

Selective Explanations

Lucas Monteiro Paes¹ *, Dennis Wei², Flavio P. Calmon¹

¹Harvard University

²IBM Research

Abstract

Feature attribution methods explain black-box machine learning (ML) models by assigning importance scores to input features. These methods can be computationally expensive for large ML models. To address this challenge, there has been increasing efforts to develop *amortized explainers*, where a machine learning model is trained to predict feature attribution scores with only one inference. Despite their efficiency, amortized explainers can produce inaccurate predictions and misleading explanations. In this paper, we propose *selective explanations*, a novel feature attribution method that (i) detects when amortized explainers generate low-quality explanations and (ii) improves these explanations using a technique called *explanations with initial guess*. Our selective explanation method allows practitioners to specify the fraction of samples that receive explanations with initial guess, offering a principled way to bridge the gap between amortized explainers and their high-quality counterparts.

1 Introduction

Large black-box models are increasingly used to support decisions in applications ranging from online content moderation [1], hiring [2], and medical diagnostics [3]. In such high-stakes settings, the need for explaining “why” a model produces a given output has led to a growing number of perturbation-based *feature attribution* methods [4–9]. Broadly speaking, these methods use input perturbations to assign numerical values to each input feature a model uses, indicating their influence on predictions. They are widely adopted in part because they work in the black-box setting with access only to model outputs (i.e., no gradients). However, existing feature attribution methods can be prohibitively expensive for the large models used in the current machine learning landscape (e.g., language models with billions of parameters) since they require a significant number of inferences for each individual explanation.

Recent literature has introduced two main strategies for speeding up feature attribution for large models: (i) employing Monte Carlo methods to approximate explanations with fewer computations [4, 5, 10, 11], and (ii) adopting an *amortized* approach, training a separate model to “mimic” the outputs of a reference explanation method [12–17]. Monte Carlo approximations can yield high-quality explanations but may converge slowly, limiting their practicality for large datasets. Amortized explainers, in turn, require only one inference per explanation, making them efficient for large black-box models and datasets. However, as demonstrated in Figure 1, amortized explainers can occasionally produce diverging explanations from the reference explainer used to train them.

We propose *selective explanation*, a method that bridges Monte Carlo and amortized explanations. By training a model that “learns to select” which method should be applied to each input, our

*Correspondence to Lucas Monteiro Paes (lucaspaes@g.harvard.edu).

<pre>[CLS] Better than most chain pizza , it's ok . \n\nWe got a thin crust , which was nice and crispy , only a little greasy , ok ingredients , not amazing on the cheese , and had kind of a bland crust \n\nI guess that doesn't sound too good . but I really promise it's better than any national chain pizza you'll find in town , It also has a really friendly . laid-back atmosphere . [SEP]</pre>	<pre>[CLS] Better than most chain pizza , it's ok . \n\nWe got a thin crust , which was nice and crispy , only a little greasy , ok ingredients , not amazing on the cheese , and had kind of a bland crust \n\nI guess that doesn't sound too good . but I really promise it's better than any national chain pizza you'll find in town , It also has a really friendly . laid-back atmosphere . [SEP]</pre>	<pre>[CLS] Better than most chain pizza , it's ok . \n\nWe got a thin crust , which was nice and crispy , only a little greasy , ok ingredients , not amazing on the cheese , and had kind of a bland crust \n\nI guess that doesn't sound too good . but I really promise it's better than any national chain pizza you'll find in town , It also has a really friendly . laid-back atmosphere . [SEP]</pre>
---	---	---

(a) Amortized (MSE = 0.31) (b) High-Quality (target) (c) Selective (MSE = 0.07)

Fig. 1: Amortized explainer (a) compared with a high-quality explainer (b) and our selective explanation method (c). All methods flag inputs that attribute why YelpLLM predicted **Negative Review** in the example. We observe that both high-quality and selective explanations attribute "not amazing" for the negative review (blue), while the amortized explainer misses this term. Similarly, the amortized explainer incorrectly the expression "Better than."

selective explanation method can produce higher-quality explanations than amortized explainers at a significantly lower average computational cost than Monte Carlo-based approaches. The key idea behind the selective explanation method is to apply Monte Carlo explanations only to points that would receive low-quality explanations from the amortized explainer; see Figure 2 for the workflow of selective explanations.

The ideas of predicting selectively and providing recourse with a more accurate but expensive method have been explored in classification and regression [18–22]. To our knowledge, however, these ideas have not been applied to explanations. We make **two contributions** in this regard that are relevant for selective prediction more generally. (1) Selective prediction uses an *uncertainty metric* to identify input points for which the predictor (the amortized explainer in our case) would produce low-quality outputs and recourse is needed. The high-dimensional nature of explanations requires us to develop new uncertainty metrics (Section 3) suitable for this setting. (2) Instead of providing recourse with a Monte Carlo explanation alone, as would be standard, we use an optimized method called *explanations with initial guess* (Section 4) that combines amortized and Monte Carlo explanations, improving explanation quality beyond that of either explanation alone.

Our **overall contribution** (3) is to combine (1) and (2) in the form of *selective explanations*, providing explanations with initial guess to improve low-quality amortized explanations. We validate our selective explanations approach on two language models as well as tabular datasets, demonstrating its ability to accurately detect low-quality explanations, enhance amortized explanations with even low-quality Monte Carlo explanations, and improve the worst explanations from the amortized model.

