to high-stake scenarios where model explanations are required, such as loan approvals [32], healthcare [2], and targeted advertisement [41]. In these fields, explainable DNN decisions are particularly important, given the practical needs of stakeholders and regulatory requirements, such as the General Data Protection Regulation (GDPR) [15].

To overcome the black-box nature of DNNs, existing work of explainable ML can be categorized into two groups. The first group of work focuses on constructing local explanation based on perturbation of the target black-box model, like LIME [29], GradCAM [31], and Integrated Gradient [35]. These pieces of work rely on resource-intensive procedures like sampling or backpropagation of the target black-box model [24], leading to undesirable trade-off between the efficiency and interpretation fidelity [7]. Another group leverages the knowledge of explanation values to train DNN-based explainers, such as FastSHAP [20], CORTX [8], and LARA [30, 37]. Such arts capable of efficiently generating explanations for an entire batch of instances through a single, streamlined feed-forward operation of the DNN-based explainer. However, they are constrained to explaining individual black box models, and often lack the ability to transfer the explainer across various models or tasks. These constraints lead to a time and resource-intensive process in practical scenarios, as they require the development and training of separate explainers for each specific task.

To address the lack of transferability in explainers, we introduce a **T**ransferable **V**ision **E**xplainer (TVE). The primary goal of TVE is to achieve transferability through a pre-training process on large-scale image datasets, such that it can seamlessly explain various downstream tasks, as long as such tasks are within the scope of pre-training data distribution. The construction of such transferable explainers introduces two non-trivial challenges: **CH1**. Without task-specific exposure during the pre-training, how to ensure the universal effectiveness of explainer for various downstream tasks? **CH2**. How to adapt the explainer to a specific task without fine-tuning on the task-specific data?

Our work effectively tackles these challenges. To address CH1, we introduce a novel concept, named *meta-attribution*, as a foundation for explaining various downstream tasks. Specifically, the meta-attribution

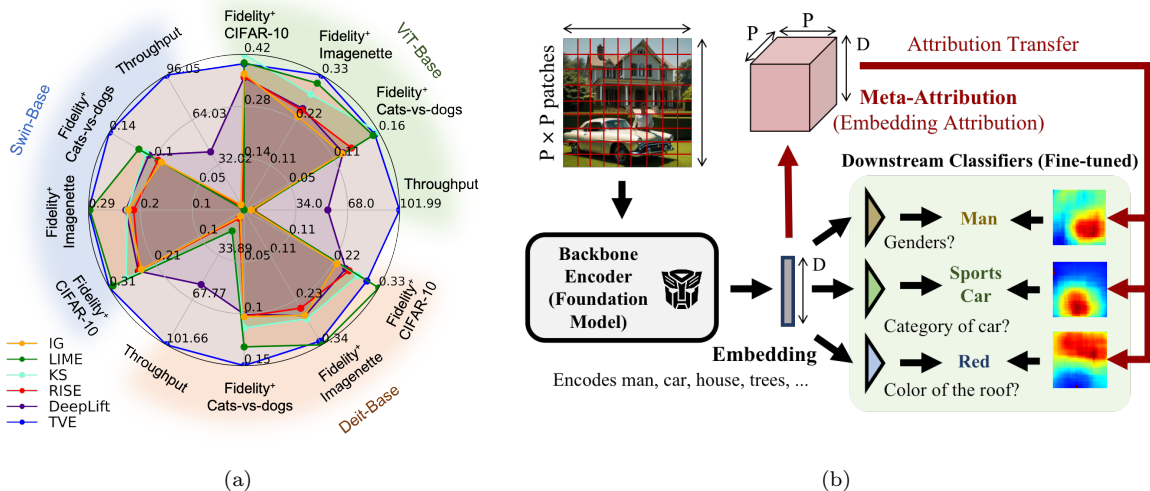


Figure 1: Performance of TVE in explaining ViT-B, Swin-B, and Deit-B on the Cats-vs-dogs, Imagenette, and CIFAR-10 datasets. $Fidelity^+$ score refers to the area under $Fidelity^+$ -sparsity curve. (b) Illustration of attribution transfer. In this framework, the backbone can be a ViT encoder; and the downstream classifiers can be MLPs. The embedding vector comprehensively encodes the features of input image. Motivated by this, the meta-attribution comprehensively encapsulates the importance of each input patch to each element of the embedding vector. This enables it to seamlessly transfer for explaining various downstream tasks.

versatilely encodes the attribution knowledge for the input instance via exhaustively attributing each dimension of instance embedding. This knowledge is reusable for explaining various downstream tasks. It guides the pre-training of TVE on large-scale image datasets, ensuring the universal effectiveness of TVE. After the pre-training, in response to CH2, we propose a *transfer rule* to adapt the meta-attribution to explaining downstream tasks, without the need for additional training on task-specific data. Figure 1(a) shows the comprehensive performance of TVE pre-trained on the *ImageNet* dataset and transferred to the *Cats-vs-dogs*, *Imagenette*, and *CIFAR-10* datasets, where TVE shows competitive fidelity and efficiency compared with state-of-the-art methods. To summarize, our work makes the following contributions:

- **Attribution transfer.** We propose a framework of attribution transfer, with a meta-attribution as foundations, and a transfer rule for explaining the downstream tasks.
- **Transferable explainer.** We build a transferable explainer TVE that explains various downstream tasks without the need for training on the task-specific data.
- **Theoretical foundation.** We validate the pre-training of TVE can minimize the explanation error bound aligned with the \mathcal{V} -information-based explanation.
- **Competitive performance in explaining various downstream tasks.** The pre-trained TVE shows promising results in explaining three architectures of vision Transformer across three downstream datasets. Significantly, the strong transferability of TVE facilitates efficient and flexible deployment to various downstream scenarios.

2 Notations

We introduce the notations for the problem formulation.

Target model. We focus on the explanation of vision models: $\mathcal{X} \rightarrow \mathcal{Y}_t$ in this work, where $\mathcal{X} = \mathbb{N}^{W \times W}$ denote the spatial space of $W \times W$ pixels; \mathbb{N} denote the space of a single pixel with three channels; and \mathcal{Y}_t denotes the label space. Moreover, we follow most of existing work [16] and implementation of DNNs [39] to consider the target model as $f_t = H_t \circ G$, where the backbone encoder $G(\bullet): \mathbb{N}^{W \times W} \rightarrow \mathbb{R}^D$ is pre-trained on

large-scale datasets; and the classifier $H_t(\bullet): \mathbb{R}^D \rightarrow \mathcal{Y}_t$ is finetuned on a specific task t . It maybe worth noting that although we follow the transfer learning setting [6, 5] to freeze the backbone encoder $G(\bullet)$ during the finetuning of f_t . Our experiment results in Section 6.3 further show that the proposed transferable explanation framework also shows effectiveness in the scenario where the target model is fully fine-tuned on downstream data.

Image Patching. We follow existing work [26] to consider the patch-wise attribution of model prediction, i.e. the importance of each patch. Specifically, we follow existing work [8, 20] to split each image \mathbf{x}_k into $P \times P$ patches in a grid pattern, where each patch has $C \times C$ pixels; and $W = CP$. Let $\mathcal{Z}(\mathbf{x}_k) = \{z_{i,j} | 1 \leq i, j \leq P\}$ denote the patches of an image $\mathbf{x}_k \in \mathcal{X}$, where a patch $z \in \mathbb{N}^{C \times C}$ aligns with continuous $C \times C$ pixels of the image. Moreover, we define $\mathcal{N}(z) \subseteq \mathcal{Z}(\mathbf{x}_k)$ as the neighbors of a patch z within the grid space, because a patch together with its neighbors have richer semantic content for model explanation. In this work, we follow the vision transformer [11] to split the image patches with $P = 14$ for 224×224 input images from the ImageNet dataset; and we consider $\mathcal{N}(z)$ as the zero-, one-, and two-hop neighbors of the patch z .

Model Perturbation. $f_t(\mathcal{S}; \mathbf{x}_k, y)$ represents the output of f_t on class y , with a perturbed instance as the input. The patch subset $\mathcal{S} \subseteq \mathcal{Z}(\mathbf{x}_k)$ controls the perturbation. Specifically, the pixels belonging to the patches $z \in \mathcal{Z}(\mathbf{x}_k) \setminus \mathcal{S}$ are removed and take 0, which is approximately the average value of normalized pixels. For example, $f_t(\mathcal{N}(z); \mathbf{x}_k, y)$ defines the output of f_t based on the perturbed input, where the pixels not belonging to the neighbors of patch z take 0.

Feature Attribution. This work focuses on the feature attribution of target models f_t for providing explanations. The feature attribution process involves generating importance scores, denoted as $\phi_{k,y,z}$ for each patch $z \in \mathcal{Z}(\mathbf{x}_k)$ of the input image $\mathbf{x}_k \in \mathcal{X}$, to indicate its importance to the model prediction $f_t(\mathcal{Z}(\mathbf{x}_k); \mathbf{x}_k, y)$ on class y .

3 Feature Attribution can Transfer

The motivation behind attribution transfer stems from model transfer in vision tasks [6, 17, 16]. Specifically, it arises from the observation that a generic backbone encoder possesses the capability to capture essential features of input images and represent them as embedding vectors. This versatility enables the backbone to effectively adapt to a wide range of downstream tasks within the scope of pre-training data distribution. As shown in Figure 1(b), information of **man**, **car**, **house** encoded in the embedding vector enables the detection of **gender**, **car**, and **building** in three different downstream scenarios, respectively. Despite the demonstrated transferability of the backbone encoder, existing research has challenges in achieving ‘transferable explainer’ across different tasks. To bridge this gap and streamline the explanation process, we propose a *meta-attribution* that can be applied across various tasks, resulting in a significant reduction in the cost associated with generating explanations.

The *meta-attribution* is defined as a tensor that versatilely encodes the reusable attribution knowledge for explaining downstream tasks. As shown in Figure 1(b), we illustrate the meta-attribution as a three-dimensional tensor. A simple and effective method in this work is attributing the importance of input patches to each element of the embedding vector for the meta-attribution. As shown in Figure 1(b), each $P \times P$ slice of this tensor corresponds to $P \times P$ patches within the input image, encoding their importance to a specific dimension of the embedding vector. In this way, the meta-attribution inherits the adaptability of the embedding vector, making it versatile enough to adapt various explanation tasks in downstream scenarios. For instance, the meta-attribution encodes the attribution knowledge for the **man** and **car** components encoded in the embedding vector, such that it can transfer to explain the **car** classification and **gender** detection in downstream scenarios. The versatility of meta-attribution can effectively address the CH1 described in Section 1. We formalize the attribution transfer in Sections 4.

4 Meta-attribution Transfer

In this section, we begin with the explanation definition by following the \mathcal{V} -information theory [40, 18, 3]. Then, we introduce the definition of meta-attribution in Definition 1. Finally, we propose a transfer rule to adapt the meta-attribution to explaining specific downstream tasks in Definition 2.

