

# FastSHAP: Real-Time Shapley Value Estimation

**Neil Jethani\***  
 New York University  
 nj594@nyu.edu

**Mukund Sudarshan\***  
 New York University  
 ms7490@nyu.edu

**Ian Covert\***  
 University of Washington  
 icovert@cs.uw.edu

**Su-In Lee**  
 University of Washington  
 suinlee@cs.uw.edu

**Rajesh Ranganath**  
 New York University  
 rajeshr@cims.nyu.edu

## Abstract

Shapley values are widely used to explain black-box models, but they are costly to calculate because they require many model evaluations. We introduce FastSHAP, a method for estimating Shapley values in a single forward pass using a learned explainer model. FastSHAP amortizes the cost of explaining many inputs via a learning approach inspired by the Shapley value’s weighted least squares characterization, and it can be trained using standard stochastic gradient optimization. We compare FastSHAP to existing estimation approaches, revealing that it generates high-quality explanations with orders of magnitude speedup.

## 1 Introduction

Model explanation is becoming increasingly important with the proliferation of black-box models, and Shapley values [34] have emerged as a popular approach due to their strong theoretical properties [25, 40, 12, 27]. However, Shapley values have not seen widespread adoption in certain settings (e.g., with large vision and language models) due to computational challenges. Shapley values are inherently difficult to calculate, with a brute force calculation requiring an exponential number of model evaluations [42], so approximations are needed to make them useful in practice.

Recent work has addressed the computational challenges with Shapley values using two main approaches. First, many works have proposed stochastic estimators [3, 40, 39, 27, 10, 44] that rely on sampling either feature subsets or permutations; these are often consistent estimators, but they require many model evaluations and involve a trade-off between run-time and accuracy. Second, some works have proposed model-specific approximations, e.g., for trees [26] or neural networks [35, 6, 2, 43]; these are generally faster, but they sometimes require many model evaluations, often induce bias, and typically lack flexibility regarding how to handle held-out features when generating explanations—a subject of continued debate in the field [1, 19, 9, 14].

Here, we pursue a different approach to accelerating Shapley value calculation: to achieve the fastest possible run-time, we learn a separate explainer model that outputs precise Shapley value estimates in a single forward pass. Our approach is inspired by the Shapley value’s weighted least squares characterization [4], which allows us to train a model that predicts an approximate solution to a carefully chosen optimization problem—the same one that KernelSHAP must solve separately for each data example [27].

**Our contributions.** We introduce FastSHAP, a flexible, amortized approach for generating real-time Shapley value explanations.<sup>1</sup> We derive a learning objective from the Shapley value’s weighted least squares characterization and investigate several ways to reduce gradient variance during training.

---

\*Equal contribution

<sup>1</sup><https://git.io/JCqFV> (PyTorch), <https://git.io/JCqbP> (TensorFlow)

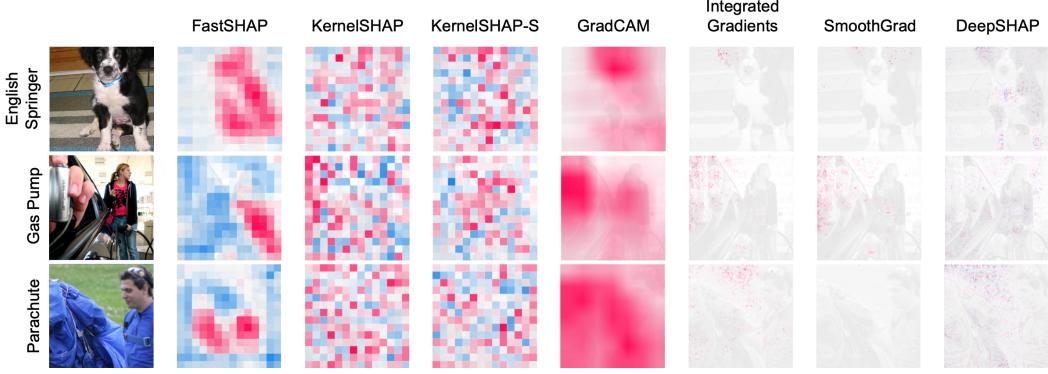


Figure 1: Explanations generated by each method for Imagenette images.

Our experiments show that FastSHAP produces accurate Shapley value estimates with orders of magnitude speedup relative to non-amortized estimation approaches, and we find that FastSHAP generates high-quality image explanations (fig. 1) that outperform gradient-based methods (e.g., IntGrad [41] and GradCAM [33]) on quantitative inclusion and exclusion metrics.

## 2 Background

In this section, we introduce notation used throughout the paper and provide an overview of Shapley values and their weighted least squares characterization. Let  $\mathbf{x} \in \mathcal{X}$  be a random vector consisting of  $d$  features, or  $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_d)$ , and let  $\mathbf{y} \in \mathcal{Y} = \{1, \dots, K\}$  be the response variable for a classification problem. We use  $\mathbf{s} \in \{0, 1\}^d$  to denote subsets of the indices  $\{1, \dots, d\}$  and define  $\mathbf{x}_s := \{\mathbf{x}_i\}_{i:s_i=1}$ . The symbols  $\mathbf{x}, \mathbf{y}, \mathbf{s}$  are random variables and  $x, y, s$  denote possible values. We use  $\mathbf{1}$  and  $\mathbf{0}$  to denote vectors of ones and zeros in  $\mathbb{R}^d$ , so that  $\mathbf{1}^\top s$  is the subset  $s$ 's cardinality, and we use  $e_i$  to denote the  $i$ th standard basis vector. Finally,  $f(\mathbf{x}; \eta) : \mathcal{X} \mapsto \Delta^{K-1}$  represents a model that outputs a probability distribution over  $\mathbf{y}$  given  $\mathbf{x}$ , and  $f_y(\mathbf{x}; \eta)$  denotes the probability for the  $y$ th class.

### 2.1 Shapley values

Shapley values were originally developed as a credit allocation method in cooperative game theory [34], but they have since been adopted as a tool for explaining predictions from black-box machine learning models [40, 12, 27]. For any value function (or set function)  $v : 2^d \mapsto \mathbb{R}$ , the Shapley values  $\phi(v) \in \mathbb{R}^d$ , or  $\phi_i(v)$  for each feature  $i = 1, \dots, d$  are given by the formula

$$\phi_i(v) = \frac{1}{d} \sum_{s_i \neq 1} \binom{d-1}{\mathbf{1}^\top s}^{-1} (v(s + e_i) - v(s)). \quad (1)$$

The difference  $v(s + e_i) - v(s)$  represents the  $i$ th feature's contribution to the subset  $s$ , and the summation represents a weighted average across all subsets that do not include  $i$ . In the model explanation context,  $v$  is chosen to be a value function that represents how an individual prediction varies as different subsets of features are removed. For example, given an input-output pair  $(x, y)$ , the value function  $v_{x,y}$  may be specified as

$$v_{x,y}(s) = \text{link} \left( \mathbb{E}_{p(\mathbf{x}_{1-s})} [f_y(x_s, \mathbf{x}_{1-s}; \eta)] \right), \quad (2)$$

where the held out features  $\mathbf{x}_{1-s}$  are marginalized out using their joint marginal distribution  $p(\mathbf{x}_{1-s})$  and a link function (e.g., logit) is applied to the model output. Recent work has debated the properties of different value function formulations, particularly the choice of how to remove features [1, 19, 9, 14]. Regardless of the formulation, this approach to model explanation enjoys several useful theoretical properties due to the use of Shapley values: for example, they are zero for irrelevant features and are guaranteed to sum to the model's prediction. We direct readers to prior work for a detailed discussion of these properties [27, 9].

Unfortunately, Shapley values also introduce computational challenges: the summation in eq. (1) involves an exponential number of subsets, making it infeasible to calculate for large  $d$ . Fast approximations are therefore required to make Shapley values practical, as we discuss next.