2 Problem Setup & Background

We aim to explain the predictions of a fixed probabilistic black-box model h that predicts $h(\mathbf{x}) = (h_1(\mathbf{x}), \dots, h_{|y|}(\mathbf{x}))$ and outputs $\operatorname{argmax}_{j \in y} h_j(\mathbf{x}) \in \mathcal{Y}$ using a vector of features $\mathbf{x} = (x_1, \dots, x_d) \in \mathbb{R}^d$. The user specifies an output of interest $\mathbf{y} \in \mathcal{Y}$ (usually $\mathbf{y} = \operatorname{argmax}_{j \in \mathcal{Y}} h_j(\mathbf{x})$) and our goal is to explain *Why would h output \mathbf{y} for a given \mathbf{x} ?* We consider a dataset $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$ comprised of $N > 0$ samples divided into three parts: $\mathcal{D}_{\text{train}}$ for training h and the explainers, \mathcal{D}_{cal} for calibration and validation, and $\mathcal{D}_{\text{test}}$ for testing. Thus, $\mathcal{D} = \mathcal{D}_{\text{train}} \cup \mathcal{D}_{\text{cal}} \cup \mathcal{D}_{\text{test}}$. Moreover, for a subset $S = \{i_1, \dots, i_{|S|}\} \subset [d]$ we write $\mathbf{x}_S \triangleq (x_{i_1}, \dots, x_{i_{|S|}})$.

Feature Attribution Methods, also called *explainers*, are functions $\mathbb{R}^d \times \mathcal{Y} \rightarrow \mathbb{R}^d$ that assess

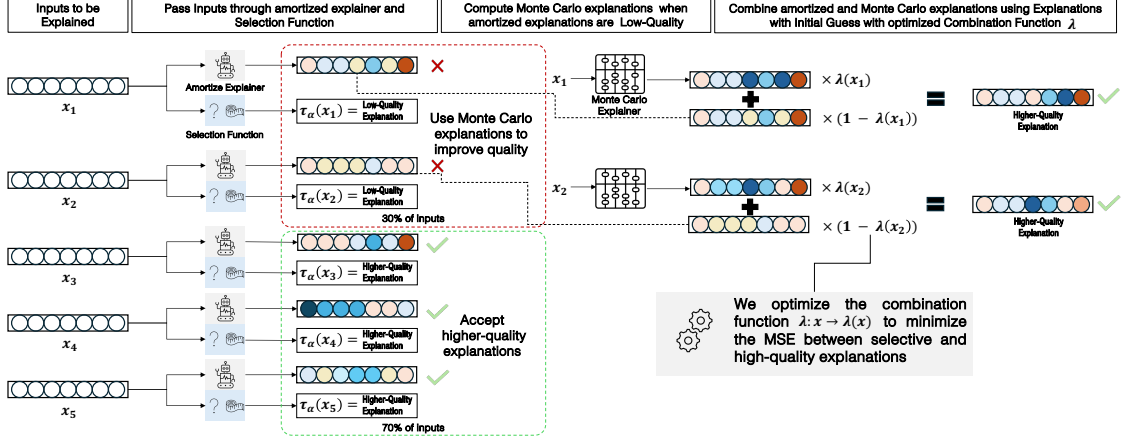


Fig. 2: Workflow of selective explanations.

the importance of each feature for the model’s (h) prediction to be \mathbf{y} for a given input vector \mathbf{x} . We consider three types of explainers:

- (i) **High-quality explainers** that use a large number of computations to provide explanations (e.g., SHAP with 2^d inferences from model h) [4, 5], denoted by $\text{HQ}(\mathbf{x}, \mathbf{y})$;
- (ii) **Monte Carlo explainers** that approximate high-quality explainers using n inferences from model h per explanation [4, 11], denoted by $\text{MC}^n(\mathbf{x}, \mathbf{y})$;
- (iii) **Amortized explainer** trained to approximate the high-quality explanations using only one inference [13, 14], denoted by $\text{Amor}(\mathbf{x}, \mathbf{y})$.

We measure the difference between two competing explanations using a loss (or distortion) function $\ell: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$, e.g., mean square error (MSE). The goal of selective explanations (SE) is to approximate high-quality explanations while minimizing the number of computations, i.e., to minimize $\|\text{SE}(\mathbf{x}, \mathbf{y}) - \text{HQ}(\mathbf{x}, \mathbf{y})\|_2^2$. However, computing high-quality explanations for large models h can be prohibitively expensive. To address this issue, we define *selective explainers* below.

Definition 1 (Selective Explainer). For a given model h , an amortized explainer Amor , a Monte Carlo explainer MC^n , a *combination function* $\lambda_h: \mathbb{R}^d \rightarrow \mathbb{R}$, and a *selection function* $\tau_\alpha: \mathbb{R}^d \rightarrow \{0, 1\}$ (parametrized by α), we define the *selective explainer* $\text{SE}(\mathbf{x}, \mathbf{y})$ as

$$\text{SE}(\mathbf{x}, \mathbf{y}) \triangleq \begin{cases} \text{Amor}(\mathbf{x}, \mathbf{y}) & , \text{ if } \tau_\alpha(\mathbf{x}) = 1, \\ \lambda_h(\mathbf{x})\text{Amor}(\mathbf{x}, \mathbf{y}) + (1 - \lambda_h(\mathbf{x}))\text{MC}^n(\mathbf{x}, \mathbf{y}) & , \text{ if } \tau_\alpha(\mathbf{x}) = 0. \end{cases} \quad (1)$$

When $\tau_\alpha = 0$, selective explanations output *explanations with initial guess* (Definition 2). Explanations with initial guess optimally linearly combine amortized and Monte Carlo explanations to leverage information from both and provide higher-quality explanations than either explainer alone. Selective explanations heavily depend on three objects that we define in this work: (i) an uncertainty metric (Section 3), (ii) a selection function (Section 3), and (iii) a combination function (Section 4).

- **Uncertainty metrics** (s_h) output the likelihood of the amortized explainer producing a low-quality explanation for an input. Lower $s_h(\mathbf{x})$ indicates a higher-quality explanation for \mathbf{x} . We propose two uncertainty metrics: Deep and Learned Uncertainty (Section 3).