4.1 \mathcal{V} -Information-based Explanation

The importance of a patch $z \in \mathcal{Z}(\mathbf{x}_k)$ to downstream model $f_t(\mathbf{x}_k)$ is formulated into the conditional mutual information $I(\mathcal{N}(z); Y_t | B)$ between $\mathcal{N}(z)$ and Y_t , given the state of remaining patches $B \subseteq \mathcal{Z}(\mathbf{x}_k) \setminus \mathcal{N}(z)$ [3]. Here, $Y_t \sim f_t(\mathbf{x}_k)$ denotes the variable corresponding to the model output. However, estimating this mutual information accurately poses a challenge due to the unknown distribution of $\mathcal{N}(z)$ and B . To address this challenge, we adopt an information-theoretic framework introduced in works by [40, 18], known to as conditional \mathcal{V} -information $I_{\mathcal{V}}(\mathcal{N}(z) \rightarrow Y_t)$. In particular, it redirects the computation of mutual information to a certain predictive model within function space \mathcal{V} , as defined by:

$$I_{\mathcal{V}}(\mathcal{N}(z) \rightarrow Y_t | B) = H_{\mathcal{V}}(Y_t | B) - H_{\mathcal{V}}(Y_t | \mathcal{N}(z), B),$$

where \mathcal{V} -entropy $H_{\mathcal{V}}(Y_t | B)$ takes the lowest entropy over the function space \mathcal{V} , which is given by

$$H_{\mathcal{V}}(Y_t | B) = \inf_{f \in \mathcal{V}} \mathbb{E}_{y \sim \mathcal{Y}_t} [-\log f(B; \mathbf{x}_k, y)]. \quad (1)$$

Note that the \mathcal{V} -information explanation should align with a pre-trained target model $f_t \in \mathcal{V}$ and a specific class label y . The \mathcal{V} -entropy should take its value at f_t and y , instead of the infimum expectation value, for the explanation. Therefore, we relax the \mathcal{V} -entropy terms $H_{\mathcal{V}}(Y_t | B)$ and $H_{\mathcal{V}}(Y_t | \mathcal{N}(z), B)$ into $-\log f_t(B; \mathbf{x}_k, y)$ and $-\log f_t(\mathcal{N}(z) \cup B; \mathbf{x}_k, y)$, respectively [14], for aligning the explanation with the target model $f_t \in \mathcal{V}$ and class label y . In this way, the attribution of patch z aligned with class y is defined as follows:

$$\phi_{k,y,z} = \mathbb{E}_{B \subseteq \mathcal{Z}(\mathbf{x}_k) \setminus \mathcal{N}(z)} [-\log f_t(B; \mathbf{x}_k, y) + \log f_t(\mathcal{N}(z) \cup B; \mathbf{x}_k, y)]. \quad (2)$$

It is impossible to enumerate the state of B over $B \subseteq \mathcal{Z}(\mathbf{x}_k) \setminus \mathcal{N}(z)$ in Equation (2). We follow existing work [27] to approximate it into two antithetical states to simplify the computation [27]. These cases involve considering the state of B to be entirely remaining patches $\mathcal{Z}(\mathbf{x}_k) \setminus \mathcal{N}(z)$ or empty set \emptyset , narrowing down the enumeration of $B \subseteq \mathcal{Z}(\mathbf{x}_k) \setminus \mathcal{N}(z)$ to $B \sim \{\mathcal{Z}(\mathbf{x}_k) \setminus \mathcal{N}(z), \emptyset\}$ in Equation (2). Based on our numerical studies in Appendix B, the approximate attribution shows positive correlation with the exact value, which indicates the approximation does not affect the quality of attribution. To summarize, we approximate the attribution value of patch z aligned with class y as follows:

$$\phi_{k,y,z} \approx \mathbb{E}_{B \sim \{\mathcal{Z}(\mathbf{x}_k) \setminus \mathcal{N}(z), \emptyset\}} [-\log f_t(B; \mathbf{x}_k, y) + \log f_t(\mathcal{N}(z) \cup B; \mathbf{x}_k, y)], \quad (3)$$

$$\sim \log f_t(\mathcal{N}(z); \mathbf{x}_k, y) - \log f_t(\mathcal{Z}(\mathbf{x}_k) \setminus \mathcal{N}(z); \mathbf{x}_k, y), \quad (4)$$

where the terms $f_t(\mathcal{Z}(\mathbf{x}_k); \mathbf{x}_k, y)$ and $f_t(\emptyset; \mathbf{x}_k, y)$ in Equation (3) are constant given \mathbf{x}_k and y , thus being omitted in Equation (4). Intuitively, the explanation of patch z depends on the gap of logit values, where $\mathcal{N}(z)$ and background patches $\mathcal{Z}(\mathbf{x}_k) \setminus \mathcal{N}(z)$ are taken as the input.

4.2 Definition of Meta-attribution

We introduce the concept of meta-attribution, formally defined in Definition 1. Note that Equation (4) relies on the downstream target model f_t , which is task-related. The purpose of meta-attribution is to disentangle the task-specific aspect of the attribution from Equation (4). This disentanglement renders the meta-attribution to be task-independent, as a foundation for explaining various tasks.

Definition 1 (Meta-attribution). *Given a backbone encoder G , the meta-attribution for a patch $z \in \mathcal{Z}(\mathbf{x}_k)$, $\mathbf{x}_k \in \mathcal{X}$, is represented by two tensors $\mathbf{g}_{k,z}$ and $\mathbf{h}_{k,z}$ as follows:*

$$\begin{aligned} \mathbf{g}_{k,z} &= G(\mathcal{N}(z); \mathbf{x}_k), \\ \mathbf{h}_{k,z} &= G(\mathcal{Z}(\mathbf{x}_k) \setminus \mathcal{N}(z); \mathbf{x}_k). \end{aligned} \quad (5)$$

Following Definition 1, the meta-attribution is defined as the input tensors of the logarithmic functions in Equation (4), where the task-specific model f_t is replaced into the backbone encoder G to disentangle the meta-attribution with specific tasks. This disentanglement enables the meta-attribution to transfer across various downstream tasks.

4.3 Transfer to Task-aligned Explanation

To explain the downstream tasks, we propose a transfer rule in Definition 2 to adapt the meta-attribution to explaining downstream tasks. This rule-based transfer method can effectively address the CH2 described in Section 1, without the need for additional training on task-specific data.

Definition 2 (Attribution Transfer). *If the task-specific function is given by $f_t = H_t \circ G$, then the explanation of $f_t(\mathbf{x}_k)$ on class y is generated by*

$$\phi_{k,y,z} = \log H_t(\mathbf{g}_{k,z}; y) - \log H_t(\mathbf{h}_{k,z}; y), \quad (6)$$

where $\mathbf{g}_{k,z}$ and $\mathbf{h}_{k,z}$ are the meta-attribution given by Equation (5); and G and H_t represent the backbone encoder and fine-tuned classifier on task t , respectively.

Following Definition 2, we can straightforwardly achieve the solution of $\phi_{k,y,z}$ to be consistent with Equation (4)¹. This alignment to $\phi_{k,y,z}$ can effectively explain downstream task t following the definition of conditional \mathcal{V} -information $I_{\mathcal{V}}(\mathcal{N}(z) \rightarrow Y_t | B)$, as described in Section 4.1.

5 Learning Meta-attribution

In this section, we introduce the details of **T**ransferable **V**ision **E**xplainer (TVE). Specifically, TVE pre-trains a DNN-based transferable explainer $E(\bullet | \theta)$ on large-scale image dataset to comprehensively learn the knowledge of meta-attribution. After the pre-training, TVE can transfer to various downstream tasks for end-to-end generating task-aligned explanation. To assess its performance, we theoretically analyze the explanation error in Theorem 1.

5.1 Explainer Pre-training

TVE employs a DNN-based explainer $E(\bullet | \theta)$ to generate the meta-attribution tensors. Specifically, the explainer $E(\bullet | \theta)$ produces two tensors for the meta-attribution, denoted as $[\hat{\mathbf{g}}_k, \hat{\mathbf{h}}_k] = E(\mathbf{x}_k | \theta)$, where $\hat{\mathbf{g}}_k = [\hat{\mathbf{g}}_{k,z} \in \mathbb{R}^D | z \in \mathcal{Z}(\mathbf{x}_k)]$ and $\hat{\mathbf{h}}_k = [\hat{\mathbf{h}}_{k,z} \in \mathbb{R}^D | z \in \mathcal{Z}(\mathbf{x}_k)]$ represent collections of meta-attribution for an instance \mathbf{x}_k . Each pair of elements $(\hat{\mathbf{g}}_{k,z}, \hat{\mathbf{h}}_{k,z})$ contribute to predicting the meta-attribution $(\mathbf{g}_{k,z}, \mathbf{h}_{k,z})$ defined in Definition 1. Pursuant to this objective, TVE updates the parameters of explainer $E(\bullet | \theta)$ to minimize the following loss function:

$$\mathcal{L}_{\theta}(\mathbf{x}_k) = \mathbb{E}_{z \sim \mathcal{Z}(\mathbf{x}_k)} [\|\hat{\mathbf{g}}_{k,z} - \mathbf{g}_{k,z}\|_2^2 + \|\hat{\mathbf{h}}_{k,z} - \mathbf{h}_{k,z}\|_2^2], \quad (7)$$

where $\mathbf{g}_{k,z}$ and $\mathbf{h}_{k,z}$ are defined in Definition 1.

Algorithm 1 summarizes one epoch of pre-training the transferable explainer $E(\bullet | \theta)$. Specifically, TVE first samples a mini-batch of image patches (lines 2); then follows Definition 1 to generate the meta-attribution (lines 3); finally updates the parameters of $E(\bullet | \theta)$ to minimize the loss function given by Equation (7) (line 4). The iteration ends with the convergence of $E(\bullet | \theta)$. Notably, the pre-training of $E(\bullet | \theta)$ is guided by the meta-attribution instead of specific tasks. This empowers the trained $E(\bullet | \theta)$ to remain impartial towards specific tasks, providing the flexibility for seamless adaptation across various downstream tasks.

Algorithm 1 One epoch of TVE pre-training

Input: Pre-training dataset \mathcal{D} .

Output: Transferable explainer $E(\bullet | \theta^*)$.

- 1: **for** $\mathbf{x}_k \sim \mathcal{D}$ **do**
 - 2: Sample patches $z \sim \mathcal{Z}(\mathbf{x}_k)$.
 - 3: Generate $\mathbf{g}_{k,z}$ and $\mathbf{h}_{k,z}$ following Definition 1.
 - 4: Update $E(\bullet | \theta)$ to minimize Equation (7).
 - 5: **end for**
-

¹We follow Definition 2 to have $\phi_{k,y,z} = \log H_t(\mathbf{g}_{k,z}; y) - \log H_t(\mathbf{h}_{k,z}; y) = \log f_t(\mathcal{N}(z); \mathbf{x}_k, y) - \log f_t(\mathcal{Z}(\mathbf{x}_k) \setminus \mathcal{N}(z); \mathbf{x}_k, y)$ that is consistent with Equation (4).

5.2 Generating Task-aligned Explanation

TVE follows Definition 2 to generate the task-aligned explanation. Specifically, to explain the inference process $(H_t \circ G)(\mathbf{x}_k)$ in task t , TVE first adopts the pre-trained transferable explainer to generate the meta-attribution $[\hat{\mathbf{g}}_k, \hat{\mathbf{h}}_k] = E(\mathbf{x}_k | \theta)$; then takes the value of $\hat{\mathbf{g}}_{k,z}$ and $\hat{\mathbf{h}}_{k,z}$ into Equation (6) to estimate the importance of each patch $z \in \mathcal{Z}(\mathbf{x}_k)$ to the inference result on class y . To summarize, TVE generates the attribution of a patch $z \in \mathcal{Z}(\mathbf{x}_k)$ by

$$\hat{\phi}_{k,y,z} = \log H_t(\hat{\mathbf{g}}_{k,z}; y) - \log H_t(\hat{\mathbf{h}}_{k,z}; y). \quad (8)$$

Let $\hat{\phi}_{k,y} = [\hat{\phi}_{k,y,z} | z \in \mathcal{Z}(\mathbf{x}_k)]$ denote the $P \times P$ explanation heatmap for the image \mathbf{x}_k , indicating the importance of all patches in \mathbf{x}_k to class y . TVE can efficiently generate the entire heatmap $\hat{\phi}_{k,y}$ for the image \mathbf{x}_k through a single feed forward pass: $\hat{\phi}_{k,y} = \log H_t(\hat{\mathbf{g}}_k; y) - \log H_t(\hat{\mathbf{h}}_k; y)$, where $\hat{\mathbf{g}}_k$ and $\hat{\mathbf{h}}_k$ are generated by $[\hat{\mathbf{g}}_k, \hat{\mathbf{h}}_k] = E(\mathbf{x}_k | \theta)$.