## 2.2 KernelSHAP

KernelSHAP [27] is a popular Shapley value implementation that relies on an alternative Shapley value interpretation. Given a value function  $v_{x,y}$ , eq. (1) shows that the Shapley values  $\phi(v_{x,y})$  are the features' weighted average contributions; equivalently, their weighted least squares characterization [4] says that they are the solution to the following optimization problem over  $\phi_{x,y} \in \mathbb{R}^d$ ,

$$\begin{aligned} \arg \min_{\phi_{x,y}} & \mathbb{E}_{p(\mathbf{s})} \left[ (v_{x,y}(\mathbf{s}) - v_{x,y}(\mathbf{0}) - \mathbf{s}^\top \phi_{x,y})^2 \right] \\ \text{s.t. } & \mathbf{1}^\top \phi_{x,y} = v_{x,y}(\mathbf{1}) - v_{x,y}(\mathbf{0}), \end{aligned} \quad (3)$$

(Efficiency constraint)

where the distribution  $p(\mathbf{s})$  is defined as

$$p(s) \propto \frac{d-1}{\binom{d}{\mathbf{1}^\top s} \cdot \mathbf{1}^\top s \cdot (d - \mathbf{1}^\top s)} \quad (\text{Shapley kernel})$$

for  $s$  such that  $0 < \mathbf{1}^\top s < d$ . Based on this view of the Shapley value, KernelSHAP is a stochastic estimator that solves an approximate version of eq. (3) given some number of subsets sampled from  $p(\mathbf{s})$ . While the estimator is consistent and empirically unbiased [8], KernelSHAP often requires many samples to achieve an accurate estimate, and it must solve eq. (3) separately for each input-output pair  $(x, y)$ ; as a result, it is unacceptably slow for some use-cases, particularly with models that are slow to evaluate. Our approach builds on KernelSHAP's motivation, leveraging the weighted least squares characterization to design a faster, amortized estimation approach.

## 3 FastSHAP

We now introduce FastSHAP, a method that amortizes the cost of generating Shapley values across many data samples. FastSHAP has two main advantages over existing approaches: (a) it avoids solving separate optimization problems for each input to be explained, and (b) it can use similar data points to efficiently learn the Shapley value function  $\phi$ .

### 3.1 Amortizing Shapley values

FastSHAP avoids calculating Shapley values separately for each  $(x, y)$  by using a single parametric function  $\phi_{\text{fast}}(\mathbf{x}, \mathbf{y}; \theta) : \mathcal{X} \times \mathcal{Y} \mapsto \mathbb{R}^d$ . Rather than solving separate optimization problems for each sample, we propose training FastSHAP to predict Shapley values in a single forward pass, where the predictions are penalized according to the weighted least squares objective in eq. (3). Specifically, we train  $\phi_{\text{fast}}(\mathbf{x}, \mathbf{y}; \theta)$  to minimize the following loss function:

$$\mathcal{L}(\theta) = \mathbb{E}_{p(\mathbf{x})} \mathbb{E}_{\text{Unif}(\mathbf{y})} \mathbb{E}_{p(\mathbf{s})} \left[ (v_{\mathbf{x}, \mathbf{y}}(\mathbf{s}) - v_{\mathbf{x}, \mathbf{y}}(\mathbf{0}) - \mathbf{s}^\top \phi_{\text{fast}}(\mathbf{x}, \mathbf{y}; \theta))^2 \right]. \quad (4)$$

Given a large enough dataset and a sufficiently expressive model class for  $\phi_{\text{fast}}$ , if the model's predictions are forced to satisfy the **Efficiency constraint**, then the global optimizer  $\phi_{\text{fast}}(\mathbf{x}, \mathbf{y}; \theta^*)$  will output exact Shapley values. Formally, the global optimizer satisfies the following:

$$\phi_{\text{fast}}(x, y; \theta^*) = \phi(v_{x,y}) \quad \forall x \in \mathcal{X}, y \in \mathcal{Y}. \quad (5)$$

To address the efficiency requirement, we explore two approaches. First, we can enforce efficiency by adjusting the Shapley value predictions using their *additive efficient normalization* [30], which applies the following operation to the model's outputs before evaluating the loss:

$$\phi_{\text{fast}}^{\text{eff}}(x, y; \theta) = \phi_{\text{fast}}(x, y; \theta) + \underbrace{\frac{1}{d} \left( v_{x,y}(\mathbf{1}) - v_{x,y}(\mathbf{0}) - \mathbf{1}^\top \phi_{\text{fast}}(x, y; \theta) \right)}_{\text{Efficiency gap}}. \quad (6)$$

The normalization step can be applied during training or only at inference time. In [appendix A](#), we prove that the additive efficient normalization is guaranteed to make the estimates closer to the true Shapley values. Second, we can relax the efficiency property by augmenting  $\mathcal{L}(\theta)$  with a penalty on the efficiency gap; the penalty requires a parameter  $\gamma > 0$ , and as  $\gamma \rightarrow \infty$  we can guarantee that efficiency holds. Below, [algorithm 1](#) summarizes our learning approach.

**Empirical considerations.** Optimizing  $\mathcal{L}(\theta)$  using a single set of samples  $(x, y, s)$  is problematic because of high variance in the gradients, which can lead to poor optimization. We therefore consider several steps to reduce the variance in the gradients. First, as is conventional in deep learning, we minibatch across multiple samples from  $p(\mathbf{x})$ . Next, in some cases we calculate the loss jointly across all classes  $y \in \{1, \dots, K\}$ . Then, we experiment with using multiple samples  $s \sim p(\mathbf{s})$  for each input sample  $x$ . Finally, we explore *paired sampling*, where each sample  $s$  is paired with its complement  $1 - s$ , which has been shown to reduce KernelSHAP's variance [\[8\]](#). We provide more information on these variance reduction steps in [appendix B](#) and show a more detailed algorithm in [appendix C](#).

### 3.2 A default value function for FastSHAP

FastSHAP has the flexibility to work with any value function  $v_{x,y}$ . Here, we describe a default value function that is useful for explaining predictions from a classification model.

The value function's aim is to assess, for each subset  $s$ , the classification probability when only the features  $\mathbf{x}_s$  are observed.

Because most models  $f(\mathbf{x}; \eta)$  do not support making predictions without all the features, we require an approximation that simulates the inclusion of only  $\mathbf{x}_s$  [\[9\]](#). To this end, we use a supervised surrogate model [\[14, 20\]](#) to approximate marginalizing out the remaining features  $\mathbf{x}_{1-s}$  using their conditional distribution.

Separate from the original model  $f(\mathbf{x}; \eta)$ , the *surrogate* model  $p_{\text{surr}}(\mathbf{y} \mid m(\mathbf{x}, \mathbf{s}); \beta)$  takes as input a vector of masked features  $m(\mathbf{x}, \mathbf{s})$ , where the masking function  $m$  replaces features  $x_i$  where  $s_i = 0$  with a `[mask]` value that is not in the support of  $\mathcal{X}$ . Similar to prior work [\[14, 20, 9\]](#), the parameters  $\beta$  are learned by minimizing the following loss function:

$$\mathcal{L}(\beta) = \mathbb{E}_{p(\mathbf{s})} \left[ D_{\text{KL}}(f(\mathbf{x}; \eta) \parallel p_{\text{surr}}(\mathbf{y} \mid m(\mathbf{x}, \mathbf{s}); \beta)) \right]. \quad (7)$$

It has been shown that the global optimizer to [eq. \(7\)](#),  $p_{\text{surr}}(\mathbf{y} \mid m(\mathbf{x}, \mathbf{s}); \beta^*)$  is equivalent to marginalizing out features from  $f(\mathbf{x}; \eta)$  with their conditional distribution [\[9\]](#):

$$p_{\text{surr}}(\mathbf{y} \mid m(\mathbf{x}, \mathbf{s}); \beta^*) = \mathbb{E}[f_y(\mathbf{x}; \eta) \mid \mathbf{x}_s = x_s]. \quad (8)$$

The choice of distribution over  $p(\mathbf{s})$  does not affect the global optimizer of [eq. \(7\)](#), but we use the *Shapley kernel* to put more weight on subsets likely to be encountered when training FastSHAP. We use the surrogate model as a default choice for two reasons: first, it requires a single prediction for each evaluation of  $v_{x,y}(s)$  (unlike marginalizing out features using their marginal distribution [\[27, 19\]](#)), which permits faster training. Second, it yields explanations that reflect the model's dependence on the *information* from each feature, rather than its *algebraic* dependence [\[9, 14\]](#).