- **Selection function** (τ_α) is a binary rule that outputs 1 for high-quality amortized explanations and 0 for low-quality ones based on the uncertainty metric. We define τ_α to ensure a fraction α of inputs receive amortized explanations. Smaller α implies higher-quality selective explanations but also more computations (Section 3).
- **Combination function** (λ_h) optimally linearly combines amortized and Monte Carlo explanations to minimize MSE from high-quality explanations (Theorem 1). We propose explanations with initial guess and fit λ_h to optimize their quality (Section 4).

Algorithm 1 describes the procedure to compute the uncertainty metric, selection function, and combination function using the results we describe in Section 3 and 4. Although selective explanations can be applied to any feature attribution method, we focus on Shapley values since they are widely used and most amortized explainers are tailored for them [12–14]. We discuss how selective explanations can be applied to LIME and provide more details on feature attribution methods in Appendix B. Next, we describe specific feature attribution methods that we use as building blocks for selective explainers of the form (1).

Shapley Values (SHAP) [4] is a **high-quality** explainer that attributes a value ϕ_i for each feature x_i in $\mathbf{x} = (x_1, \dots, x_d)$ which is the marginal contribution of feature x_i if the model was to predict \mathbf{y}

$$\phi_i(\mathbf{x}, \mathbf{y}) = \frac{1}{d} \sum_{S \subset [d] \setminus \{i\}} \binom{d-1}{|S|}^{-1} (h_{\mathbf{y}}(\mathbf{x}_{S \cup \{i\}}) - h_{\mathbf{y}}(\mathbf{x}_S)). \quad (2)$$

SHAP has several desirable properties and is widely used. However, as (2) indicates, computing Shapley values and the attribution vector $\text{HQ}(\mathbf{x}, \mathbf{y}) = (\phi_1(\mathbf{x}, \mathbf{y}), \dots, \phi_d(\mathbf{x}, \mathbf{y}))$ requires 2^d inferences from h , making SHAP impractical for large models where inference is costly. This has motivated several approximation methods for SHAP, discussed next.

Shapley Value Sampling (SVS) [11] is a **Monte Carlo** explainer that approximates SHAP by restricting the sum in (2) to m uniformly sampled permutations of features performing $n = md + 1$ inferences. We denote SVS that samples m feature permutations by SVS- m .

Kernel Shap (KS) [4] is a **Monte Carlo** explainer that approximate Shapley values using the fact that SHAP can be computed by solving a weighted linear regression problem using n input perturbations resulting in n inferences. We refer to Kernel Shap using n inferences as KS- n .

Stochastic Amortization [13] is a **Amortized** explainer that uses noisy Monte Carlo explanations to learn high-quality explanations. Covert et al. [13] trained an amortized explainer in a model class \mathcal{F} (multilayer perceptrons) $\text{Amor} \in \mathcal{F}$ to take (\mathbf{x}, \mathbf{y}) and predicts an explanation $\text{Amor}(\mathbf{x}, \mathbf{y}) \approx \text{HQ}(\mathbf{x}, \mathbf{y})$ by minimize the L_2 norm from Monte Carlo explanations $\text{MC}^n(\mathbf{x}, \mathbf{y})$.

Algorithm 1 Building a Selective Explainer

Require: Datasets: $\mathcal{D}_{\text{train}}, \mathcal{D}_{\text{cal}}$. Explainers: $\text{Amor}, \text{MC}^n, \text{MC}^{n'}$. Coverage: α .

Ensure: Selection function: τ_α . combination function : λ_h .

- 1: Fit the uncertainty metric s_h using $\mathcal{D}_{\text{train}}, \text{Amor}$, and MC^n (using (4) or (5))
 - 2: Compute t_α using \mathcal{D}_{cal} (7)
 - 3: Define the selection function τ_α using s_h and t_α (6)
 - 4: Define bins $Q_i = [t_{\alpha_i}, t_{\alpha_{i+1}})$ for partition $\alpha_i = \frac{i-1}{k}$ for $i \in [k+1]$ (9)
 - 5: For $i \in [k+1]$ Compute λ_i as in (12) using $\mathcal{D}_{\text{cal}}, \text{Amor}, \text{MC}^n$, and $\text{MC}^{n'}$.
 - 6: Define $\lambda_h(\mathbf{x}) = \sum_{i=1}^{k+1} \lambda_i \mathbf{1}[s_h(\mathbf{x}) \in Q_i]$ as in (9)
 - 7: **return** $\tau_\alpha, \lambda_h(\mathbf{x})$
-

Specifically, the amortized explainer is given by

$$\text{Amor} \in \underset{f \in \mathcal{F}}{\operatorname{argmin}} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{train}}} \|f(\mathbf{x}, \mathbf{y}) - \text{MC}^n(\mathbf{x}, \mathbf{y})\|_2^2. \quad (3)$$

Amortized Shap for LLMs [14] is a **Amortized** explainer similar to stochastic amortization but tailored for LLMs. Yang et al. [14] train a linear regression on the LLM embeddings $[e_1(\mathbf{x}), \dots, e_{|\mathbf{x}|}(\mathbf{x})]$ to minimize the L_2 norm from Monte Carlo explanations $\text{MC}^n(\mathbf{x}, \mathbf{y})$ and define the amortized explainer as $\text{Amor}(\mathbf{x}, \mathbf{y}) = (W_{\mathbf{y}}e_1(\mathbf{x}) + b_{\mathbf{y}}, \dots, W_{\mathbf{y}}e_{|\mathbf{x}|}(\mathbf{x}) + b_{\mathbf{y}})$, $W_{\mathbf{y}}$ is a matrix and $b_{\mathbf{y}} \in \mathbb{R}$.