In particular, $H_t(\bullet; y)$ in Equation (8) encodes the knowledge of downstream task t . This knowledge significantly enables the explanation to align with the task t *without the need for additional training on the task-specific data*.

5.3 Theoretical Analysis

The theoretical analysis focuses on understanding the behavior of estimation error $|\hat{\phi}_{k,y,z} - \phi_{k,y,z}|$ during the TVE pre-training, where $\phi_{k,y,z}$ takes the \mathcal{V} -Information-aligned explanation defined in Section 4.1. Specifically, we examine the following two distinct cases to understand how the reduction in the pre-training loss function $\mathcal{L}_\theta(\mathbf{x}_k)$ diminishes the estimation error $|\hat{\phi}_{k,y,z} - \phi_{k,y,z}|$.

Ideal Case. We ideally consider $\mathcal{L}_\theta(\mathbf{x}_k) \rightarrow 0$ in this case. According to Equation (7), we have that $\hat{\mathbf{g}}_{k,z} \rightarrow \mathbf{g}_{k,z}$ and $\hat{\mathbf{h}}_{k,z} \rightarrow \mathbf{h}_{k,z}$. Then, the relations $\frac{H_t(\hat{\mathbf{g}}_{k,z}; y)}{H_t(\mathbf{g}_{k,z}; y)} \rightarrow 1$ and $\frac{H_t(\hat{\mathbf{h}}_{k,z}; y)}{H_t(\mathbf{h}_{k,z}; y)} \rightarrow 1$ are established. In this context, we have $|\hat{\phi}_{k,y,z} - \phi_{k,y,z}| \rightarrow 0$ according to Equations (6) and (8). This indicates $\hat{\phi}_{k,y,z}$ exactly converges to $\phi_{k,y,z}$ in the ideal scenario.

Practical Case. Without loss of generality, we consider $\mathcal{L}_\theta(\mathbf{x}_k)$ is not reduced to zero in this case. Specifically, Equation (7) indicates the reduction of $\mathcal{L}_\theta(\mathbf{x}_k)$ leads to $\hat{\mathbf{g}}_{k,z}$ and $\hat{\mathbf{h}}_{k,z}$ gradually approach $\mathbf{g}_{k,z}$ and $\mathbf{h}_{k,z}$, respectively. As a result, the values of $\frac{H_t(\hat{\mathbf{g}}_{k,z}; y)}{H_t(\mathbf{g}_{k,z}; y)}$ and $\frac{H_t(\hat{\mathbf{h}}_{k,z}; y)}{H_t(\mathbf{h}_{k,z}; y)}$ gradually converge to a narrower range around 1. We formulate this trend by assuming their values to be bounded within a range of $1 - \epsilon \leq \frac{H_t(\hat{\mathbf{g}}_{k,z}; y)}{H_t(\mathbf{g}_{k,z}; y)}, \frac{H_t(\hat{\mathbf{h}}_{k,z}; y)}{H_t(\mathbf{h}_{k,z}; y)} \leq 1 + \epsilon$, where $0 \leq \epsilon \ll 1$. Under these assumptions, we establish the upper bound of $|\hat{\phi}_{k,y,z} - \phi_{k,y,z}|$ in Theorem 1, with a detailed proof in Appendix C. This allows us to understand the behavior of estimation error in practical cases where $\mathcal{L}_\theta(\mathbf{x}_k)$ is not reduced to zero.

Theorem 1 (Explanation Error Bound). *Given the classifier $H_t(\bullet; \bullet)$ of the downstream task, if the output of classifier $H_t(\hat{\mathbf{g}}_{k,z}; y)$ and $H_t(\hat{\mathbf{h}}_{k,z}; y)$ fall within the range of $1 - \epsilon \leq \frac{H_t(\hat{\mathbf{g}}_{k,z}; y)}{H_t(\mathbf{g}_{k,z}; y)}, \frac{H_t(\hat{\mathbf{h}}_{k,z}; y)}{H_t(\mathbf{h}_{k,z}; y)} \leq 1 + \epsilon$, then, the upper bound of explanation error is given by*

$$\mathbb{E}_{\mathbf{x}_k \sim \mathcal{D}_t, y \sim \mathcal{Y}_t, z \sim \mathcal{Z}(\mathbf{x}_k)} |\hat{\phi}_{k,y,z} - \phi_{k,y,z}| \leq \frac{2\epsilon}{1 - \epsilon}, \quad (9)$$

where $\hat{\phi}_{k,y,z}$ and $\phi_{k,y,z}$ are given by Equation (8) and (6), respectively; and \mathcal{D}_t denotes the downstream dataset.

Intuition of Theorem 1. The value of ϵ reduces as the pre-training loss function $\mathcal{L}_\theta(\mathbf{x}_k)$ decreases. This reduction in ϵ explicitly lowers the estimation error bound $\frac{2\epsilon}{1 - \epsilon}$ aligned with the \mathcal{V} -Information-aligned explanation $\phi_{k,y,z}$ on downstream tasks. This underscores the TVE pre-training can significantly enhance the explanations for downstream tasks.

6 Experiment Results

In this section, we conduct experiments to evaluate TVE by answering the following research questions: **RQ1:** How does TVE perform compared with state-of-the-art baseline methods in terms of the fidelity? **RQ2:** How

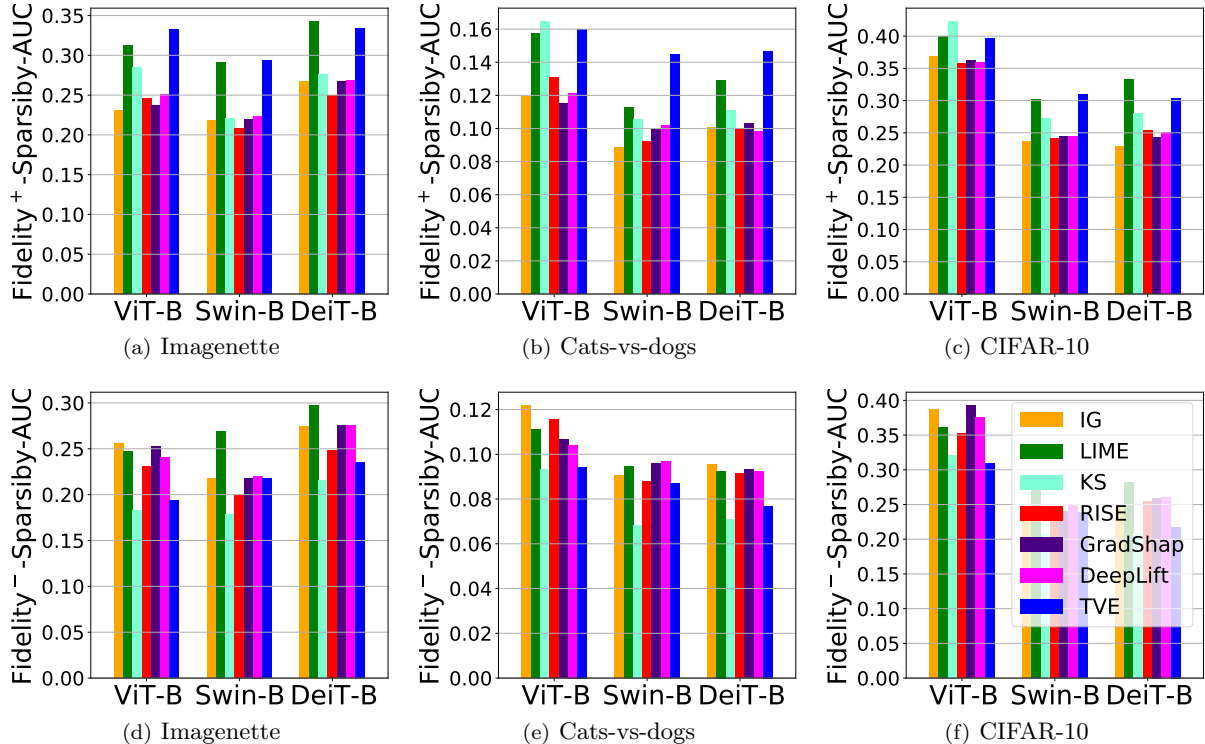


Figure 2: Fidelity⁺-Sparsity-AUC(↑) on the Imagenette (a), Cat-vs-dogs (b), and CIFAR-10 (c) datasets. Fidelity⁻-Sparsity-AUC(↓) on the Imagenette (d), Cat-vs-dogs (e), and CIFAR-10 (f) datasets.

does TVE perform in explaining fully fine-tuned target model on down-stream datasets? **RQ3:** How is the transferability of TVE across different downstream datasets? **RQ4:** Do both pre-training and attribution transfer in TVE contribute to explaining downstream tasks?

6.1 Experiment Setup

We clarify the datasets, target models, hyper-parameter settings in this section. More details about the baseline methods, evaluation metrics and implementation details are given in Appendixes F, G, and H, respectively.

Datasets. We consider the large-scale ImageNet dataset for TVE pre-training; and the Cats-vs-dogs [13], CIFAR-10 [23], and Imagenette [19] datasets for the downstream explaining tasks. Further details about the datasets are given in Appendix D.

Target Models. We comprehensively consider three architectures of vision transformers for downstream classification tasks, including the ViT-Base [11], Swin-Base [25], Deit-Base [36] transformers. We consider two settings of fine-tuning target models: *classifier-tuning* and *full-fine-tuning*. More details about the target model are given in Appendix E.

Hyper-parameter Settings. The experiment follows the pipeline of TVE pre-training, explanation generation and evaluation on multiple downstream datasets. Specifically, TVE adopts the Mask-AutoEncoder [16] as the backbone, followed by multiple Feed-Forward (FFN) layers² to generate the meta-attribution. More details about the explainer architecture and hyper-parameters of pre-training TVE are given in Appendix H. When deploying TVE to explaining downstream tasks, the explanation aligns with the prediction class given by the target model.

²A Mask-AutoEncoder consists of a ViT encoder followed by a ViT decoder; and an FFN layer consists of Linear layers, Layer-norm, and activation function, which are widely used in the Transformer structure. More details about the architecture are given in Appendix H.