## 4 Related Work

Recent work on Shapley value explanations has largely focused on how to remove features [\[1, 19, 9, 14\]](#) and how to approximate Shapley values efficiently [\[6, 2, 26, 8, 42\]](#). Model-specific approximations are relatively fast, but they often introduce bias and are entangled with specific

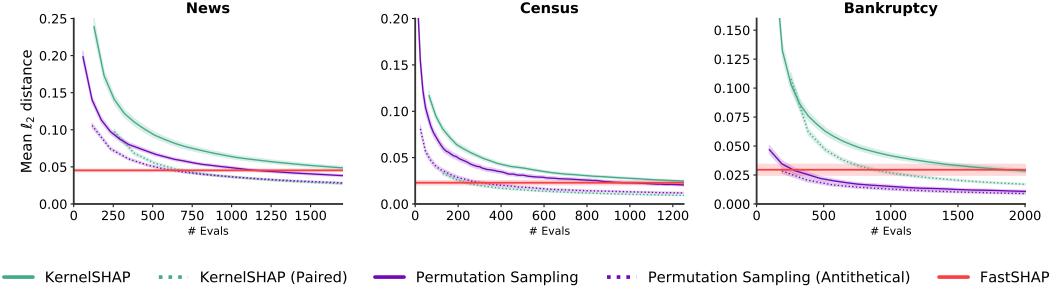


Figure 2: **To achieve Shapley value estimates comparable to FastSHAP, the baselines require between  $200 - 1200 \times$  more model evaluations.** Here we show how close each Shapley value approximator is to the ground truth as a function of the number of model evaluations. FastSHAP is represented by a horizontal line as it only requires a single forward pass of the  $\phi_{\text{fast}}$  network.

feature removal strategies [35, 27, 2, 26]. In contrast, model-agnostic stochastic approximations are more flexible, but they must trade off run-time and accuracy in the explanation. For example, KernelSHAP samples subsets to approximate the solution to a weighted least squares problem [27, 8], while other approaches sample marginal contributions [3, 40] or feature permutations [18, 28]. FastSHAP trains a model to output an estimate that would otherwise require hundreds of model evaluations, and unlike other fast approximations, it is agnostic to the model class and feature removal approach.

Other methods have been proposed to generate explanations using separate explainer models. These have been referred to as *amortized explanation methods* [9, 20], and they include several approaches [11, 5, 45, 32, 31, 20] that are comparable to gradient-based methods [37, 41, 38, 33] in terms of compute time. Notably, one approach generates a training dataset of ground truth explanations and then learns an explainer model to output explanations directly [32]—a principle that can be applied with any attribution method, at least in theory. In the case of Shapley values, generating a large training set would be very costly, so FastSHAP sidesteps the need for a training set using a stochastic optimization approach. Finally, FastSHAP can also be understood within the broader context of amortized optimization—an approach that has proven effective for variational inference [23, 36], style transfer [21], and in our case, model explanation.

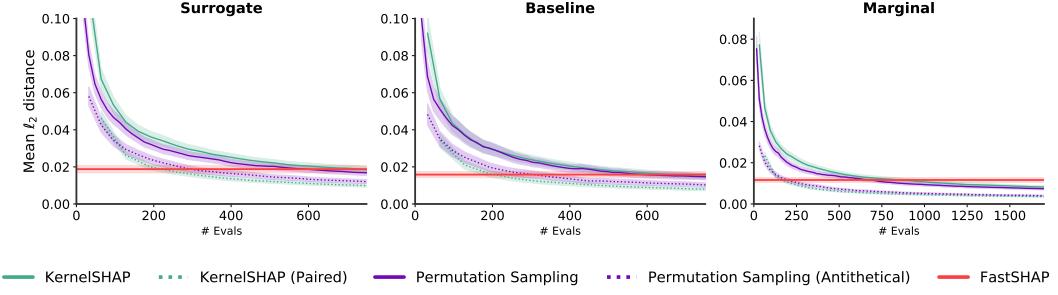
## 5 Structured Data Experiments

We analyze the performance of FastSHAP on multiple tabular datasets and compare it to several well-studied baselines. We first evaluate the accuracy of FastSHAP by comparing its outputs to ground truth Shapley values. Then, to disentangle the benefits of amortizing from using an in-distribution value function, we make the same comparisons when using different value function formulations  $v_{x,y}(\mathbf{s})$ . Unless otherwise stated, we use the value function introduced in section 3.2. Next, we discuss how to choose training hyperparameters for FastSHAP. Finally, in section 6 we test FastSHAP’s ability to explain images on three popular datasets.

**Baseline methods.** To contextualize FastSHAP’s approximation accuracy, we compare it to several non-amortized stochastic estimators. First, we compare to KernelSHAP [27] and its acceleration that uses paired sampling [8]. Next, we compare to a permutation sampling approach and its acceleration that uses antithetical sampling [28]. As a performance metric, we calculate the proximity to Shapley values that were obtained by running KernelSHAP to convergence; we view these as our ground truth because KernelSHAP is known to converge to the true Shapley values given infinite samples [8]. These baselines were all run using an open-source implementation.<sup>2</sup>

**Implementation details.** We use either neural networks or tree-based models for each of  $f(\mathbf{x}; \eta)$  and  $p_{\text{surr}}(\mathbf{y} \mid m(\mathbf{x}, \mathbf{s}); \beta)$ . The model  $\phi_{\text{fast}}(\mathbf{x}, \mathbf{y}; \theta)$  is implemented with a network  $g(\mathbf{x}; \theta) : \mathcal{X} \rightarrow \mathbb{R}^d \times \mathcal{Y}$  that outputs a vector of Shapley values for every  $y \in \mathcal{Y}$ . For brevity, we refer the reader to appendix D for exact implementation details including model classes, network architectures, and training hyperparameters.

<sup>2</sup><https://github.com/iancovert/shapley-regression/> (License: MIT)



**Figure 3: Given a limited computational budget, amortization improves the quality of Shapley value approximation regardless of the choice of value function.** Using the marketing dataset, we evaluate each Shapley value approximator using three different value functions: surrogate, baseline, and marginal.

We also perform a series of experiments to determine several training hyperparameters for FastSHAP in [appendix F](#) exploring (a) whether or not to use paired sampling, (b) the number of samples of  $s$  to take, and (c) how to best enforce the efficiency constraint. Based on the results, we use the following setting for all experiments in this section: we use paired sampling, we tune the number of samples  $s$  per  $x$  (typically using between 32 and 64 samples), we set  $\gamma = 0$ , and use the additive efficient normalization during both training and inference.

### 5.1 Accuracy of FastSHAP explanations

The estimated Shapley values must be close to the ground truth Shapley values for each data sample. We use four different datasets to compare each Shapley value estimation method in terms of their closeness to the ground truth.

**Setup.** Our experiments include data from a 1994 United States census, a bank marketing campaign, `bankruptcy` statistics, and numerical features about online news articles [13] (License: CC BY). The census data contains 49K samples with 12 input features like demographics, income, and information about a person’s job. The binary label for each sample indicates whether a person makes over \$50K a year. The marketing dataset contains 45K samples with 17 input features describing details of phone calls with customers, and the labels indicate whether the customer subscribed to a term deposit. The `bankruptcy` dataset contains 7K samples and 96 features describing various companies and whether or not the company went bankrupt. The news dataset contains 40K samples of 60 numerical features about various articles published on Mashable, and a label indicating the number of times an article was shared on social media; we binarize the label to be 1 if the share count is greater than the median number (1400), and 0 otherwise. The datasets are each split 80/20 for training and testing respectively.

**Results.** In [fig. 2](#), we show the distance of each method’s estimates to the ground truth as a function of the number of model evaluations required to compute Shapley values, for the news, census, and `bankruptcy` datasets. The first column of [fig. 3](#) shows the same plot for the marketing dataset. Each sample of  $s$  requires a forward pass of the  $f(\mathbf{x}; \eta)$  model, and since FastSHAP only requires a single forward pass of the  $\phi_{\text{fast}}$  network to output Shapley values, we show it as a horizontal line.

The permutation sampling baselines require around 200 model evaluations for the `census`, `marketing`, and `bankruptcy` datasets, but require over 500 for the `news` dataset to catch up to the performance of FastSHAP. With paired sampling KernelSHAP takes around 200 model evaluations for the `marketing` dataset, but without it requires close to 500. For the other three datasets, KernelSHAP takes over 1000 model evaluations to achieve performance comparable to that of FastSHAP. This equates to a roughly 200-1200x speed-up when explaining new instances of data with FastSHAP.

### 5.2 Disentangling choice of value function from amortization

In this experiment, we observe the effect of amortization across different choices of  $v_{x,y}(s)$ .

**Setup.** We re-use the `marketing` dataset for this experiment with the following three choices for the value function  $v_{x,y}(s)$ :

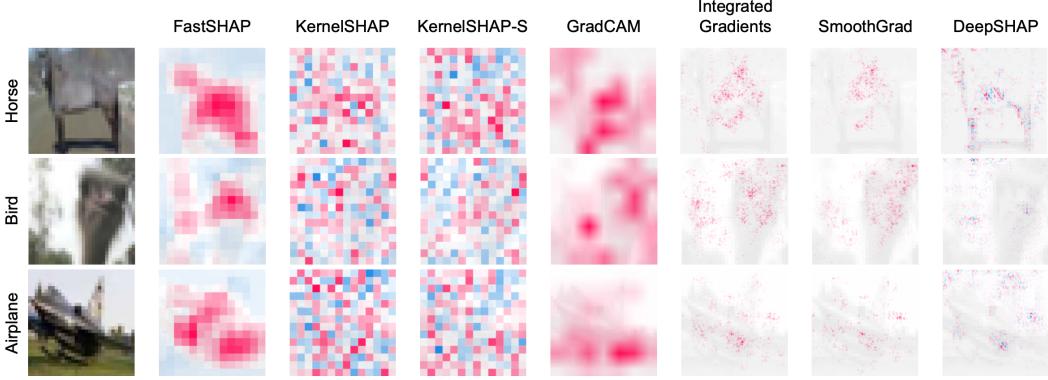


Figure 4: Explanations generated by each method for CIFAR-10 images.