We use stochastic amortization to produce amortized explainers for tabular datasets and Amortized Shap for LLMs to produce explainers for LLM predictions. Both explainers are trained using SVS-12 as MC^n . High-quality and Monte Carlo explanations are computed using the Captum library [23].

3 Selecting Explanations

In this section, we define key concepts for selective explainers: (i) uncertainty metrics s_h to quantify the likelihood of an explanation being low-quality and (ii) selection functions (τ_α) to predict when amortized explanations are high-quality based on the value of an uncertainty metric.

Uncertainty Metrics for High-Dimensional Regression: An uncertainty metric is a function tailored for the model h that takes \mathbf{x} and outputs a real number $s_h(\mathbf{x})$ that encodes information about the uncertainty of the model h in the prediction for \mathbf{x} . Generally, if $s_h(\mathbf{x}) < s_h(\mathbf{x}')$ then the model is more confident about the prediction $h(\mathbf{x})$ than $h(\mathbf{x}')$ [18, 19]. Existing uncertainty metrics cater to (i) classification [18–22] and (ii) one-dimensional regression [22, 24–26], but none specifically address high-dimensional regression – which is our case of interest (d -dimensional explanations). Next, we propose two uncertainty metrics tailored to high-dimensional outputs: (i) Deep uncertainty and (ii) Learned uncertainty.

Deep Uncertainty is inspired by deep ensembles [27], a method that uses an ensemble of models to provide confidence intervals for the predictions of one model. We run the training pipeline for the amortized explainer described in (3) k times, each with a different random seed, resulting in k different amortized explainers $\text{Amor}^1, \dots, \text{Amor}^k$. We define the deep uncertainty as

$$s_h^{\text{Deep}}(\mathbf{x}) \triangleq \frac{1}{dk} \sum_{i=1}^d \text{Var}(\text{Amor}^1(\mathbf{x})_i, \dots, \text{Amor}^k(\mathbf{x})_i). \quad (4)$$

Here, $\text{Var}(a_1, \dots, a_k)$ is the variance of the sample $\{a_1, \dots, a_k\}$ and $\text{Amor}^j(\mathbf{x})_i$ indicates the i -th entry of the feature attribution vector $\text{Amor}^j(\mathbf{x})$. Hence, deep uncertainty is the average (across entries) of the variance (across all trained amortized explainers) for the predicted attributions.

If the deep uncertainty for a point \mathbf{x} is zero, then the amortized explainers produce the same feature attribution. On the other hand, if the deep uncertainty is high, then the feature attributions vary widely across the amortized explainers. Intuitively, the points with a higher deep uncertainty are more affected by a random seed change, implying more uncertainty in the explanation.

Learned Uncertainty uses data to predict the amortized explainer uncertainty at an input point \mathbf{x} . We choose ℓ (the loss function) between two explanations to be MSE. The learned uncertainty metric is a function in the class \mathcal{F} (multilayer perceptron in our experiments) such that

$$s_h^{\text{Learn}} \in \underset{s \in \mathcal{F}}{\operatorname{argmin}} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{train}}} |s(\mathbf{x}) - \ell(\text{Amor}(\mathbf{x}; \mathbf{y}), \text{MC}^n(\mathbf{x}; \mathbf{y}))|^2. \quad (5)$$

Ideally, instead of using the Monte Carlo explanation MC^n as the reference in (5), we would like to use high-quality explanations, i.e., $\ell(\text{Amor}(\mathbf{x}; \mathbf{y}), \text{HQ}(\mathbf{x}; \mathbf{y}))$. However, these computationally expensive explanations are usually not available. Thus, we resort to using Monte Carlo explanations.

For large language models, the textual input \mathbf{x} is encoded in a sequence of token embedding $[e_1(\mathbf{x}), \dots, e_{|\mathbf{x}|}(\mathbf{x})]$ such that $e_i(\mathbf{x}) \in \mathbb{R}^d$ for $i \in [|\mathbf{x}|]$. In this case, we use the mean (i.e., “mean-pooling”) of the token embeddings to train the learned uncertainty metric instead of \mathbf{x} .

We analyze the performance of the proposed uncertainty metrics in Section 5, showing that it can be used to detect low-quality explanations from the amortized explainer. Our results indicate that these functions closely approximate the best possible uncertainty measure – the Oracle with knowledge of high-quality explanations (Figure 3). Next, we define the selection function that allows practitioners to set a coverage (percentage of points) α that will receive amortized explanations.

Selection functions: a selection function is the binary qualifier (τ_α) that thresholds the uncertainty metric by $t_\alpha \in \mathbb{R}$ given by

$$\tau_\alpha(\mathbf{x}) \triangleq \begin{cases} 1 & \text{if } s_h(\mathbf{x}) \leq t_\alpha \text{ (high-quality explanations)} \\ 0 & \text{if } s_h(\mathbf{x}) > t_\alpha \text{ (low-quality explanations)} \end{cases}. \quad (6)$$

Intuitively, t_α is the maximum uncertainty level tolerated by the user. In practice, if the output of the selection function is 1 (high-quality explanation), we use the explanations from the amortized model; if the output of the selection function is 0 (low-quality explanation), we use explanations with initial guess (see Definition 2 below) to improve the explanation provided to the user. The threshold t_α is chosen to be the α -quantile of the uncertainty metric to ensure that at least a fraction α of points receive a computationally cheap explanation – we call α the *coverage*. Specifically, given α , we calibrate t_α in the calibration dataset \mathcal{D}_{cal} and compute it as

$$t_\alpha \triangleq \min_{t \in \mathbb{R}} t, \text{ such that } \Pr_{\text{cal}}[s_h(\mathbf{x}) \leq t] \geq \alpha, \quad (7)$$

where \Pr_{cal} is the empirical distribution of the calibration dataset. For discussions on selecting coverage with guarantees on the number of inferences for selective explanations, see Appendix C.