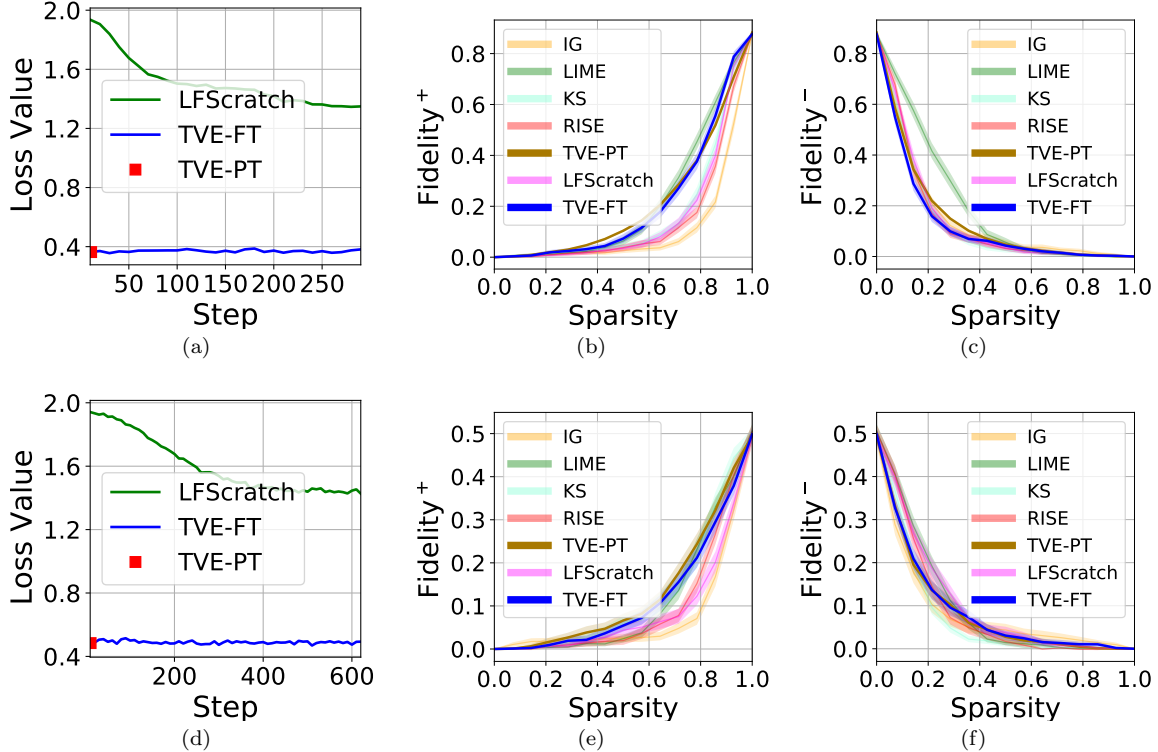


Figure 3: Fine-tuning loss versus epoch (a), Fidelity⁺ \uparrow versus Sparsity (b), and Fidelity⁻ \downarrow versus Sparsity (c) on the Imagenette dataset. Fine-tuning loss versus epoch (d), Fidelity⁺ \uparrow versus Sparsity (e), and Fidelity⁻ \downarrow versus Sparsity (f) on the cats-vs-dogs dataset.

6.2 Evaluation of Fidelity (RQ1)

In this section, we evaluate the fidelity of TVE under the classifier-tuning setting. Due to the space constraints, we present 18 figures illustrating the Fidelity⁺-sparsity curve(\uparrow) and the Fidelity⁻-sparsity curve(\downarrow) for explaining the ViT-Base, Swin-Base, and Deit-Base models on the Cats-vs-dogs, Imagenette, and CIFAR-10 datasets in Appendix I. To streamline our evaluation, we simplify the assessment of fidelity-sparsity curves by calculating its Area Under the Curve (AUC) over the sparsity from zero to one, which aligns with the average fidelity value. Intuitively, a higher Fidelity⁺-sparsity-AUC(\uparrow) indicates superior Fidelity⁺(\uparrow) across most sparsity levels, reflecting a more faithful explanation. Similarly, a lower Fidelity⁻-sparsity-AUC(\downarrow) signifies a more faithful explanation. More details about the fidelity-sparsity-AUC are given in Appendix G. On the Cats-vs-dogs, Imagenette, and CIFAR-10 datasets, we present the Fidelity⁺-sparsity-AUC(\uparrow) for explanations in Figures 2 (a)-(c), respectively, as well as the Fidelity⁻-sparsity-AUC(\downarrow) in Figures 2 (d)-(f), respectively. We have the following observations:

- TVE consistently exhibits promising performance in terms of both Fidelity⁺ (\uparrow) and Fidelity⁻ (\downarrow), outperforming the majority of baseline methods. This underscores TVE faithfully explains various downstream tasks within the scope of pre-training data distribution.
- TVE exhibits significant strengths in both Fidelity⁺ (\uparrow) and Fidelity⁻ (\downarrow), highlighting its effectiveness in identifying both important and non-important features. In contrast, the baseline methods fail to simultaneously achieve high Fidelity⁺ and low Fidelity⁻. For example, consider LIME’s performance when explaining the Deit-Base model on the CIFAR-10 dataset. While LIME excels in Fidelity⁺, it falls short in Fidelity⁻.

6.3 Explaining Fully Fine-tuned Models (RQ2)

In this section, we evaluate the fidelity of TVE under the full-fine-tuning setting to demonstrate its generalization ability. Notably, the ViT-Base classification model including both the backbone and classifier are fine-tuned

Table 1: Explanation Fidelity⁺-Sparsity-AUC(↑) and Fidelity⁻-Sparsity-AUC(↓) for DeiT-Base, Swin-Base, and DeiT-Base target models on the Cat-vs-dogs, Imagenette, and CIFAR-10 datasets.

	Datasets	Cats-vs-dogs		Imagenette		CIFAR-10	
Target model	Method	Fidelity ⁺ (↑)	Fidelity ⁻ (↓)	Fidelity ⁺ (↑)	Fidelity ⁻ (↓)	Fidelity ⁺ (↑)	Fidelity ⁻ (↓)
ViT-Base	ViTShapley	0.11±0.09	0.13±0.10	0.25±0.13	0.25±0.14	0.36±0.17	0.36±0.17
	TVE- H_g	0.14±0.11	0.10±0.08	0.29±0.14	0.18 ±0.10	0.39±0.18	0.34±0.17
	TVE	0.16 ±0.13	0.09 ±0.07	0.33 ±0.16	0.19±0.12	0.40 ±0.18	0.31 ±0.16
Swin-Base	ViTShapley	0.09±0.05	0.11±0.07	0.24±0.07	0.24±0.09	0.25±0.11	0.28±0.14
	TVE- H_g	0.14 ±0.09	0.10±0.07	0.29 ±0.08	0.24±0.07	0.26±0.12	0.27±0.13
	TVE	0.14 ±0.10	0.09 ±0.05	0.29 ±0.10	0.22 ±0.06	0.31 ±0.14	0.24 ±0.12
DeiT-Base	ViTShapley	0.12±0.08	0.1±0.07	0.22±0.09	0.29±0.11	0.28±0.13	0.24±0.13
	TVE- H_g	0.13±0.08	0.09±0.06	0.33 ±0.10	0.25±0.08	0.32 ±0.14	0.24±0.13
	TVE	0.15 ±0.10	0.08 ±0.06	0.33 ±0.10	0.24 ±0.08	0.30±0.13	0.22 ±0.12

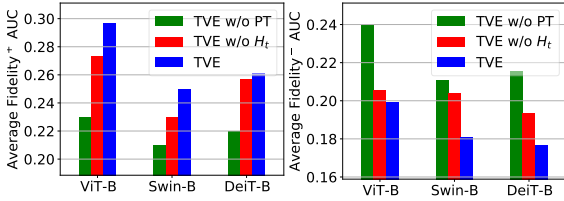


Figure 4: Fidelity of ablation studies.

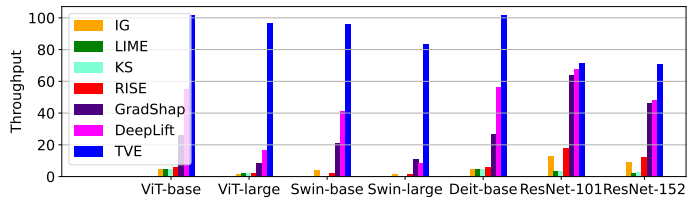


Figure 5: Throughput of explaining different architectures.

on downstream data, which are not available to TVE pre-training. The explanation considers three methods: learning from scratch (LFScratch), TVE pre-training (TVE-PT), and TVE fine-tuning (TVE-FT). To adapt to the fully fine-tuned target model, LFScratch trains the explainer on the downstream dataset for one epoch; TVE-PT simply transfers the pre-trained explainer to explaining the down-stream tasks; TVE-FT follows Algorithm 1 to fine-tune the explainer using the fine-tuned backbone encoder on the downstream dataset for one epoch. Here, we consider the Imagenette and Cat-vs-dogs datasets for the downstream tasks. Further details about fine-tuning the target models and explainers are given in Appendixes E and H, respectively. The loss value of LFScratch and TVE-FT versus the fine-tuning steps are shown in Figures 3 (a) and (d). The fidelity-sparsity curves of all methods are given in Figures 3 (b), (c), (e), and (f). We have the following observations:

- *TVE pre-training provides a good initial explainer for adaption to fully fine-tuned encoders.* According to Figures 3 (a,d), the TVE pre-trained explainer shows lower training loss than learning from scratch in the early epochs. This indicates the pre-training provides a good initial explainer for explaining downstream tasks.
- *TVE-PT can effectively explain the fully fine-tuned target model, even without fine-tuning the explainer on downstream datasets.* According to Figures 3 (b,c,e,f), TVE-PT shows competitive fidelity when comparing with TVE-FT and other baseline methods, and a significant improvement over LFScratch. This indicates the strong generalization ability of TVE, acquired through pre-training on the large-scale ImageNet dataset.
- *The pre-training of transferable explainer and fine-tuning of backbone encoder can be executed independently and parallelly.* Specifically, TVE pre-trains the transferable explainer based on open-sourced pre-trained backbone encoders and large-scale ImageNet dataset; meanwhile, the encoder can be fine-tuned in parallel on downstream datasets. This can significantly improve the efficiency and flexibility of deploying TVE to practical scenarios.

6.4 Evaluation of Transferability (RQ3)

We evaluate the transferability of TVE compared with ViT-Shapley [9], a state-of-the-art DNN-based explainer for vision models. Specifically, ViT-Shapley pre-trains the explainer on the large-scale ImageNet dataset, and deploys it to the Cat-vs-dogs, Imagenette, and CIFAR-10 datasets to generate the explanations. Different from ViT-Shapley, TVE transfers the explainer to downstream datasets via taking the task-specific classifier H_t into Equation (8). Moreover, we also consider a TVE- H_g method to study whether the pre-training of TVE contributes to explaining downstream tasks. Different from TVE, TVE- H_g takes a general classifier (pre-trained

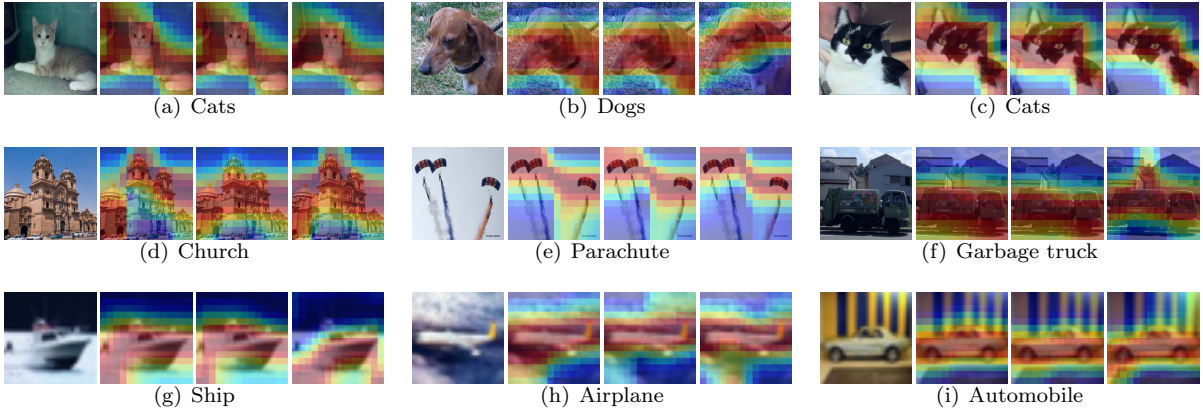


Figure 6: Visualization of explanation on the **Cats-vs-dogs** (a)-(c), **Imagenette** (d)-(f), and **CIFAR-10** (g)-(i) datasets. From the left to the right, each heatmap explains the inference of the **Swin-Base**, **DeiT-Base**, and **ViT-Base** models, respectively.

on the **ImageNet** dataset) into Equation (8) to generate the explanation. We follow Section 6.2 to adopt the fidelity-sparsity AUC to evaluate the average fidelity. Table 1 illustrates the fidelity for explaining the **ViT-Base**, **Swin-Base**, and **DeiT-Base** models on the **Cat-vs-dogs**, **Imagenette**, and **CIFAR-10** datasets. We have the following insights:

- *TVE has stronger transferability than ViT-Shapley.* Both TVE and **ViT-Shapley** are pre-trained on the large-scale **ImageNet** dataset, and transferred to the downstream datasets without additional training. Table 1 shows TVE has higher Fidelity⁺(\uparrow) and lower Fidelity⁻(\downarrow) than **ViT-Shapley**.
- *The pre-training of TVE significantly contributes to explaining downstream tasks.* TVE- H_g adopts the generally pre-trained explainer and classifier to explain downstream tasks, and achieves a reasonable fidelity on most of the datasets. This indicates the pre-training of TVE captures the transferable features across various datasets for explaining downstream tasks.
- *It is more faithful to explain downstream tasks based on the task-specific classifiers.* TVE outperforms TVE- H_g on most architectures and datasets, which indicates the attribution transfer had better take the classifier aligned with the downstream task for H_t in Definition 2.