**(Surrogate/In-distribution)** link ( $p_{\text{surr}}(y | m(x, s); \beta)$ )

**(Marginal/Out-of-distribution)** link ( $\mathbb{E}_{p(\mathbf{x}_{1-s})} [f_y(x_s, \mathbf{x}_{1-s}; \eta)]$ )

**(Baseline removal)** link ( $f_y(x_s, x_{1-s}^b; \eta)$ ), where  $x^b \in \mathcal{X}$  are fixed baseline values (the mean for continuous features and mode for discrete ones).

We compute ground truth Shapley values separately using each value function with the identity link. To generate estimates, we run both FastSHAP and KernelSHAP with each value function.

**Results.** Figure 3 highlights how amortization improves Shapley value estimation regardless of the choice of value function. Across the different value functions, KernelSHAP and permutation sampling require at least 200 model evaluations to achieve the same estimation quality as FastSHAP.

## 6 Image Experiments

Images yield a challenging benchmark for Shapley value estimators as they are high-dimensional and unstructured. We compare FastSHAP to KernelSHAP on two image datasets. Since gradient-based methods are the de-facto standard for explaining large classifiers, we also compare each Shapley value approximator to several widely-used gradient-based image explanation methods.

### 6.1 Datasets

We consider two popular image datasets for our experiments.

**CIFAR-10.** The CIFAR-10 dataset [24] (License: MIT) has 60,000,  $32 \times 32$  images across 10 classes. 50,000 samples are used for training, and 5,000 samples each are used for validation and testing. Each image is upscaled to  $224 \times 224$  using bilinear interpolation to interface with the ResNet50 architecture [15]. Examples are shown in fig. 4.

**Imagenette.** The Imagenette dataset [17] (License: Apache 2.0) is a subset of 10 classes from the Imagenet dataset and has 13,394 images. The data is split 9,469/1,963/1,962 for training, validating, and testing. Each image is cropped to keep only the  $224 \times 224$  central region. These images are divided into superpixels, just like the CIFAR-10 experiment. Examples are shown in fig. 1.

### 6.2 Explanations methods

We test three Shapley value estimators: FastSHAP, KernelSHAP, and DeepSHAP [27], a Shapley value computation method designed to work with images. With reference to our choices of value function in section 5.2, we implement KernelSHAP with the zeros baseline value function, which we refer to as KernelSHAP, and with the in-distribution surrogate value function, which we refer to as KernelSHAP-S. DeepSHAP applies DeepLIFT [35] to compute Shapley values. We also compare these estimators to popular gradient-based explanation methods including GradCAM [33], SmoothGrad [38], IntGrad [41]. Gradient-based methods generate explanations relatively quickly, and have therefore been widely adopted to explain image classifiers in practice.

**Implementation details.** First, each image is divided into a grid of  $14 \times 14$  superpixels. The models  $f(\mathbf{x}; \eta)$  and  $p_{\text{surr}}(\mathbf{y} | m(\mathbf{x}, \mathbf{s}); \beta)$  are both ResNet-50 models [15] pretrained on Imagenet and fine-tuned on the corresponding imaging dataset. For the  $\phi_{\text{fast}}(\mathbf{x}, \mathbf{y}; \theta)$  model, we use an identical pretrained model, but replace the last few layers with a convolutional layer so the output is  $14 \times 14 \times K$  such that super-pixel explanations are generated. We detail this procedure and other training hyperparameters in appendix E. To train  $\phi_{\text{fast}}(\mathbf{x}, \mathbf{y}; \theta)$  quickly, we only use a single paired sample of  $\mathbf{s}$  per  $\mathbf{x}$  and do not enforce efficiency. KernelSHAP and KernelSHAP-S are implemented using the `shap`<sup>3</sup> package’s default parameters. Implementations for GradCAM, SmoothGrad, and IntGrad are adopted from the `tf-explain`<sup>4</sup> package using the default hyperparameter settings.

### 6.3 Quantitative evaluation

**Setup.** Evaluating the quality of Shapley value estimates requires access to ground truth Shapley values, which is computationally infeasible for images. Instead, we introduce two quantitative metrics for evaluating image explanations, similar in spirit to several recent proposals [29, 16, 20]. These proposed metrics evaluate the classification accuracy of a model after including or excluding pixels according to their estimated importance.

Our evaluation procedure resembles that of Jethani et al. [20], who propose the EVAL-X metric as an alternative to retraining separate models on each set of image explanations [16]. We train a single evaluator model  $p_{\text{eval}}$  to approximate the  $f(\mathbf{x}; \eta)$  model’s output given only a subset of features. This procedure is analogous to the  $p_{\text{surr}}$  training procedure in section 3.2, except with the subset distribution set to  $p(\mathbf{s}) = \text{Uniform}(\{0, 1\}^d)$  to ensure all subsets are equally weighted.

A set of images is first labeled by the original model  $f(\mathbf{x}; \eta)$  using the most likely predicted class. We then generate explanations for this set using each method. Next, we use the explanations as feature rankings and compute the average top-1 accuracy (a measure of agreement with the original model) as a function of the number of pixels included/excluded. The area under each curve is termed the *Inclusion AUC* and *Exclusion AUC* respectively. Prior works have suggested [29, 16] that an accurate image explanation should both (a) maximally degrade the performance of  $p_{\text{eval}}$  when important features are excluded, and (b) maximally improve the performance of  $p_{\text{eval}}$  when important features are included.

Further, to use the same  $p_{\text{eval}}$  for each explanation method, explanations are provided only at the superpixel-level when possible. Otherwise, pixel-level explanations are coarsened using the sum total importance within a superpixel.

**Results.** Table 1 shows the Inclusion AUC and the Exclusion AUC achieved by each metric for both the CIFAR-10 and Imagenette datasets. We further present the curves used to generate these AUCs for Imagenette in fig. 5. Lower Exclusion AUCs and higher Inclusion AUCs are better. These results show that FastSHAP outperforms all baseline methods when evaluated with Exclusion AUC: when the pixels identified as important by FastSHAP are removed from the images, the sharpest decline in top-1 accuracy is observed. Further, FastSHAP performs well when evaluated on the basis

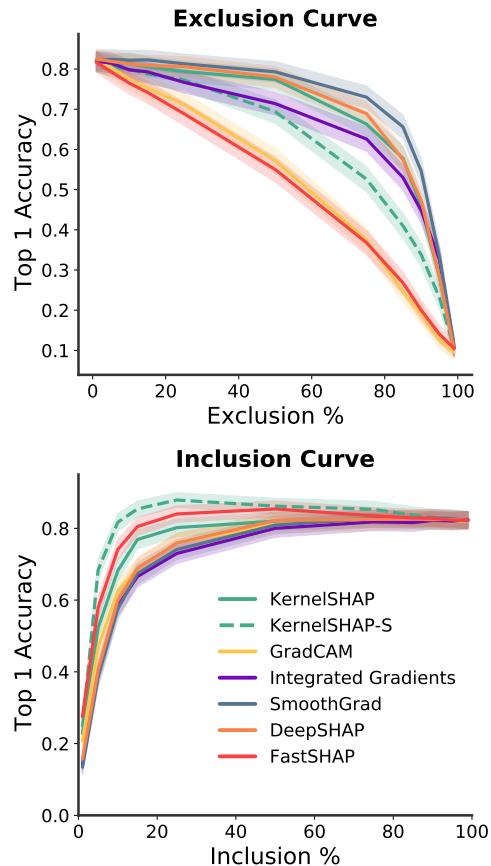


Figure 5: **Imagenette inclusion and exclusion curves.** The change in top-1 accuracy as an increasing percentage of the pixels estimated to be important are excluded (right) or included (left) from the set of 1000 images.

<sup>3</sup><https://shap.readthedocs.io/en/latest/> (License: MIT)

<sup>4</sup><https://tf-explain.readthedocs.io/en/latest/> (License: MIT)

Table 1: **Exclusion and Inclusion AUCs.** Evaluation of each method on the basis of Exclusion AUC (lower is better) and Inclusion AUC (higher is better).