Remark 1. A property of selective predictions [19], that is transferred to selective explanations is that it is possible to control the explainer’s performance via the threshold t_α with guaranteed performance but without providing predictions for all points. This result is displayed in Figure 3.

4 Explanations with Initial Guess

In the previous section, we introduced methods to detect points likely to receive low-quality explanations from amortized explainers. This raises the question: *How can we improve the explanations for these points?* One approach is to simply use Monte Carlo (MC) explanations instead of amortized explanations. However, this ignores potentially valuable information already computed by the amortized explainer. In this section, we propose a more effective solution called *explanations with initial guess*, which combines amortized and Monte Carlo explanations to improve explanation quality.

Explanation with Initial Guess uses an optimized linear combination of the amortized explanation with a more computationally expensive method – the Monte Carlo explainer – to improve the quality of the explanation. We formally define *explanations with initial guess* next.

Definition 2 (Explanation with Initial Guess). Given a Monte Carlo explainer $MC^n(\mathbf{x}, \mathbf{y})$, and a combination function $\lambda_h : \mathbb{R}^d \rightarrow \mathbb{R}$ that reflects the quality of the amortized explanation \mathbf{Amor} , we define the explanation with initial guess as

$$IG(\mathbf{x}, \mathbf{y}) \triangleq \lambda_h(\mathbf{x})\mathbf{Amor}(\mathbf{x}, \mathbf{y}) + (1 - \lambda_h(\mathbf{x}))MC^n(\mathbf{x}, \mathbf{y}). \quad (8)$$

Recall that when $\tau_\alpha(\mathbf{x}) = 0$, selective explanations use the explanation with initial guess (1) to improve low-quality amortized explanations, i.e., $SE(\mathbf{x}, \mathbf{y}) = IG(\mathbf{x}, \mathbf{y})$.

Defining explanations with initial guess as the linear combination between the amortized and the Monte Carlo explanations is inspired by the literature on shrinkage estimators [28, 29] that use an initial guess ($\mathbf{Amor}(\mathbf{x}, \mathbf{y})$ in our case) to improve the estimation MSE in comparison with only using the empirical average (a role played by $MC^n(\mathbf{x}, \mathbf{y})$ in our case). Next, we tune λ_h to minimize the MSE from high-quality explanations.

Optimizing the Explanation Quality: Our goal is for explanations with initial guess to approximate the high-quality explanations from HQ, i.e., $\|IG(\mathbf{x}, \mathbf{y}) - HQ(\mathbf{x}, \mathbf{y})\|$ to be minimized. To achieve this, we optimize the function λ_h as follows.

First, since high-quality explanations HQ are unavailable, we use another Monte Carlo explanation $MC^{n'}$ that closely approximates HQ. $MC^{n'}$ is different from MC^n and potentially more computationally expensive. Importantly, $MC^{n'}$ is only needed beforehand when computing λ_h , not at prediction time. In our experiments, we use SVS-12 for $MC^{n'}$.

Second, we quantize the range of the uncertainty metric s_h into bins to aggregate points with similar uncertainty and define the bins Q_i by a partition $0 = \alpha_1 < \alpha_2 < \dots < \alpha_m = 1$ of $[0, 1]$:

$$Q_i \triangleq [t_{\alpha_i}, t_{\alpha_{i+1}}), \quad \forall i \in [m - 1] \quad (9)$$

where t_{α_i} is defined as in (7). We then define the combination function to be

$$\lambda_h(\mathbf{x}) = \lambda_i \text{ if } s_h(\mathbf{x}) \in Q_i, \quad (10)$$

λ_h is chosen to optimize the explanation-quality for points with similar uncertainty, λ_i is given by:

$$\lambda_i \triangleq \operatorname{argmin}_{\lambda \in \mathbb{R}} \sum_{\substack{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{cal}} \\ s_h(\mathbf{x}) \in Q_i}} \left\| IG(\mathbf{x}, \mathbf{y}) - MC^{n'}(\mathbf{x}, \mathbf{y}) \right\|_2^2. \quad (11)$$

We only compute λ_i once per bin we provide explanations with initial guess (8), i.e., when $\tau_\alpha(\mathbf{x}) = 0$.

Theorem 1 provides a closed-form solution for λ_i .

Theorem 1 (Optimal λ_h). *Let $0 = \alpha_1 < \alpha_2 < \dots < \alpha_m = 1$ and define Q_i as in (9). Then the solution to the optimization problem in (11) is given by*

$$\lambda_i = \frac{\sum_{\substack{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{cal}} \\ s_h(\mathbf{x}) \in Q_i}} \langle MC^n(\mathbf{x}, \mathbf{y}) - MC^{n'}(\mathbf{x}, \mathbf{y}), MC^n(\mathbf{x}, \mathbf{y}) - \mathbf{Amor}(\mathbf{x}, \mathbf{y}) \rangle}{\sum_{\substack{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{cal}} \\ s_h(\mathbf{x}) \in Q_i}} \| \mathbf{Amor}(\mathbf{x}, \mathbf{y}) - MC^n(\mathbf{x}, \mathbf{y}) \|_2^2}. \quad (12)$$

The range of uncertainty functions is **quantized** for two main reasons. First, the uncertainty metric s_h encodes the amortized explainer’s uncertainty for each point \mathbf{x} . This uncertainty quantification should be reflected in the choice of λ_h . Quantizing the range of s_h allows us to group points

with similar uncertainty levels and optimize λ_h for each group separately. Second, quantizing the range of s_h enables us to have multiple point per bin Q_i allowing us to compute λ_i to minimize the MSE in each bin.