6.5 Ablation Studies (RQ4)

We ablatedly study the contribution of the key steps in TVE to explaining downstream tasks, including the pre-training of transferable explainer and attribution transfer aligned to each task. For our evaluation, we consider three methods: TVE w/o Pre-training (PT), TVE w/o H_t , and TVE. Specifically, for TVE w/o PT, the explainer is randomly initialized without pre-training, and attribution transfer follows Definition 2. For TVE w/o H_t , the transferable explainer is pre-trained following Algorithm 1, and the explanation for each task is generated by $\hat{\phi}_k = \log H_g(\hat{\mathbf{g}}_k; y) - \log H_g(\hat{\mathbf{h}}_k; y)$, where H_g takes a general classifier pre-trained on the **ImageNet** datasets, instead of being fine-tuned corresponding to the task. Figure 4 illustrates the results of Fidelity⁺-Sparsity-AUC(\uparrow) and Fidelity⁻-Sparsity-AUC(\downarrow) for each method, where the fidelity score represents the averaged value on the **Cats-vs-dogs**, **Imagenette**, and **CIFAR-10** datasets. Other configurations remain consistent with Appendix H. Overall, we have the following observations:

- *TVE pre-training significantly contributes to explaining the downstream tasks.* This can be verified by the fidelity degradation observed from TVE w/o PT in Figure 4.
- *The classifier H_t for attribution transfer should align with the explaining task t .* It is observed in Figure 4 that TVE outperforms TVE w/o H_t . This indicates the task-aligned H_t is better than general classifiers for the attribution transfer to a specific task t .

6.6 Evaluation of Latency

In this section, we evaluate the latency of TVE compared with baseline methods. Specifically, we adopt the metric $\text{Throughput} = \frac{N_{\text{test}}}{T} (\uparrow)$ to evaluate the explanation latency, where N_{test} takes the number of testing instances and T signifies the total time consumed during the explanation process. Details about our computational infrastructure are given in Appendix J. Figure 5 shows the throughput of different methods explaining the ViT-Base/Large, Swin-Base/Large, Deit-Base, and ResNet-101/152 models on the ImageNet dataset. Overall, we observe:

- *TVE is more efficient than state-of-the-art baseline methods*, by generating explanations through a single feed-forward pass of the explainer. In contrast, the baseline methods rely on intensive samplings of the forward or backward passes of the target model, resulting in a considerably slower explanation process. For example, although KernelSHAP exhibits comparable Fidelity $^-$ (\downarrow) with TVE, as shown in Figure 2, its significantly lower throughput limits its practicality in real-world scenarios.
- *TVE exhibits the most negligible decrease in throughput as the size of the target model grows*, as seen when transitioning from ViT-Base to ViT-Large. This advantage stems from the fact that TVE’s latency is contingent upon the explainer’s model size, rather than the target model. In contrast, the baseline methods suffer from notable performance slowdown as the size of the target model increases, due to the necessity of sampling the target model to generate explanations.

6.7 Case Studies

In this section, we visualize the explanations generated by TVE, demonstrating its power in helping human users understand vision models. Specifically, we randomly sample three instances from the Cats-vs-dogs, Imagenette, and CIFAR-10 datasets, and visualize the explanations of Swin-Base, Deit-Base, and ViT-Base models in Figure 6, where sub-figures (a)-(c), (d)-(f), and (g)-(i) correspond to the Cats-vs-dogs, Imagenette, and CIFAR-10 datasets, respectively. In each sub-figure, from the left-side to the right-side, the three heatmaps explain the inference of the Swin-Base, Deit-Base, and ViT-Base model, respectively. Notably, TVE generates the explanation heatmap in an end-to-end manner *without pre- or post-processing*. More case studies on the ImageNet dataset are shown in Appendix K. According to the case study, we observe:

- *The salient patches emphasized by TVE’s explanation reveal semantically meaningful patterns*. For example, as depicted in Figures 6 (d), (e), and (g), the Swin-Base model concentrates on the tower, canopy and bow, respectively, to identify a church, parachute, and ship.
- *TVE does not rely on pre-processing of the image or post-processing of the explanation heatmap*. In contrast, existing work EAC [33] requires SAM [21] to segment the input image before explaining, which is less flexible than TVE.
- *Different model architectures make predictions based on distinct image elements*. For instance, as illustrated in Figure 6 (g), the Swin-Base and Deit-Base models primarily emphasize the ship’s bow for identification. In contrast, the ViT-Base model takes into account the ship’s keel for its prediction.

7 Conclusion

In this work, we propose a framework of attribution transfer, incorporating a meta-attribution to extract the foundation knowledge and a transfer rule to utilize this knowledge for explaining various downstream tasks. Building upon this framework, we introduce TVE, a transferable explainer pre-trained on large-scale image datasets. Notably, TVE shows strong transferability to effectively explain various downstream tasks without the need for training on task-specific data. Experiment results validate the promising performance of TVE in explaining three architectures of vision Transformer across three downstream datasets. Significantly, the strong transferability of TVE facilitates efficient and flexible deployment to various downstream scenarios.

References

- [1] Marco Ancona, Enea Ceolini, Cengiz Öztireli, and Markus Gross. Towards better understanding of gradient-based attribution methods for deep neural networks. *arXiv preprint arXiv:1711.06104*, 2017.
- [2] Chia-Yuan Chang, Jiayi Yuan, Sirui Ding, Qiaoyu Tan, Kai Zhang, Xiaoqian Jiang, Xia Hu, and Na Zou. Towards fair patient-trial matching via patient-criterion level fairness constraint. *arXiv preprint arXiv:2303.13790*, 2023.
- [3] Hanjie Chen, Faeze Brahman, Xiang Ren, Yangfeng Ji, Yejin Choi, and Swabha Swayamdipta. Rev: information-theoretic evaluation of free-text rationales. *arXiv preprint arXiv:2210.04982*, 2022.
- [4] Lu Chen, Siyu Lou, Keyan Zhang, Jin Huang, and Quanshi Zhang. Harsanyinet: Computing accurate shapley values in a single forward propagation. *arXiv preprint arXiv:2304.01811*, 2023.
- [5] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020.
- [6] Sasank Chilamkurthy. Transfer learning for computer vision tutorial. *PyTorch Tutorials*, 2017.
- [7] Yu-Neng Chuang, Guanchu Wang, Fan Yang, Zirui Liu, Xuanning Cai, Mengnan Du, and Xia Hu. Efficient xai techniques: A taxonomic survey. *arXiv preprint arXiv:2302.03225*, 2023.
- [8] Yu-Neng Chuang, Guanchu Wang, Fan Yang, Quan Zhou, Pushkar Tripathi, Xuanning Cai, and Xia Hu. Cortx: Contrastive framework for real-time explanation. *arXiv preprint arXiv:2303.02794*, 2023.
- [9] Ian Covert, Chanwoo Kim, and Su-In Lee. Learning to estimate shapley values with vision transformers. *arXiv preprint arXiv:2206.05282*, 2022.
- [10] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [11] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [12] Mengnan Du, Ninghao Liu, and Xia Hu. Techniques for interpretable machine learning. *Communications of the ACM*, 63(1):68–77, 2019.
- [13] Jeremy Elson, John (JD) Douceur, Jon Howell, and Jared Saul. Asirra: A captcha that exploits interest-aligned manual image categorization. In *Proceedings of 14th ACM Conference on Computer and Communications Security (CCS)*. Association for Computing Machinery, Inc., October 2007.
- [14] Kawin Ethayarajh, Yejin Choi, and Swabha Swayamdipta. Understanding dataset difficulty with v-usable information. pages 5988–6008, 2022.
- [15] Bryce Goodman and Seth Flaxman. European union regulations on algorithmic decision-making and a “right to explanation”. *AI magazine*, 38(3):50–57, 2017.
- [16] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16000–16009, 2022.
- [17] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9729–9738, 2020.
- [18] John Hewitt, Kawin Ethayarajh, Percy Liang, and Christopher D Manning. Conditional probing: measuring usable information beyond a baseline. *arXiv preprint arXiv:2109.09234*, 2021.
- [19] Jeremy Howard. Imagenette: A smaller subset of 10 easily classified classes from imagenet, March 2019.

- [20] Neil Jethani, Mukund Sudarshan, Ian Connick Covert, Su-In Lee, and Rajesh Ranganath. Fastshap: Real-time shapley value estimation. In *International Conference on Learning Representations*, 2021.
- [21] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. *arXiv preprint arXiv:2304.02643*, 2023.
- [22] Narine Kokhlikyan, Vivek Miglani, Miguel Martin, Edward Wang, Bilal Alsallakh, Jonathan Reynolds, Alexander Melnikov, Natalia Kliushkina, Carlos Araya, Siqi Yan, et al. Captum: A unified and generic model interpretability library for pytorch. *arXiv preprint arXiv:2009.07896*, 2020.
- [23] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [24] Yang Liu, Sujay Khandagale, Colin White, and Willie Neiswanger. Synthetic benchmarks for scientific research in explainable machine learning. 2021.
- [25] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10012–10022, 2021.
- [26] Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. *Advances in neural information processing systems*, 30, 2017.
- [27] Rory Mitchell, Joshua Cooper, Eibe Frank, and Geoffrey Holmes. Sampling permutations for shapley value estimation. *The Journal of Machine Learning Research*, 23(1):2082–2127, 2022.
- [28] Vitali Petsiuk, Abir Das, and Kate Saenko. Rise: Randomized input sampling for explanation of black-box models. *arXiv preprint arXiv:1806.07421*, 2018.
- [29] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. ” why should i trust you?” explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1135–1144, 2016.
- [30] Yao Rong, Guanchu Wang, Qizhang Feng, Ninghao Liu, Zirui Liu, Enkelejda Kasneci, and Xia Hu. Efficient gnn explanation via learning removal-based attribution. *arXiv preprint arXiv:2306.05760*, 2023.
- [31] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, pages 618–626, 2017.
- [32] Emily Steel and Julia Angwin. On the web’s cutting edge, anonymity in name only. *The Wall Street Journal*, 4, 2010.
- [33] Ao Sun, Pingchuan Ma, Yuanyuan Yuan, and Shuai Wang. Explain any concept: Segment anything meets concept-based explanation. *arXiv preprint arXiv:2305.10289*, 2023.
- [34] Yanpeng Sun, Qiang Chen, Xiangyu He, Jian Wang, Haocheng Feng, Junyu Han, Errui Ding, Jian Cheng, Zechao Li, and Jingdong Wang. Singular value fine-tuning: Few-shot segmentation requires few-parameters fine-tuning. *Advances in Neural Information Processing Systems*, 35:37484–37496, 2022.
- [35] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *International conference on machine learning*, pages 3319–3328. PMLR, 2017.
- [36] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *International conference on machine learning*, pages 10347–10357. PMLR, 2021.
- [37] Guanchu Wang, Yu-Neng Chuang, Mengnan Du, Fan Yang, Quan Zhou, Pushkar Tripathi, Xuanting Cai, and Xia Hu. Accelerating shapley explanation via contributive cooperator selection. *arXiv preprint arXiv:2206.08529*, 2022.