	CIFAR-10		Imagenette	
	Exclusion AUC	Inclusion AUC	Exclusion AUC	Inclusion AUC
<b>FastSHAP</b>	<b>0.42 (0.41, 0.43)</b>	0.78 (0.77, 0.79)	<b>0.51 (0.49, 0.52)</b>	<b>0.79 (0.78, 0.80)</b>
KernelSHAP	0.64 (0.63, 0.65)	0.78 (0.77, 0.79)	0.68 (0.67, 0.70)	0.77 (0.75, 0.78)
KernelSHAP-S	0.54 (0.52, 0.55)	<b>0.86 (0.85, 0.87)</b>	0.61 (0.60, 0.62)	<b>0.82 (0.80, 0.83)</b>
GradCAM	0.52 (0.51, 0.53)	0.76 (0.75, 0.77)	<b>0.52 (0.50, 0.53)</b>	0.74 (0.73, 0.76)
Integrated Gradients	0.55 (0.54, 0.56)	0.74 (0.73, 0.75)	0.65 (0.64, 0.67)	0.73 (0.71, 0.74)
SmoothGrad	0.70 (0.69, 0.71)	0.72 (0.71, 0.73)	0.72 (0.71, 0.73)	0.73 (0.72, 0.75)
DeepSHAP	0.65 (0.64, 0.66)	0.79 (0.78, 0.80)	0.69 (0.68, 0.71)	0.74 (0.73, 0.75)

of Inclusion AUC, second only to KernelSHAP-S. We show more explanations for each method in [appendix G](#).

GradCAM performs competitively with FastSHAP on the basis of Exclusion AUC and KernelSHAP-S marginally outperforms FastSHAP on the basis of Inclusion AUC; this is perhaps surprising, given that KernelSHAP-S’s explanations have a high level of noise ([fig. 4](#)). However, it is unclear which metric is a more important measure of interpretability. While, GradCAM achieves low Exclusion AUC, it does not do as well when evaluated using Inclusion AUC, and the opposite relationship holds for KernelSHAP-S. Instead, the quality of the explanations returned by a method should be judged holistically on the basis of both metrics. Based on this intuition, FastSHAP produces the most versatile explanations among these methods.

#### 6.4 Speed evaluation

**Setup.** All imaging experiments were run using 8 cores of an Intel Xeon Gold 6148 Processor and a single NVIDIA Tesla V100. We recorded the clock time required to train each model and to explain the 1000 test set images for each explanation method. For FastSHAP and KernelSHAP-S we also report the time required to train the  $\phi_{\text{fast}}$  and  $p_{\text{surr}}$  models. We do not include the time used to train  $f(\mathbf{x}; \eta)$ , as it is the same across methods.

**Results.** [Table 2](#) shows the time (in minutes) required to train each method and the time taken to explain 1000 images. FastSHAP incurs a fixed training cost and very low marginal cost for each explanation. These results suggest that FastSHAP is well suited for real-time applications, where it is crucial to keep explanation times as low as possible. Further, when users wish to explain a large quantity of data, the low explanation cost can quickly make up for training times.

## 7 Discussion

In this work, we introduced FastSHAP, a method for estimating Shapley values in a single forward pass using a learned explainer model. We derive a learning objective from the Shapley value’s weighted least squares characterization and investigate several ways to reduce gradient variance during training. We observe that FastSHAP produces accurate Shapley value estimates with respect to the ground truth, while achieving a significant speedup over non-amortized approaches. We also find FastSHAP produces high-quality image explanations, outperforming popular gradient-based explanation methods on quantitative inclusion and exclusion metrics.

While Shapley values provide a useful theoretical grounding for image explanations, they are largely avoided by practitioners to explain image-scale models due to computational cost. Amortized Shapley value estimation can solve this problem by providing fast and high-quality explanations, as we observe in our experiments. Further, FastSHAP stands to benefit disproportionately compared to existing Shapley value methods as the state of machine learning advances, because its performance is directly tied to the quality of available predictive models and optimization techniques.

Table 2: Training and explaining times in minutes for 1000 images.

	CIFAR-10	Imagenette
Explain	FastSHAP	<b>0.04</b>
	KernelSHAP	453.69
	KernelSHAP-S	460.10
	GradCAM	0.38
	IntGrad	0.91
	SmoothGrad	1.00
	DeepSHAP	5.39
Train	FastSHAP	693.57
	KernelSHAP-S	362.03
		146.49
		73.22

## References

- [1] Aas, K., Jullum, M., and Løland, A. (2019). Explaining individual predictions when features are dependent: More accurate approximations to Shapley values. *arXiv preprint arXiv:1903.10464*.
- [2] Ancona, M., Oztireli, C., and Gross, M. (2019). Explaining deep neural networks with a polynomial time algorithm for Shapley value approximation. In *International Conference on Machine Learning*, pages 272–281. PMLR.
- [3] Castro, J., Gómez, D., and Tejada, J. (2009). Polynomial calculation of the Shapley value based on sampling. *Computers & Operations Research*, 36(5):1726–1730.
- [4] Charnes, A., Golany, B., Keane, M., and Rousseau, J. (1988). Extremal principle solutions of games in characteristic function form: core, Chebychev and Shapley value generalizations. In *Econometrics of Planning and Efficiency*, pages 123–133. Springer.
- [5] Chen, J., Song, L., Wainwright, M., and Jordan, M. (2018a). Learning to explain: An information-theoretic perspective on model interpretation. In *International Conference on Machine Learning*, pages 883–892. PMLR.
- [6] Chen, J., Song, L., Wainwright, M. J., and Jordan, M. I. (2018b). L-Shapley and C-Shapley: Efficient model interpretation for structured data. *arXiv preprint arXiv:1808.02610*.
- [7] Chen, T. and Guestrin, C. (2016). Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 785–794.
- [8] Covert, I. and Lee, S.-I. (2021). Improving KernelSHAP: Practical Shapley value estimation using linear regression. In *International Conference on Artificial Intelligence and Statistics*, pages 3457–3465. PMLR.
- [9] Covert, I., Lundberg, S., and Lee, S.-I. (2020a). Explaining by removing: A unified framework for model explanation. *arXiv preprint arXiv:2011.14878*.
- [10] Covert, I., Lundberg, S., and Lee, S.-I. (2020b). Understanding global feature contributions with additive importance measures. *Advances in Neural Information Processing Systems*, 33.
- [11] Dabkowski, P. and Gal, Y. (2017). Real time image saliency for black box classifiers. In *Advances in Neural Information Processing Systems*, pages 6967–6976.
- [12] Datta, A., Sen, S., and Zick, Y. (2016). Algorithmic Transparency via Quantitative Input Influence: Theory and Experiments with Learning Systems. In *Proceedings - 2016 IEEE Symposium on Security and Privacy, SP 2016*, pages 598–617. Institute of Electrical and Electronics Engineers Inc.
- [13] Dua, D. and Graff, C. (2017). UCI machine learning repository.
- [14] Frye, C., de Mijolla, D., Begley, T., Cowton, L., Stanley, M., and Feige, I. (2020). Shapley explainability on the data manifold. *arXiv preprint arXiv:2006.01272*.
- [15] He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778.
- [16] Hooker, S., Erhan, D., Kindermans, P.-J., and Kim, B. (2018). A benchmark for interpretability methods in deep neural networks. *arXiv preprint arXiv:1806.10758*.
- [17] Howard, J. and Gugger, S. (2020). Fastai: A layered API for deep learning. *Information*, 11(2):108.
- [18] Illés, F. and Kerényi, P. (2019). Estimation of the Shapley value by ergodic sampling. *arXiv preprint arXiv:1906.05224*.
- [19] Janzing, D., Minorics, L., and Blöbaum, P. (2020). Feature relevance quantification in explainable AI: A causal problem. In *International Conference on Artificial Intelligence and Statistics*, pages 2907–2916. PMLR.