We use the **Monte Carlo** explainer $MC^{n'}$ because: (i) as mentioned above, we assume we don't have access to high-quality explanations due to their computational cost and (ii) even when using this Monte Carlo explainer, we show that in all bins λ_i approximates well the optimal combination function computed assuming access to high-quality explanations from HQ defined as

$$\lambda_i^{\text{opt}} = \underset{\lambda \in [0,1]}{\operatorname{argmin}} \sum_{\substack{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{cal}} \\ s_h(\mathbf{x}) \in Q_i}} \|\mathbb{G}(\mathbf{x}, \mathbf{y}) - \text{HQ}(\mathbf{x}, \mathbf{y})\|_2^2.$$

Specifically, Theorem 2 shows that $\lambda_i \approx \lambda_i^{\text{opt}}$, Appendix E shows the formal version of the Theorem along with the proofs for all results in this section.

Theorem 2 (Informal $\lambda_i \approx \lambda_i^{\text{opt}}$). *If (i) MC^n is sufficiently different from the amortized explainer Amor and (ii) $MC^{n'}$ approximates the high-quality explanations HQ then λ_i and λ_i^{opt} are close with high-probability for all bins Q_i , i.e.,*

$$|\lambda_i - \lambda_i^{\text{opt}}| \leq \epsilon \text{ with probability at least } 1 - e^{-C|Q_i|}.$$

for a $C > 0$ and $|Q_i|$ is the number of points in the validation dataset \mathcal{D}_{cal} that are in the bin Q_i .

5 Experimental Results

This section analyzes the performance of selective explanations. We show that (i) uncertainty metrics accurately identify low-quality explanations (Figure 3), (ii) explanations with initial guess have a higher quality than amortized and Monte Carlo explanations (Figure 4), (iii) selective explanations improve the performance of the lowest-quality explanations (Figure 5), and (iv) selective explanations improve local fidelity (Figure 6). We also (a) analyze how the quality of Monte Carlo explanations impact explanations with initial guess (Appendix D.4) and (b) show that selective explanations can be used to improve the inference vs. MSE trade-off of Monte Carlo explanations (Appendix D.5).

Experimental Setup: We generate selective explanations and evaluate their MSE and Spearman's correlation to the high-quality explanation computed using a large number of inferences¹. Although our results hold for any feature attribution method, in this section, we focus on Shapley values due to their frequent use and their prevalence in the literature on amortized explainers [12–14]. Seaborn [30] is used to compute 95% confidence intervals using bootstrap.

Datasets & Tasks: We show results for four datasets: two tabular datasets UCI-Adult [31] and UCI-News [32], and two text classification datasets Yelp Review [33] and Toxigen [34]. In the UCI-Adult dataset, the task is to predict if a given individual makes more than \$50k a year from a vector with 12 features; in UCI-News, the task is to predict if a news article will be shared more than 1400 (median sharing count) times from a vector with 58 features. In the Yelp Review dataset, the task is to predict whether a given Yelp review is positive or not, and in the Toxigen dataset, the task is to predict whether a given input text is toxic or not. We use 4000 samples from each dataset due to the computational cost of computing high-quality explanations for this evaluation. **Models:**

¹We provide details on how high-quality explanations were computed in Appendix D.1

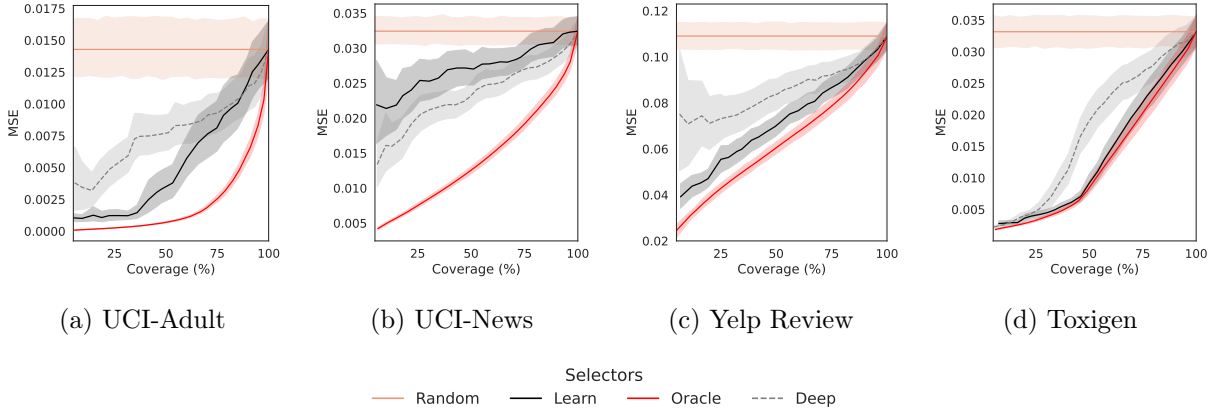


Fig. 3: Coverage (α) vs. test MSE from the high-quality explanation. The MSE is computed over the points such that $\tau_\alpha(\mathbf{x}) = 1$, i.e., predicted to be high-quality for a given coverage (x-axis). When coverage is 100%, the MSE is the average performance for the amortized explainer.

For the tabular datasets, we train a multilayer perceptron [35] to learn the desired task. We use the HuggingFace Bert-based model `textattack/bert-base-uncased-yelp-polarity` [36] for the Yelp dataset and the Roberta-based model `tomh/toxigen_roberta` [34] for the Toxigen dataset ².