- [38] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations*, pages 38–45, 2020.
- [39] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online, October 2020. Association for Computational Linguistics.
- [40] Yilun Xu, Shengjia Zhao, Jiaming Song, Russell Stewart, and Stefano Ermon. A theory of usable information under computational constraints. *arXiv preprint arXiv:2002.10689*, 2020.
- [41] Fan Yang, Ninghao Liu, Suhang Wang, and Xia Hu. Towards interpretation of recommender systems with sorted explanation paths. In *2018 IEEE International Conference on Data Mining (ICDM)*, pages 667–676. IEEE, 2018.
- [42] Zhou Yang, Ninghao Liu, Xia Ben Hu, and Fang Jin. Tutorial on deep learning interpretation: A data perspective. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*, pages 5156–5159, 2022.

Appendix

A Related Work

Explainable machine learning (ML) has made significant advancements, leading to model transparency and better human understanding of deep neural networks (DNNs) [12]. Specifically, existing work of explainable ML can be categorized into two groups: local explainers and DNN-based explainers [7].

Local Explainer. Local explainer focuses on constructing local explanation based on perturbation of the target black-box model, like KernelSHAP [26], LIME [29], GradCAM [31], and Integrated Gradient [35]. Specifically, KernelSHAP approximates the Shapleyvalue by learning an explainable surrogate (linear) model based on the DNN output of reference input for each feature; LIME generates the explanation by sampling points around the input instance and using DNN output at these points to learn a surrogate (linear) model; Integrated Gradients estimates the explanation by the integral of the gradients of DNN output with respect to the inputs, along the pathway from specified references to the inputs. These pieces of work rely on resource-intensive procedures like sampling or backpropagation of the target black-box model [24], leading to undesirable trade-off between the efficiency and interpretation fidelity [7].

DNN-based Explainer. This branch of work leverages the training process to acquire proficiency in constructing a DNN-based explainer, utilizing explanation values as training labels [7]. This innovative strategy empowers the simultaneous generation of explanations for an entire batch of instances through a single, streamlined feedforward operation of the DNN-based explainer. Exemplifying this progress are innovative approaches like FastSHAP [20], ViT-Shapley [9], CoRTX [8], LARA [30, 37], and HarsanyiNet [4]. To be concrete, FastSHAP and ViT-Shapley adopt a DNN as the explainer to learn the Shapley value, which relies on task-specific training and cannot be transferred across different tasks; and CoRTX argues the training of DNN-based explainer through a contrastive pre-training framework, and adopt the true Shapley value to fine-tune the explainer. The DNN-based explainer have played a pivotal role in significantly streamlining the deployment of DNN explanations within real-time applications. However, they are constrained to explaining individual black box models, and they lack the ability to transfer the explanation across various models and tasks. This limitation results in the explanation of various tasks in practical scenarios becoming time- and resource-consuming due to the necessity of training different explainers for each task.

B Approximation of Attribution

We conduct experiments to study the relationship between the approximate attribution $\mathbb{E}_{B \sim \{\mathcal{Z}(\mathbf{x}_k) \setminus \mathcal{N}(z), \emptyset\}}[\dots]$ and its exact value $\mathbb{E}_{B \sim \text{Subset of } \mathcal{Z}(\mathbf{x}_k) \setminus \mathcal{N}(z)}[\dots]$ on the **ImageNet** dataset, where \dots is the abbreviation of $-\log f_t(B; \mathbf{x}_k, y) + \log f_t(\mathcal{N}(z) \cup B; \mathbf{x}_k, y)$. Specifically, we collect the samples of $\mathbb{E}_{B \sim \{\mathcal{Z}(\mathbf{x}_k) \setminus \mathcal{N}(z), \emptyset\}}[\dots]$ and $\mathbb{E}_{B \sim \text{Subset of } \mathcal{Z}(\mathbf{x}_k) \setminus \mathcal{N}(z)}[\dots]$, where \mathbf{x}_k take 100 instances randomly sampled from the **ImageNet** dataset; and the target models f_t take the **ViT-Base**(a, d), **Swin-Base**(b, e), and **DeiT-Base**(c, f) models trained on the **ImageNet** dataset. The samples of $\mathbb{E}_{B \sim \text{Subset of } \mathcal{Z}(\mathbf{x}_k) \setminus \mathcal{N}(z)}[\dots]$ versus $\mathbb{E}_{B \sim \{\mathcal{Z}(\mathbf{x}_k) \setminus \mathcal{N}(z), \emptyset\}}[\dots]$ is plotted in Figure 7. It is observed that the value of $\mathbb{E}_{B \sim \{\mathcal{Z}(\mathbf{x}_k) \setminus \mathcal{N}(z), \emptyset\}}[\dots]$ after the approximation shows positive linear correlation with $\mathbb{E}_{B \sim \text{Subset of } \mathcal{Z}(\mathbf{x}_k) \setminus \mathcal{N}(z)}[\dots]$. This indicates the approximate value $\mathbb{E}_{B \sim \{\mathcal{Z}(\mathbf{x}_k) \setminus \mathcal{N}(z), \emptyset\}}[\dots]$ can take the place of $\mathbb{E}_{B \sim \text{Subset of } \mathcal{Z}(\mathbf{x}_k) \setminus \mathcal{N}(z)}[\dots]$ for the function of attribution.

C Proof of Theorem 1

We prove Theorem 1 in this section.

Theorem 1 (Explanation Error Bound). *Given the classifier $H_t(\bullet)$ of the downstream task, if the output of classifier $H_t(\hat{\mathbf{g}}_{k,z}; y)$ and $H_t(\hat{\mathbf{h}}_{k,z}; y)$ fall within the range of $1 - \epsilon \leq \frac{H_t(\hat{\mathbf{g}}_{k,z}; y)}{H_t(\mathbf{g}_{k,z}; y)}, \frac{H_t(\hat{\mathbf{h}}_{k,z}; y)}{H_t(\mathbf{h}_{k,z}; y)} \leq 1 + \epsilon$, then, the upper bound of explanation error is given by*

$$\mathbb{E}_{\mathbf{x}_k \sim \mathcal{D}_t, y \sim \mathcal{Y}_t, z \sim \mathcal{Z}(\mathbf{x}_k)} |\hat{\phi}_{k,y,z} - \phi_{k,y,z}| \leq \frac{2\epsilon}{1 - \epsilon}, \quad (10)$$

where $\hat{\phi}_{k,y,z}$ and $\phi_{k,y,z}$ are given by Equation (8) and (6), respectively; and \mathcal{D}_t denotes the downstream dataset.

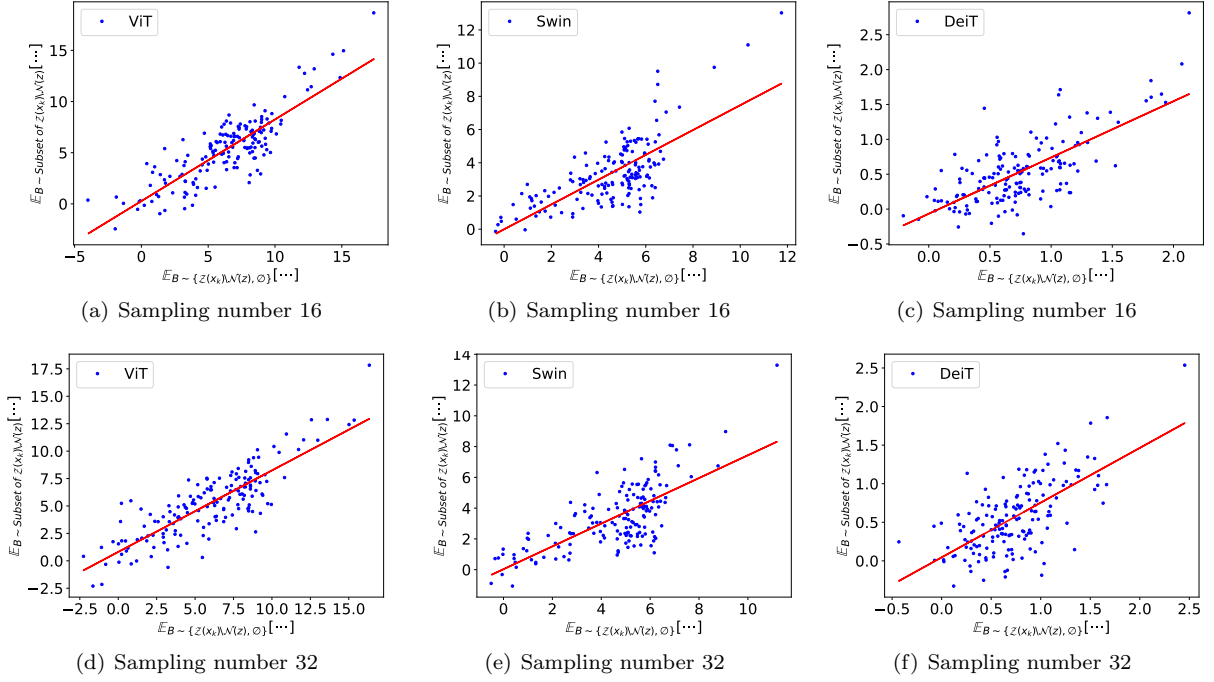


Figure 7: $\mathbb{E}_{B \sim \text{Subset of } \mathcal{Z}(\mathbf{x}_k) \setminus \mathcal{N}(z), \emptyset}[\dots]$ versus $\mathbb{E}_{B \sim \{\mathcal{Z}(\mathbf{x}_k) \setminus \mathcal{N}(z), \emptyset\}}[\dots]$, where \dots is the abbreviation of $-\log f_t(B; \mathbf{x}_k, y) + \log f_t(\mathcal{N}(z) \cup B; \mathbf{x}_k, y)$; and f_t takes the trained ViT-Base(a, d), Swin-Base(b, e), and DeiT-Base(c, f) models on the ImageNet dataset. The sampling number of $B \sim \text{Subset of } \mathcal{Z}(\mathbf{x}_k) \setminus \mathcal{N}(z)$ is 16 and 32 for Sub-figures (a)-(c) and (d)-(f), respectively.