- [20] Jethani, N., Sudarshan, M., Aphinyanaphongs, Y., and Ranganath, R. (2021). Have we learned to explain?: How interpretability methods can learn to encode predictions in their interpretations. In *International Conference on Artificial Intelligence and Statistics*, pages 1459–1467. PMLR.
- [21] Johnson, J., Alahi, A., and Fei-Fei, L. (2016). Perceptual losses for real-time style transfer and super-resolution. In *European Conference on Computer Vision*, pages 694–711. Springer.
- [22] Ke, G., Meng, Q., Finley, T., Wang, T., Chen, W., Ma, W., Ye, Q., and Liu, T.-Y. (2017). Lightgbm: A highly efficient gradient boosting decision tree. *Advances in neural information processing systems*, 30:3146–3154.
- [23] Kingma, D. P. and Welling, M. (2013). Auto-encoding variational Bayes. *arXiv preprint arXiv:1312.6114*.
- [24] Krizhevsky, A., Hinton, G., et al. (2009). Learning multiple layers of features from tiny images.
- [25] Lipovetsky, S. and Conklin, M. (2001). Analysis of regression in game theory approach. *Applied Stochastic Models in Business and Industry*, 17(4):319–330.
- [26] Lundberg, S. M., Erion, G., Chen, H., DeGrave, A., Prutkin, J. M., Nair, B., Katz, R., Himmelfarb, J., Bansal, N., and Lee, S.-I. (2020). From local explanations to global understanding with explainable AI for trees. *Nature Machine Intelligence*, 2(1):56–67.
- [27] Lundberg, S. M. and Lee, S.-I. (2017). A unified approach to interpreting model predictions. *Advances in Neural Information Processing Systems*, 30:4765–4774.
- [28] Mitchell, R., Cooper, J., Frank, E., and Holmes, G. (2021). Sampling permutations for Shapley value estimation. *arXiv preprint arXiv:2104.12199*.
- [29] Petsiuk, V., Das, A., and Saenko, K. (2018). RISE: Randomized input sampling for explanation of black-box models. *arXiv preprint arXiv:1806.07421*.
- [30] Ruiz, L. M., Valenciano, F., and Zarzuelo, J. M. (1998). The family of least square values for transferable utility games. *Games and Economic Behavior*, 24(1-2):109–130.
- [31] Schulz, K., Sixt, L., Tombari, F., and Landgraf, T. (2020). Restricting the flow: Information bottlenecks for attribution. *arXiv preprint arXiv:2001.00396*.
- [32] Schwab, P. and Karlen, W. (2019). CXPlain: Causal explanations for model interpretation under uncertainty. In *Advances in Neural Information Processing Systems*, pages 10220–10230.
- [33] Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., and Batra, D. (2017). Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 618–626.
- [34] Shapley, L. S. (1953). A value for n-person games. *Contributions to the Theory of Games*, 2(28):307–317.
- [35] Shrikumar, A., Greenside, P., and Kundaje, A. (2017). Learning important features through propagating activation differences. In *International Conference on Machine Learning*, pages 3145–3153. PMLR.
- [36] Shu, R., Bui, H. H., Zhao, S., Kochenderfer, M. J., and Ermon, S. (2018). Amortized inference regularization. *arXiv preprint arXiv:1805.08913*.
- [37] Simon, G. and Vincent, T. (2020). A projected stochastic gradient algorithm for estimating Shapley value applied in attribute importance. In *International Cross-Domain Conference for Machine Learning and Knowledge Extraction*, pages 97–115. Springer.
- [38] Smilkov, D., Thorat, N., Kim, B., Viégas, F., and Wattenberg, M. (2017). Smoothgrad: removing noise by adding noise. *arXiv preprint arXiv:1706.03825*.
- [39] Song, E., Nelson, B. L., and Staum, J. (2016). Shapley effects for global sensitivity analysis: Theory and computation. *SIAM-ASA Journal on Uncertainty Quantification*, 4(1):1060–1083.

- [40] Štrumbelj, E. and Kononenko, I. (2014). Explaining prediction models and individual predictions with feature contributions. *Knowledge and Information Systems*, 41(3):647–665.
- [41] Sundararajan, M., Taly, A., and Yan, Q. (2017). Axiomatic attribution for deep networks. In *International Conference on Machine Learning*, pages 3319–3328. PMLR.
- [42] Van den Broeck, G., Lykov, A., Schleich, M., and Suciu, D. (2021). On the tractability of SHAP explanations. In *Proceedings of the 35th Conference on Artificial Intelligence (AAAI)*.
- [43] Wang, R., Wang, X., and Inouye, D. I. (2021). Shapley explanation networks. In *International Conference on Learning Representations*.
- [44] Williamson, B. and Feng, J. (2020). Efficient nonparametric statistical inference on population feature importance using Shapley values. In *International Conference on Machine Learning*, pages 10282–10291. PMLR.
- [45] Yoon, J., Jordon, J., and van der Schaar, M. (2018). INVASE: Instance-wise variable selection using neural networks. In *International Conference on Learning Representations*.

## A Additive efficient normalization

Assuming a game  $v$  with Shapley values  $\phi(v) \in \mathbb{R}^d$ , consider the following optimization problem. Given Shapley values estimates  $\hat{\phi} \in \mathbb{R}^d$  that do not necessarily satisfy the efficiency property, we can project them onto the *efficient hyperplane*, or the subset of  $\mathbb{R}^d$  where the efficiency property is satisfied:

$$\min_{\phi_{\text{eff}} \in \mathbb{R}^d} \|\phi_{\text{eff}} - \hat{\phi}\|^2 \quad \text{s.t.} \quad \mathbf{1}^\top \phi_{\text{eff}} = v(\mathbf{1}) - v(\mathbf{0}).$$

We can solve the problem via its Lagrangian, denoted by  $\mathcal{L}(\phi_{\text{eff}}, \nu)$  with the Lagrange multiplier  $\nu \in \mathbb{R}$ , as follows:

$$\begin{aligned} \mathcal{L}(\phi_{\text{eff}}, \nu) &= \|\phi_{\text{eff}} - \hat{\phi}\|^2 + \nu(v(\mathbf{1}) - v(\mathbf{0}) - \mathbf{1}^\top \phi_{\text{eff}}) \\ &\Rightarrow \phi_{\text{eff}}^* = \hat{\phi} - \mathbf{1} \frac{v(\mathbf{1}) - v(\mathbf{0}) - \mathbf{1}^\top \hat{\phi}}{d} \end{aligned}$$

This transformation, where the efficiency gap is split evenly and added to each of the estimates, is known as the *additive efficient normalization* [30]. We implement it as an output layer for FastSHAP's predictions to ensure that they satisfy the efficiency property (section 3).

The normalization step is guaranteed to improve the Euclidean distance of any estimate  $\hat{\phi}$  to the true Shapley values  $\phi(v)$ , because they lie on the efficient hyperplane. This can be seen this from the following relationship:

$$\begin{aligned} \|\phi(v) - \hat{\phi}\|^2 &= \|\phi(v) - \phi_{\text{eff}}^*\|^2 + \|\phi_{\text{eff}}^* - \hat{\phi}\|^2 \\ &\geq \|\phi(v) - \phi_{\text{eff}}^*\|^2 \end{aligned}$$

## B Reducing gradient variance

Recall our learning objective  $\mathcal{L}(\theta)$ , which is inspired by the Shapley value's weighted least squares characterization:

$$\mathcal{L}(\theta) = \mathbb{E}_{p(\mathbf{x})} \mathbb{E}_{\text{Unif}(\mathbf{y})} \mathbb{E}_{p(\mathbf{s})} \left[ (v_{\mathbf{x}, \mathbf{y}}(\mathbf{s}) - v_{x, y}(\mathbf{0}) - \mathbf{s}^\top \phi_{\text{fast}}(\mathbf{x}, \mathbf{y}; \theta))^2 \right].$$

The objective's gradient is given by

$$\nabla_\theta \mathcal{L}(\theta) = \mathbb{E}_{p(\mathbf{x})} \mathbb{E}_{\text{Unif}(\mathbf{y})} \mathbb{E}_{p(\mathbf{s})} \left[ \nabla(\mathbf{x}, \mathbf{y}, \mathbf{s}; \theta) \right], \quad (9)$$

where we define

$$\nabla(x, y, s; \theta) := \nabla_\theta (v_{x, y}(s) - v_{x, y}(\mathbf{0}) - s^\top \phi_{\text{fast}}(x, y; \theta))^2.$$

When FastSHAP is trained with a single set of samples  $(x, y, s)$ , the gradient covariance is given by  $\text{Cov}(\nabla(\mathbf{x}, \mathbf{y}, \mathbf{s}; \theta))$ , which may be too large for effective optimization. We use several strategies to reduce gradient variance. First, given a model that outputs estimates for all  $y \in \{1, \dots, K\}$ , we calculate the loss for all classes jointly. This yields the gradients  $\nabla(\mathbf{x}, \mathbf{s}; \theta)$ , defined as

$$\nabla(x, s; \theta) := \mathbb{E}_{\text{Unif}(\mathbf{y})} [\nabla(x, \mathbf{y}, s; \theta)],$$

where we have the relationship

$$\text{Cov}(\nabla(\mathbf{x}, \mathbf{s}; \theta)) \preceq \text{Cov}(\nabla(\mathbf{x}, \mathbf{y}, \mathbf{s}; \theta))$$

due to the law of total covariance. Next, we consider minibatches of  $b$  independent  $x$  samples, which yields gradients  $\nabla_b(\mathbf{x}, \mathbf{s}; \theta)$  with covariance given by

$$\text{Cov}(\nabla_b(\mathbf{x}, \mathbf{s}; \theta)) = \frac{1}{b} \text{Cov}(\nabla(\mathbf{x}, \mathbf{s}; \theta)).$$

We then consider sampling  $m$  independent coalitions  $s$  for each input  $x$ , resulting in the gradients  $\nabla_b^m(\mathbf{x}, \mathbf{s}; \theta)$  with covariance given by

$$\text{Cov}(\nabla_b^m(\mathbf{x}, \mathbf{s}; \theta)) = \frac{1}{mb} \text{Cov}(\nabla(\mathbf{x}, \mathbf{s}; \theta)).$$

Finally, we consider a paired sampling approach, where each sample  $s \sim p(\mathbf{s})$  is paired with its complement  $1 - s$ . Paired sampling has been shown to reduce KernelSHAP's variance [8], and our experiments show that it helps FastSHAP achieve better Shapley value estimates.