Efficacy of Uncertainty Measures: In Figure 3, the x-axis shows the coverage (α) of the amortized explainer, while the y-axis shows the average mean square error (MSE) ³ of the selected amortized explanations from high-quality explanations, using deep uncertainty (with 20 models) and learned uncertainty to select which points should fall within the coverage. The Oracle⁴ is computed by sorting examples from smallest to highest MSE and computing the average MSE for the bottom α -fraction of points and is the best that can be done. Figure 3 shows that both deep uncertainty and learned uncertainty metrics can successfully identify examples that will receive lower and higher-quality explanations. For the large models ((c) and (d)), the learned uncertainty metric can identify points that will receive low-quality explanations almost as accurately as the Oracle. Also, we can ensure an MSE smaller than 0.003 (Adult), 0.025 (News), 0.07 (Yelp), and 0.007 (Toxigen) instead of the average MSEs that are 0.014, 0.032, 0.11, and 0.032 respectively with theoretical guarantees [19] for 50% of the points as described in Remark 1.

Explanations with Initial Guess: In Figure 4 we compare explanations with initial guess (Definition 2) to only using Monte Carlo explanations to provide improve for low-quality explanations, i.e., $\lambda_h = 0$ which we call Naive. When the MSE from the Monte Carlo is smaller than from the amortized explainer ((a) and (c)), employing explanations with initial guess results in a smaller MSE compared to naively using the Monte Carlo explainer. This suggests that despite their lower quality, the amortized explanations contain valuable information that can be used. When the Monte Carlo is larger than the amortized MSE ((b) and (d)), naive worsens the MSE while explanations with initial guess reduce the MSE, even when using poorer-quality Monte Carlo explanations. We used KS-32 for the Tabular datasets and SVS-12 for the textual datasets⁵.

²For more details on implementation, please see Appendix D.1.

³In Appendix D.2, we also show the effect of our uncertainty metrics on Spearman’s correlation.

⁴The oracle is computationally expensive because it requires access to high-quality explanations.

⁵In Appendix D.3, we also show that selective explanations improve the MSE while maintaining the same level of Spearman’s correlation as the Naive approach.

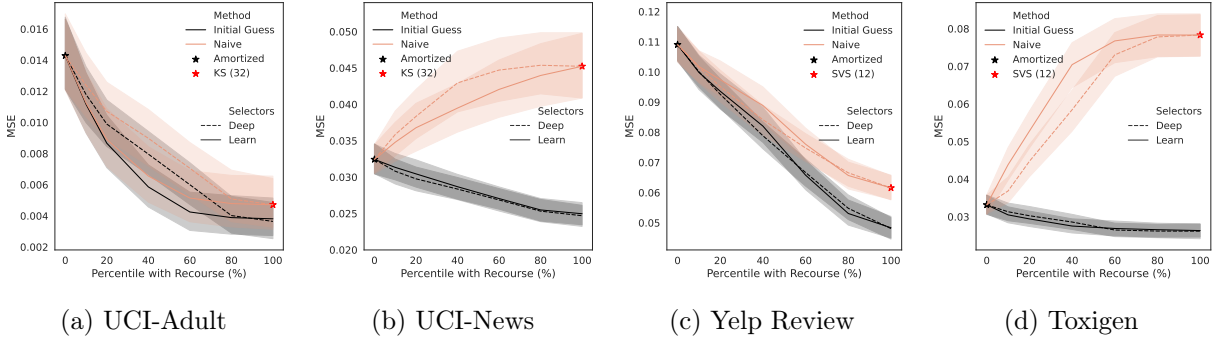


Fig. 4: Fraction $(1 - \alpha)$ of points that receive explanations with initial guess (x-axis) vs. MSE of selective explanations w.r.t. high-quality explanations (y-axis). Naive uses $\lambda_h = 0$ while Initial guess uses λ_h in (12). MSE is computed across all points in the test dataset.

Worst Case Performance Improvement: In Figure 5, we analyze the performance of selective explanations for varying coverages (both in terms of MSE and Spearman’s correlation) for the points that receive the worst-performing explanations. We observe that selective explanations, even with only 20% of points receiving explanations with initial guess, increase Spearman’s correlation and decrease MSE consistently across datasets. Remarkably, when providing explanations with initial guess for 20% of the population in the Yelp dataset (Figure 5 (c)), selective explanations result in Spearman’s correlation for the worst 4% of points that is better than that for the worst 10% from the amortized explainer – even clearer in the UCI-Adult dataset. We use SVS-3 for the tabular datasets and SVS-12 for the text datasets.

Perturbation Curve: Figure 6 shows that selective explanations increase the local fidelity of the amortized explainer and that the local fidelity increases with the percentage of points that receive explanations with initial guess (i.e., decreases with coverage). Both Yelp and Toxigen models are receiving recourse by using SVS-12. Notably, for Yelp (Figure 6 (a)), when providing explanations with initial guess for 60% of the points and using the amortized explainer the other 40% of the time, we achieve local fidelity that is close to the computationally expensive high-quality explanations for the 30% most important tokens.

6 Final Remarks

Conclusion: We propose *Selective explanations* that first identify which inputs would receive a low-quality but computationally cheap explanation (amortized) and then perform model inferences to improve the quality of these explanations. Specifically, we propose *explanations with initial guess* to improve the quality of explanations by combining computationally cheap explanations (amortized) with more expensive explanations (Monte Carlo) using an optimized combination function, improving the explanation performance beyond both explanations. We perform experiments in large language

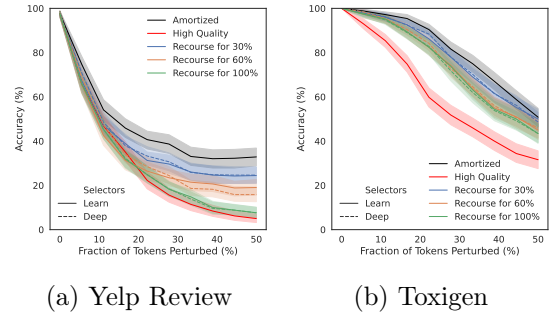


Fig. 6: Model accuracy (y-axis) when removing the tokens with the highest attribution scores according to the amortized explainer (black), selective explanations with recourse for 30% (blue), 60% (orange), and 100% (green) of points, and high-quality explanations (red).