Proof. To achieve the explanation error bound, we first have the upper bound of $\hat{\phi}_{k,y,z} - \phi_{k,y,z}$ given by

$$\hat{\phi}_{k,y,z} - \phi_{k,y,z} = \log H_t(\hat{\mathbf{g}}_{k,z}; y) - \log H_t(\mathbf{g}_{k,z}; y) + \log H_t(\mathbf{h}_{k,z}; y) - \log H_t(\hat{\mathbf{h}}_{k,z}; y), \quad (11)$$

$$= \log \frac{H_t(\hat{\mathbf{g}}_{k,z}; y)}{H_t(\mathbf{g}_{k,z}; y)} + \log \frac{H_t(\mathbf{h}_{k,z}; y)}{H_t(\hat{\mathbf{h}}_{k,z}; y)} \leq \frac{H_t(\hat{\mathbf{g}}_{k,z}; y)}{H_t(\mathbf{g}_{k,z}; y)} - 1 + \frac{H_t(\mathbf{h}_{k,z}; y)}{H_t(\hat{\mathbf{h}}_{k,z}; y)} - 1, \quad (12)$$

$$\leq \frac{H_t(\hat{\mathbf{g}}_{k,z}; y)}{H_t(\mathbf{g}_{k,z}; y)} - 1 + \frac{H_t(\mathbf{h}_{k,z}; y)}{H_t(\hat{\mathbf{h}}_{k,z}; y)} - 1 \leq \epsilon + \epsilon, \quad (13)$$

Then, we have the lower bound of $\hat{\phi}_{k,y,z} - \phi_{k,y,z}$ as follows,

$$\hat{\phi}_{k,y,z} - \phi_{k,y,z} = -\log \frac{H_t(\mathbf{g}_{k,z}; y)}{H_t(\hat{\mathbf{g}}_{k,z}; y)} - \log \frac{H_t(\hat{\mathbf{h}}_{k,z}; y)}{H_t(\mathbf{h}_{k,z}; y)} \geq 1 - \frac{H_t(\mathbf{g}_{k,z}; y)}{H_t(\hat{\mathbf{g}}_{k,z}; y)} + 1 - \frac{H_t(\hat{\mathbf{h}}_{k,z}; y)}{H_t(\mathbf{h}_{k,z}; y)}, \quad (14)$$

$$= 2 - \left(\frac{H_t(\mathbf{g}_{k,z}; y)}{H_t(\hat{\mathbf{g}}_{k,z}; y)} + \frac{H_t(\hat{\mathbf{h}}_{k,z}; y)}{H_t(\mathbf{h}_{k,z}; y)} \right) \geq 2 - \frac{1}{1-\epsilon} - \frac{1}{1-\epsilon} = \frac{-2\epsilon}{1-\epsilon} \quad (15)$$

Combining Equations (10) and (15), we achieve the upper bound of estimation error given by

$$|\hat{\phi}_{k,y,z} - \phi_{k,y,z}| \leq \max \left\{ 2\epsilon, \frac{2\epsilon}{1-\epsilon} \right\} = \frac{2\epsilon}{1-\epsilon}. \quad (16)$$

□

D Details about the Datasets

We consider the large-scale ImageNet dataset [10] for TVE pre-training; and the **Cats-vs-dogs** [13], **CIFAR-10** [23], and **Imagenette** [19] datasets for the downstream task of explanation. **ImageNet** [10]: A large scale image dataset which has over one million color images covering 1000 categories, where each image has 224×224

pixels. **Cats-vs-dogs** [13]: A dataset of cats and dogs images. It has 25000 training instances and 12500 testing instances. **CIFAR-10** [23]: An image dataset with 60,000 color images in 10 different classes, where each image has 32×32 pixels. **Imagenette** [19]: A benchmark dataset of explainable machine learning for vision models. It contains 10 classes of the images from the Imagenet.

E Details about Target Models for Downstream Classification.

E.1 Setup of Fine-tuning the Target Models

For downstream classification tasks, we comprehensively consider three architectures of vision transformers as the backbone encoders, including the **ViT-Base/Large** [11], **Swin-Base/Large** [25], **Deit-Base** [36] transformers. The classification models (to be explained) consist of one of the backbone encoders with **ImageNet** pre-trained weights and a linear classifier. For the task-specific fine-tuning of target models, we consider two mechanisms: *classifier-tuning* and *full-fine-tuning*. Specifically, the classifier-tuning follows the transfer learning setting [6, 17, 5] to freeze the parameters of backbone encoder during the fine-tuning; and the full-fine-tuning updates all parameters during the finetuning. Note that the classifier-tuning can not only be more efficient but also prevent the over-fitting problem on downstream data due to fewer trainable parameters [34]. We consider the classifier-tuning for most of our experiments including Sections 6.2, 6.4, 6.5, and 6.7; and consider the full-fine-tuning in Section 6.3; while these two mechanisms yield the same result for Section 6.6. The hyper-parameters of task-specific fine-tuning are given in Appendix E.2.

E.2 Hyper-parameter Setting of Fine-tuning the Target Models on Downstream Tasks

The downstream classification models consist of the backbones of **ViT-Base/Large**, **Swin-Base/Large**, **Deit-Base** transformers, and a linear classifier. The hyper-parameters of fine-tuning the classification models on the **Cats-vs-dogs**, **CIFAR-10**, and **Imagenette** datasets are given in Table 2. After the fine-tuning, the classification accuracy on each downstream dataset is given in Table 3.

Table 2: Hyper-parameters of fine-tuning the target model on downstream datasets.

Datasets	Cats-vs-dogs	CIFAR-10	Imagenette
Target backbone	ViT-Base, Swin-Base, and Deit-Base		
Classifier	Linear classifier		
Fine-tuning mechanism	classifier-tuning and full-fine-tuning		
Optimizer	ADAM		
Learning rate	2×10^{-4}		
Mini-batch size	256		
Scheduler	Linear		
Warm-up-ratio	0.05		
Weight-decay	0.05		
Epoch	5		

Table 3: Accuracy of the target model on downstream datasets.

Model Architecture	ViT-Base		Swin-Base		Deit-Base	
	θ_H	θ_H, θ_G	θ_H	θ_H, θ_G	θ_H	θ_H, θ_G
Cats-vs-dogs	99.6%	99.5%	99.6%	99.7%	99.4%	98.1%
Imagenette	99.3%	99.3%	99.8%	99.7%	99.8%	99.4%
CIFAR-10	92.2%	98.9%	97.0%	98.6%	94.2%	98.1%

F Details about the Baseline Methods

We consider seven baseline methods for comparison, which include general explanation methods: LIME [29], IG [35], RISE [28], and DeepLift [1]; Shapley explanation methods: KernelSHAP (KS) [26], and GradShap [26]; and DNN-based explainer: ViT-Shapley [9] in our experiment.

ViT-Shapley: This work adopts vision transformers as the explainer to learn the Shapley value. This work requires task-specific data to train the explainer. **RISE:** RISE randomly perturbs the input, and average all the masks weighted by the perturbed DNN output for the final saliency map. The sampling number takes the default value 50. **IG:** Integrated Gradients estimates the explanation by the integral of the gradients of DNN output with respect to the inputs, along the pathway from specified references to the inputs. **DeepLift:** DeepLift generates the explanation by decomposing DNN output on a specific input by backpropagating the contributions of all neurons in the network to every feature of the input. **KernelSHAP:** KernelSHAP approximates the Shapley value by learning an explainable surrogate (linear) model based on the DNN output of reference input for each feature. The sampling number takes the default value 25 for each instance according to the `captum.ai` [22]. **GradShap:** GradShap estimates the importance features by computing the expectations of gradients by randomly sampling from the distribution of references. **LIME [29]:** LIME generates the explanation by sampling points around the input instance and using DNN output at these points to learn a surrogate (linear) model. The sampling number takes the default value 25 according to the `captum.ai`. For implementation, we take the IG, DeepLift, and GradShap algorithms on the `captum.ai`, where the `multiply_by_inputs` factor takes false to achieve the local attribution for each instance.

G Evaluation Metrics

Fidelity-sparsity Curve: We consider the fidelity to evaluate the explanation following existing work [42, 8]. Specifically, the fidelity evaluates the explanation via *removing the important or trivial patches* from the input instance and *collecting the prediction difference of the target model f_t* . These two perspectives of evaluation are formalized into Fidelity⁺ and Fidelity⁻, respectively. Specifically, provided a subset of patches $\mathcal{S}^* \subseteq \mathcal{Z}(\mathbf{x}_k)$ that are important to the target model f_t by an explanation method, the Fidelity⁺ and Fidelity⁻ evaluates the explanation following

$$\begin{aligned}\uparrow \text{Fidelity}^+ &= \frac{1}{|\mathcal{D}_{\text{task}}|} \sum_{\mathbf{x} \in \mathcal{D}_{\text{task}}} f_t(\mathcal{Z}(\mathbf{x}_k); \mathbf{x}_k, y) - f_t(\mathcal{Z}(\mathbf{x}_k) \setminus \mathcal{S}^*; \mathbf{x}_k, y), \\ \downarrow \text{Fidelity}^- &= \frac{1}{|\mathcal{D}_{\text{task}}|} \sum_{\mathbf{x} \in \mathcal{D}_{\text{task}}} f_t(\mathcal{Z}(\mathbf{x}_k); \mathbf{x}_k, y) - f_t(\mathcal{S}^*; \mathbf{x}_k, y).\end{aligned}$$

Higher Fidelity⁺ indicates a better explanation for prediction y , since the truly important patches of image \mathbf{x}_k have been removed, leading to a significant difference of model prediction. Moreover, lower Fidelity⁻ implies a better explanation for prediction y , since the truly important patches have been preserved in \mathcal{S}^* to keep the prediction similar to the original one. The fidelity should be compared at the same level of sparsity $|\mathcal{S}^*|/|\mathcal{U}|$. Consequently, we consider the evaluation of fidelity versus the sparsity in most cases.

Fidelity-sparsity-AUC Metric To streamline our evaluation, we simplify the assessment of fidelity-sparsity curves by calculating its Area Under the Curve (AUC) over the sparsity from zero to one, which aligns with the average fidelity value. In the last paragraph, we have shown that higher Fidelity⁺ and lower Fidelity⁻ at the same level of sparsity indicate more faithful explanation. To streamline the evaluation, the assessment of fidelity-sparsity curves can be simplified into its Area Under the Curve (AUC) over the sparsity from zero to one, as shown in Figures 8 (a) and (b). The Fidelity-sparsity-AUC aligns with the average fidelity value. Specifically, a higher Fidelity⁺-sparsity-AUC (\uparrow) indicates better Fidelity⁺ performance across most sparsity levels, reflecting a more faithful explanation. Similarly, a lower Fidelity⁻-sparsity-AUC signifies a more faithful explanation. For the given example in Figures 8 (a) and (b), explanation A is more faithful than B.

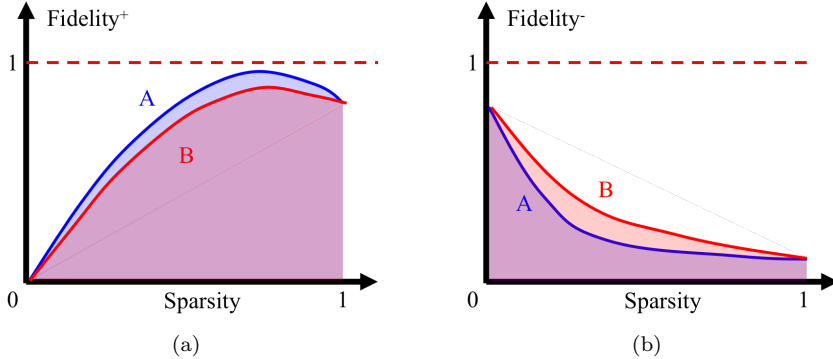


Figure 8: Illustration of Fidelity⁺-sparsity-AUC (a) and Fidelity⁻-sparsity-AUC (b)

H Implementation Details about TVE

Architecture of Generic Explainer. The architecture of the transferable explainer is shown in Figure 9 (a). Specifically, the explainer takes the `Mask-AutoEncoder-Base` [16] for the backbone. As shown in Figure 10, the `Mask-AutoEncoder-Base` architecture is a pipeline of 12-layer ViT encoder and 8-layer ViT decoder, where the input and output shape are $[BS, 3, 224, 224]$ and $[BS, P, P, 768]$, respectively. More details about the `Mask-AutoEncoder-Base` can be referred to its source code³.

Since the output shape of the `Mask-AutoEncoder-Base` is $[BS, P \times P, 768]$ is not matched with that of the meta-attribution $[BS, P \times P, D]$, where BS denotes the mini-batch size. We adopt $n \times$ FFN-layers as explainer heads to map the output tensor of the `Mask-AutoEncoder-Base` into meta-attribution, where we found $n = 17$ enables the explainer to have strong generalization ability to explain various downstream tasks. The structure of an explainer head is given in Figure 9 (b). The first explainer head does not have the skip connection due to the mismatch of tensor shapes. The last explainer head does not have the GELU activation.