## C Detailed algorithm

Here, we provide a more detailed algorithm that includes minibatching, multiple coalition samples, paired sampling, efficiency regularization, and parallelization over all output classes.

---

### Algorithm 2: FastSHAP training with all outputs

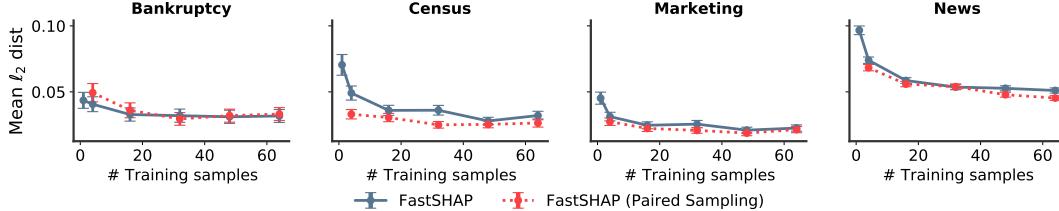
---

**Input:** Value function  $v_{x,y}$ , learning rate  $\alpha$ , batch size  $b$ , samples  $m$ , penalty parameter  $\gamma$   
**Output:** FastSHAP explainer  $\phi_{\text{fast}}(\mathbf{x}, \mathbf{y}; \theta)$   
 initialize  $\phi_{\text{fast}}(\mathbf{x}, \mathbf{y}; \theta)$   
**while** not converged **do**  
 | set  $\mathcal{R} \leftarrow 0$ ,  $\mathcal{L} \leftarrow 0$   
 | **for**  $i = 1, \dots, b$  **do**  
 | | sample  $x \sim p(\mathbf{x})$   
 | | **for**  $y = 1, \dots, K$  **do**  
 | | | predict  $\hat{\phi} \leftarrow \phi_{\text{fast}}(x, y; \theta)$   
 | | | calculate  $\mathcal{R} \leftarrow \mathcal{R} + (v_{x,y}(\mathbf{1}) - v_{x,y}(\mathbf{0}) - \mathbf{1}^\top \hat{\phi})^2$  // Pre-normalization  
 | | | **if** normalize **then**  
 | | | | set  $\hat{\phi} \leftarrow \hat{\phi} + d^{-1} (v_{x,y}(\mathbf{1}) - v_{x,y}(\mathbf{0}) - \mathbf{1}^\top \hat{\phi})$   
 | | | **end**  
 | | | **for**  $j = 1, \dots, m$  **do**  
 | | | | **if** paired sampling **and**  $i \bmod 2 = 0$  **then**  
 | | | | | set  $s \leftarrow 1 - s$  // Invert previous subset  
 | | | | **else**  
 | | | | | sample  $s \sim p(\mathbf{s})$   
 | | | | **end**  
 | | | | calculate  $\mathcal{L} \leftarrow \mathcal{L} + (v_{x,y}(s) - v_{x,y}(\mathbf{0}) - s^\top \hat{\phi})^2$   
 | | | **end**  
 | | **end**  
 | **end**  
 | update  $\theta \leftarrow \theta - \alpha \nabla_\theta (\frac{\mathcal{L}}{bmK} + \gamma \frac{\mathcal{R}}{bK})$   
**end**

---

## D FastSHAP structured data model setup

For the model  $f(\mathbf{x}; \eta)$  we use a neural network on the news datasets and gradient boosted trees on the census (LightGBM [22]) and bankruptcy (XGBoost [7]) datasets. The FastSHAP model  $\phi_{\text{fast}}$  and the value functions surrogate model  $p_{\text{surr}}(\mathbf{y} \mid m(\mathbf{x}, \mathbf{s}); \beta)$  use neural networks, consisting



**Figure 6: For a fixed computational budget, paired sampling improves the solution quality of FastSHAP.** Here we observe the mean  $\ell_2$  distance to ground truth of estimated Shapley values as a function of the number of times  $s$  is sampled per instance  $x$ .

of 2-3 fully connected layers with ReLU activations. Each fully connected layer consists of 128 hidden units. The  $p_{\text{surr}}$  models uses a Softmax output layer, while  $\phi_{\text{fast}}$  has no output activation. The models are trained using Adam with a learning rate of  $10^{-3}$ . We use a learning rate scheduler that multiplies the learning rate by a factor of 0.95 after 3 epochs of no validation loss improvement. Early stopping was triggered after the validation loss ceased to improve for 10 epochs.

## E FastSHAP image data model setup

The models  $f(\mathbf{x}; \eta)$  and  $p_{\text{surr}}$  are ResNet-50 models pretrained on Imagenet. We use these without modification to the architecture, and fine-tune them on each imaging dataset. To create the  $\phi_{\text{fast}}(\mathbf{x}, \mathbf{y}; \theta)$  model, we modify the architecture so a tensor of size  $14 \times 14 \times K$  is returned. First, the layers after the 4th convolutional block are removed. The output of this block is  $14 \times 14 \times 256$ . We append a 2D convolutional layer with  $K$  filters, each of size  $1 \times 1$ . The output is therefore  $14 \times 14 \times K$ , where the  $y$ th  $14 \times 14$  slice of this tensor corresponds to the superpixel-level Shapley values for class  $y \in gY$ . Each model is trained using Adam with a learning rate of  $10^{-3}$ . We use a learning rate scheduler that multiplies the learning rate by a factor of 0.8 after 3 epochs of no validation loss improvement. Early stopping was triggered after the validation loss ceased to improve for 20 epochs.

## F Hyperparameter choices for FastSHAP

FastSHAP allow users to trade speed for accuracy using various training hyperparameters. In this section, we explore various settings of these hyperparameters and observe their impact on the performance of FastSHAP.

**Setup.** There are two types of hyperparameters: sampling hyperparameters, which affect the number of samples of  $s$  taken during training, and efficiency hyperparameters, which enforce the **Efficiency constraint**. Sampling hyperparameters include: (a) whether or not to use paired samples, and (b) the number of samples of  $s$  per  $x$  to take during training. Efficiency hyperparameters include: (a) the choice of  $\gamma$  in [eq. \(4\)](#), (b) when to perform the additive efficient normalization. Additive efficient normalization can occur both during training and during inference, only at inference, or not at all.

To understand the effect of sampling hyperparameters, we use all of our structured data datasets. Due to computational constraints, we test the effect of efficiency hyperparameters using only the `census` and `bankruptcy` datasets. We use the in-distribution value function  $p_{\text{surr}}$  with FastSHAP. The ground truth SHAP values are computed the same was as in our previous experiments.

**Results.** [Figure 6](#) shows the mean  $\ell_2$  distance between FastSHAP and the ground truth. The red line indicates paired sampling. We notice that across all four datasets, increasing the number of training samples of  $s$  generally improves the mean  $\ell_2$  distance to ground truth. We also notice that given any fixed number of samples (greater than 1), paired sampling achieves better performance than without.

[Table 3](#) shows the results of an ablation study using the efficiency hyperparameters. Normalization refers to the additive efficient normalization step during both training and inference. The penalty refers to the efficiency penalty term  $\mathcal{R}$  in [algorithm 2](#). We notice that normalization during training uniformly achieves better results than without normalization or normalization only during inference. Choosing the efficiency regularization coefficient is difficult, and sometimes can worsen the solution

as shown in [table 3](#). As such, we avoid the penalty in our experiments, but use additive efficient normalization.

**Table 3: FastSHAP ablation results.** The distance to the ground truth SHAP values is displayed for several FastSHAP variations, showing that normalization helps and that the penalty is unnecessary.