Backbone Encoder. We comprehensively consider three backbone encoders for during the pre-training of transferable explainer, including the `ViT-Base/Large`, `Swin-Base/Large`, `DeiT-Base` transformers. Their pre-trained weights are loaded from the HuggingFace library [38]. The hyper-parameter setting of TVE pre-training is given in Table 4.

Table 4: Hyper-parameters of TVE pre-training on the ImageNet dataset.

Target Encoder	ViT-Base	Swin-Base	DeiT-Base
Explainer Architecture		Figure 9	
Pixel # per image $W \times W$		224×224	
Patch # per image $P \times P$		14×14	
Pixel # per patch $C \times C$		16×16	
Shape of \mathbf{g}_k and \mathbf{h}_k	$14 \times 14 \times 768$	$14 \times 14 \times 1024$	$14 \times 14 \times 768$
Optimizer		ADAM	
Learning rate		1×10^{-3}	
Mini-batch size		64 per GPU \times 4 GPUs	
Scheduler		CosineAnnealingLR	
Warm-up-ratio		0.05	
Weight-decay		0.05	
Training steps		2×10^5	
Neighbor patches		0-, 1-, 2-hop neighbor patches	

³https://github.com/huggingface/transformers/blob/main/src/transformers/models/vit_mae/modeling_vit_mae.py

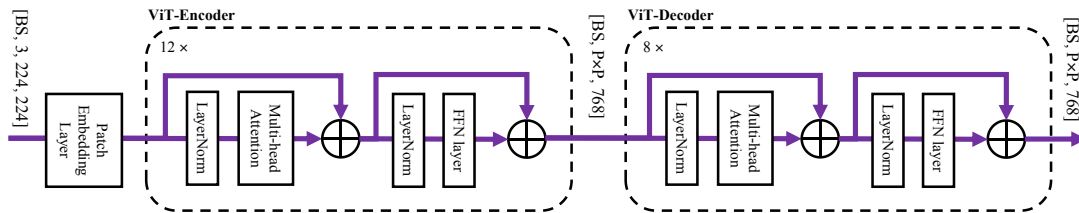
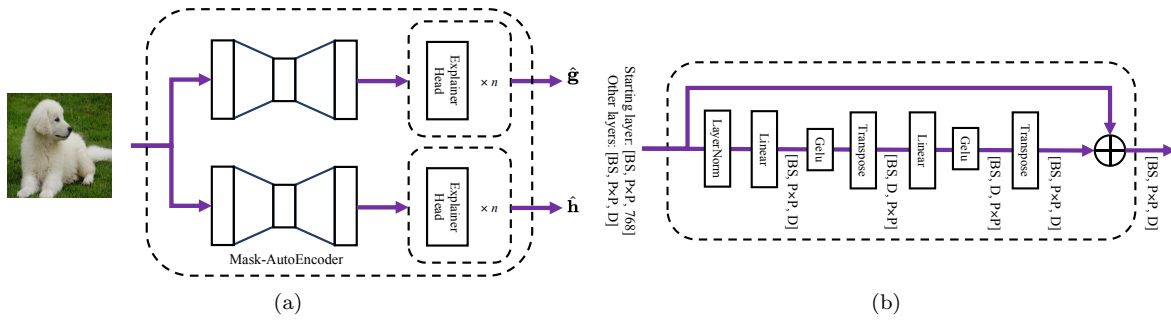


Figure 10: Structure of Mask-Autoencoder.

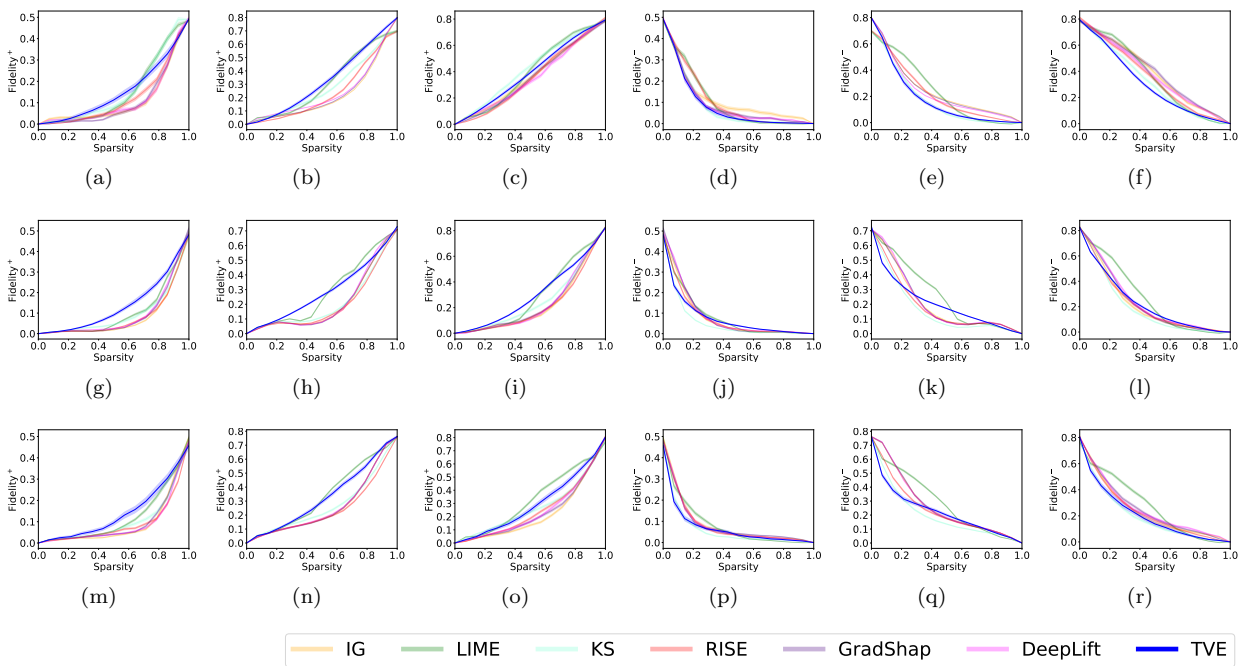


Figure 11: Fidelity⁺-sparsity curve for explaining ViT-Base on Cats-vs-dogs (a), Imagenette (a), and CIFAR-10 (c). Fidelity⁻-sparsity curve of ViT-Base on Cats-vs-dogs (d), Imagenette (e), and CIFAR-10 (f). Fidelity⁺-sparsity curve of Swin-Base on Cats-vs-dogs (g), Imagenette (h), and CIFAR-10 (i). Fidelity⁻-sparsity curve of Swin-Base on Cats-vs-dogs (j), Imagenette (k), and CIFAR-10 (l). Fidelity⁺-sparsity curve of Deit-Base on Cats-vs-dogs (m), Imagenette (n), and CIFAR-10 (o). Fidelity⁻-sparsity curve of Deit-Base on Cats-vs-dogs (p), Imagenette (q), and CIFAR-10 (r).

I Fidelity-Sparsity Curve of Section 6.2

We show the fidelity-sparsity curve for explaining ViT-Base, Swin-Base, and Deit-Base on the Cats-vs-dogs, Imagenette, and CIFAR-10 datasets in Figures 11 (a)-(r). It is observed that TVE consistently exhibits promising performance in terms of both Fidelity⁺(\uparrow) and Fidelity⁻(\downarrow), surpassing the majority of baseline methods. This indicates TVE’s ability to faithfully explain various downstream tasks.

J Computational Infrastructure

The computational infrastructure information is given in Table 5.

Table 5: Computing infrastructure for the experiments.

Device Attribute	Value
Computing Infrastructure	GPU
GPU Model	NVIDIA-A5000
GPU Memory	24564MB
GPU Number	8
CUDA Version	12.1
CPU Memory	512GB

K More Case Studies

We give more explanation heatmaps of ViT-Base on the ImageNet dataset in Figure 12, which are generated by TVE.

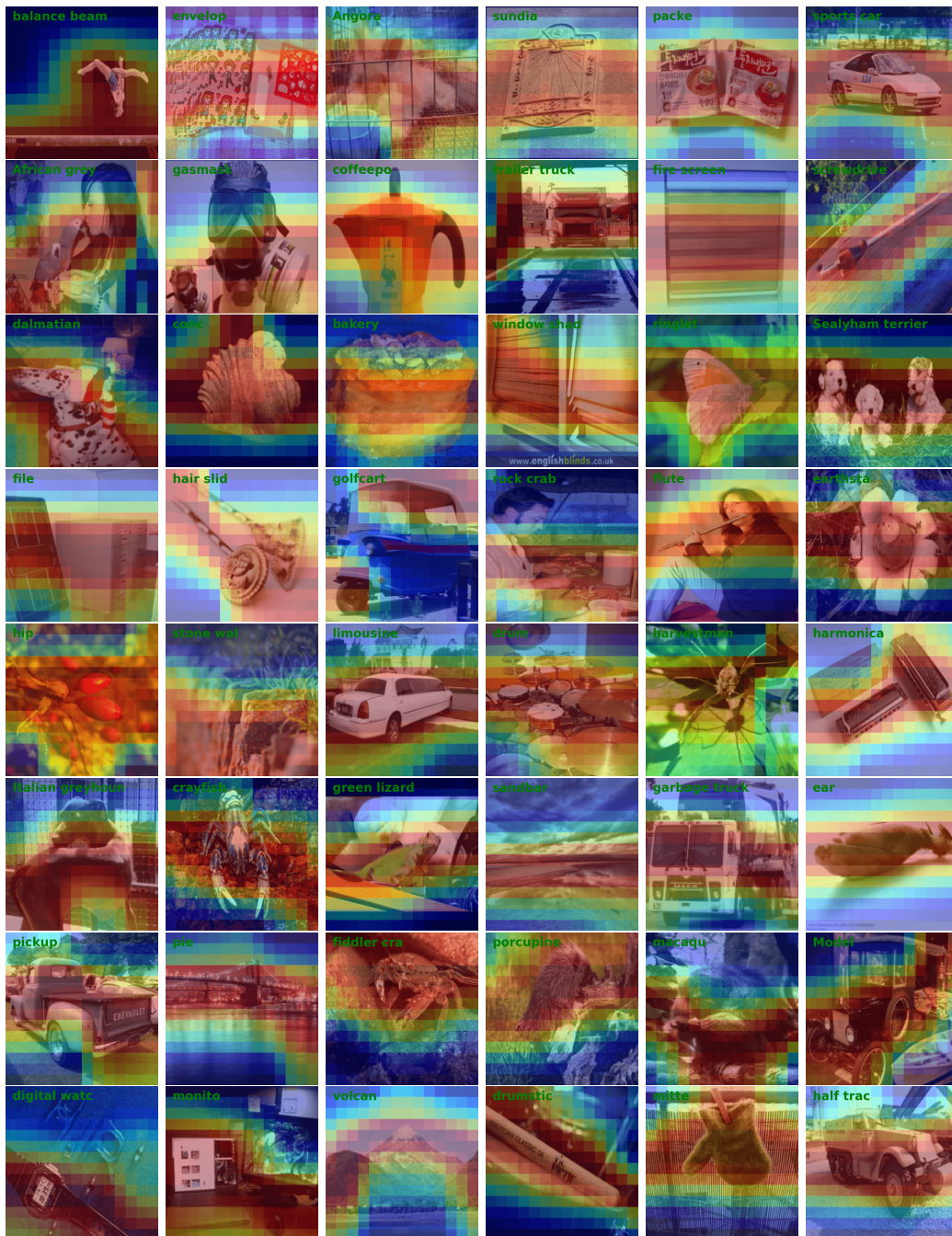


Figure 12: Explanation heatmaps of ViT-Base on the ImageNet dataset.