	Census		Bankruptcy	
	$\ell_2$	$\ell_1$	$\ell_2$	$\ell_1$
Normalization	<b>0.0229</b>	<b>0.0863</b>	<b>0.0295</b>	<b>0.2436</b>
Normalization + Penalty	0.0261	0.0971	0.0320	0.2740
Inference Norm.	0.0406	0.1512	0.0407	0.3450
Inference Norm. + Penalty	0.0452	0.1671	0.0473	0.4471
No Norm.	0.0501	0.1933	0.0408	0.3474
No Norm. + Penalty	0.0513	0.1926	0.0474	0.4490

## G FastSHAP Image Explanations

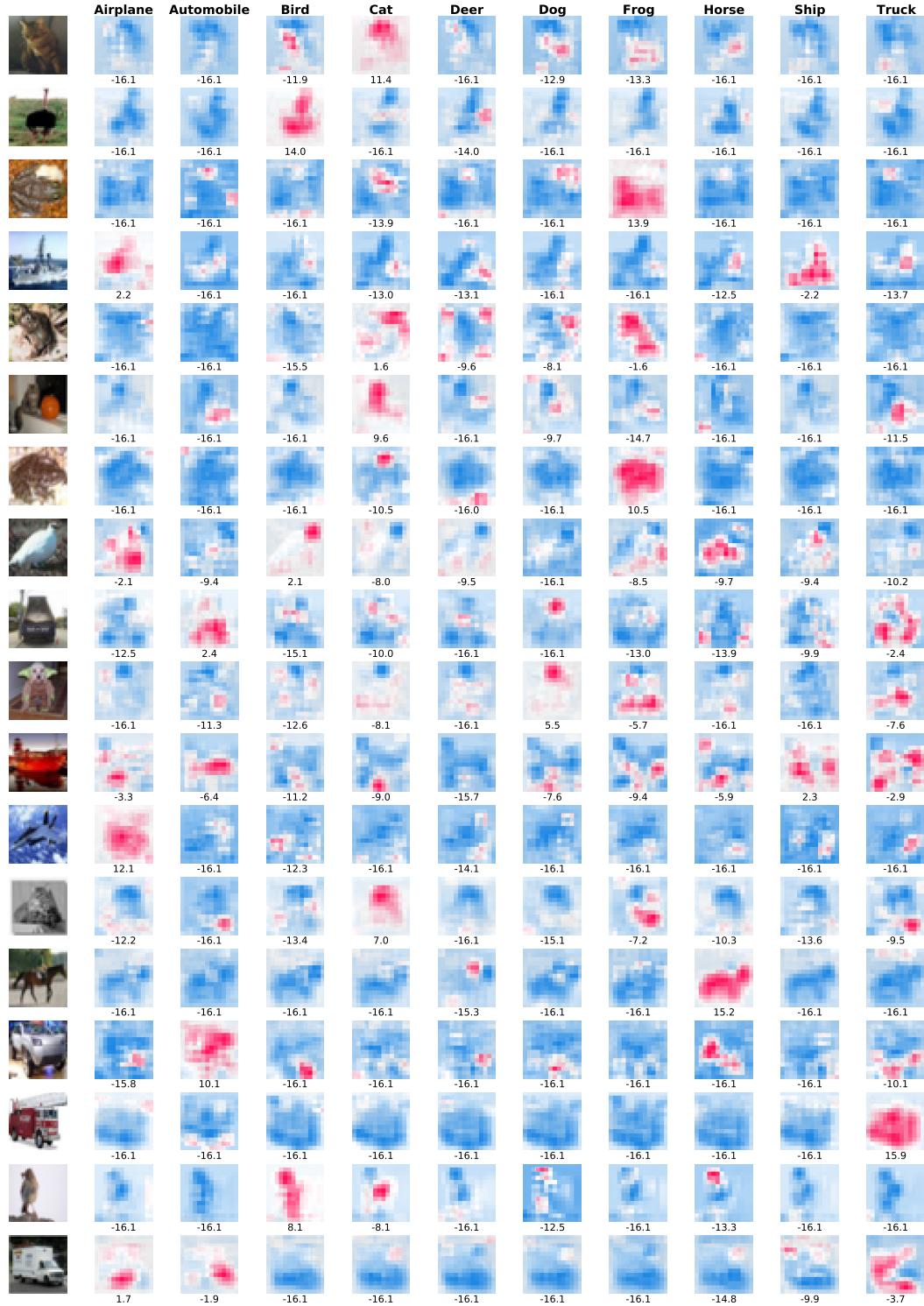


Figure 7: Explanations generated by FastSHAP for 18 randomly selected CIFAR-10 images. Shapley value explanations are plotted for each CIFAR-10 class. Each column is entitled with the corresponding class. The model's prediction (in logits) is provided below each image.

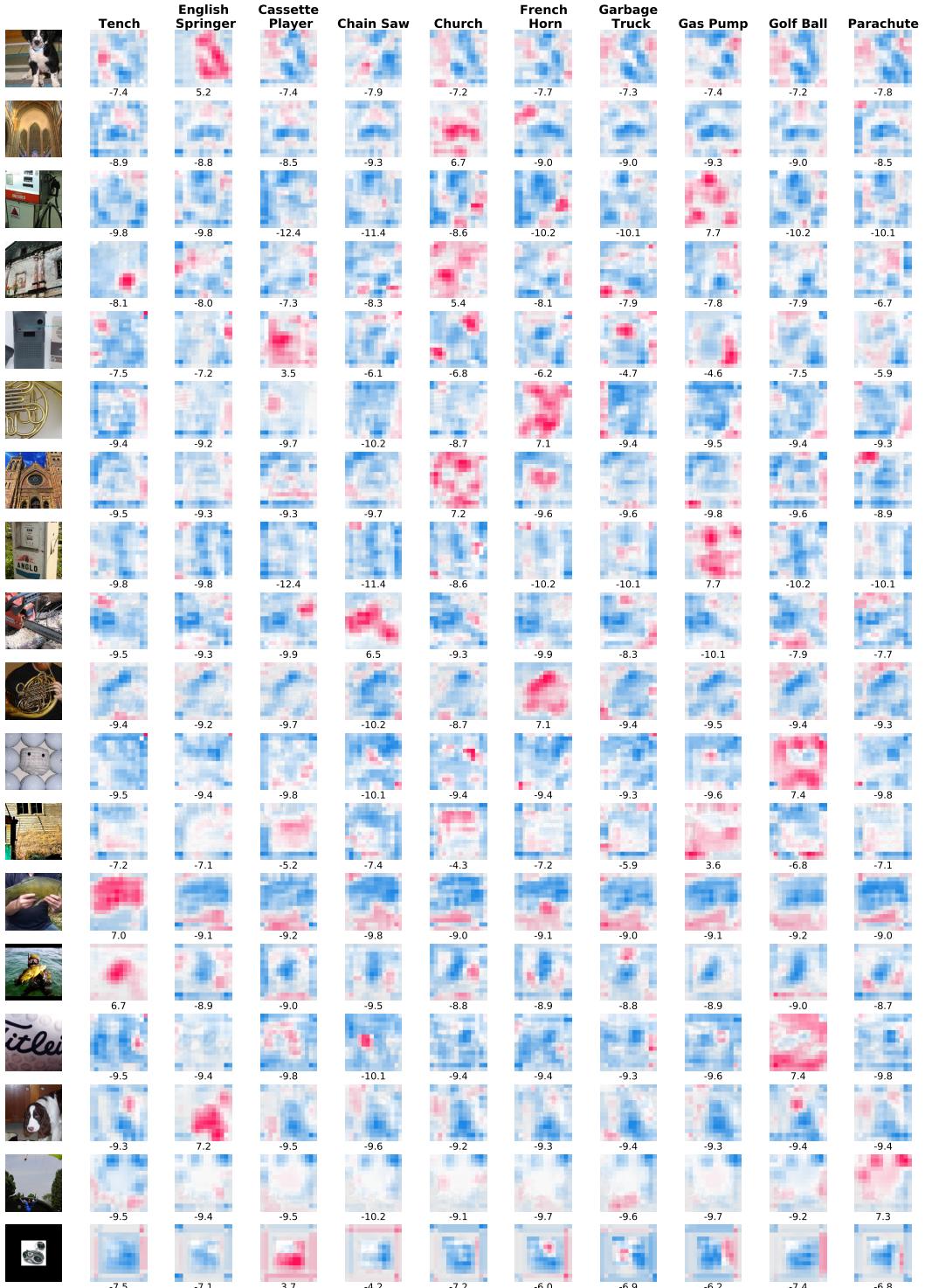


Figure 8: **Explanations generated by FastSHAP for 18 randomly selected Imagenette images.** Shapely value explanations are plotted for each CIFAR-10 class. Each column is entitled with the corresponding class. The model's prediction (in logits) is provided below each image.



**Figure 9: Explanations generated about the predicted class for 15 randomly selected CIFAR-10 images.** Each row is entitled with the corresponding class, and each column is entitled which the method used to generate the explanation.



**Figure 10: Explanations generated about the correct class for 15 randomly selected Imagenette images.** Each row is entitled with the corresponding class, and each column is entitled which the method used to generate the explanation.