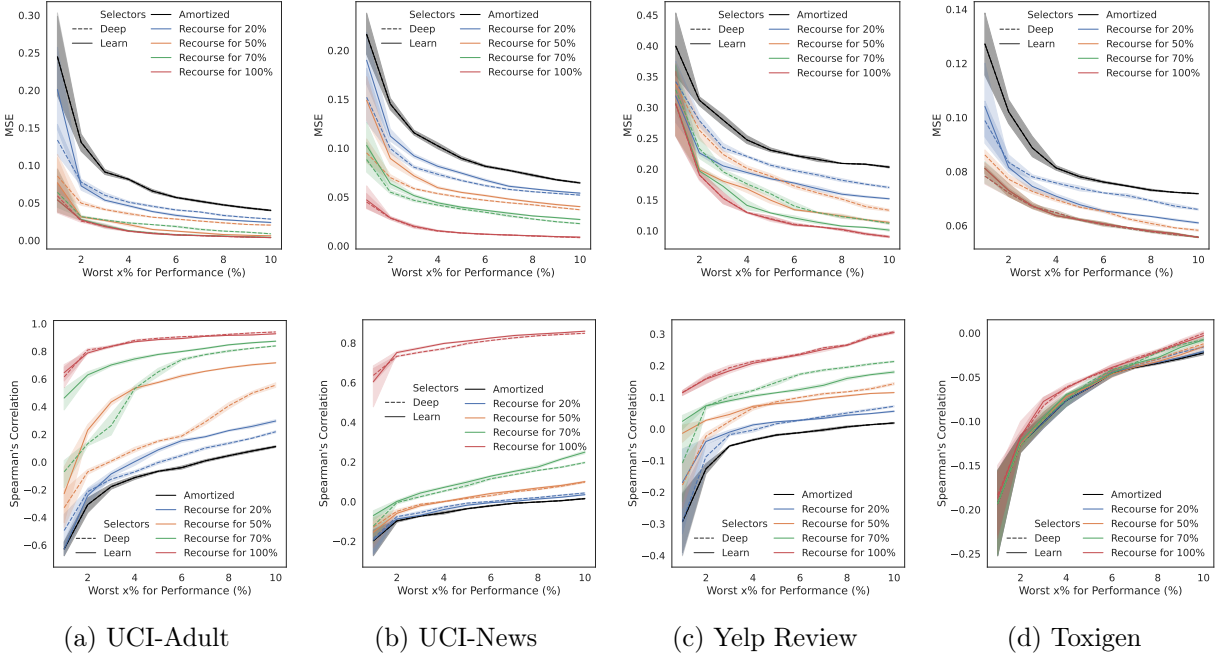


Fig. 5: MSE (top) and Spearman’s correlation (bottom) for explanations with the worst performance (highest MSE and smallest Spearman’s) in $\mathcal{D}_{\text{test}}$. Colors indicate different percentages $1 - \alpha$ of points receiving explanations with initial guess: 20% (blue), 50% (orange), 70% (green), 100% (red), and amortized explanations 0% (black). Performance is computed in each quantile.

models and tabular data classifiers empirically demonstrating the efficacy of selective explanations. Our experiments indicate that selective explanations (i) efficiently identify points that the amortized explainer would produce low-quality explanations, (ii) improve the quality of the worst-quality explanations, and (iii) improve the local fidelity of amortized explanations.

Limitations: Selective explanations can be applied to any feature attribution method for which amortized and Monte Carlo explainers were developed. However, our empirical results focus on Shapley values. We leave the application of selective explanations to other attribution methods for future work. Additionally, we focus on large language models (LLMs) used for text classification. Consequently, we do not explore image classifiers, which may also interest the interpretability community.

Acknowledgements

The authors thank Amit Dhurandhar for early discussions on the trustworthiness of amortized explainers. This material is based upon work supported by the National Science Foundation under grants CAREER 1845852, CIF 1900750, CIF 2312667, and FAI 2040880, and awards from Google Research and Amazon.

References

- [1] OpenAI. Using gpt-4 for content moderation. <https://openai.com/index/using-gpt-4-for-content-moderation>. Accessed: 2024-05-01.

- [2] Preetam Ghosh and Vaishali Sadaphal. Jobrecogpt – explainable job recommendations using llms, 2023.
- [3] Sheng Wang, Zihao Zhao, Xi Ouyang, Qian Wang, and Dinggang Shen. Chatcad: Interactive computer-aided diagnosis on medical image using large language models, 2023.
- [4] Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 4765–4774. Curran Associates, Inc., 2017. URL <http://papers.nips.cc/paper/7062-a-unified-approach-to-interpreting-model-predictions.pdf>.
- [5] 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, San Francisco, CA, USA, August 13-17, 2016*, pages 1135–1144, 2016.
- [6] Lucas Monteiro Paes, Dennis Wei, Hyo Jin Do, Hendrik Strobelt, Ronny Luss, Amit Dhurandhar, Manish Nagireddy, Karthikeyan Natesan Ramamurthy, Prasanna Sattigeri, Werner Geyer, and Soumya Ghosh. Multi-level explanations for generative language models, 2024.
- [7] Vivek Miglani, Aobo Yang, Aram H. Markosyan, Diego Garcia-Olano, and Narine Kokhlikyan. Using Captum to explain generative language models, 2023.
- [8] Jianbo Chen and Michael I. Jordan. Ls-tree: Model interpretation when the data are linguistic. *ArXiv*, abs/1902.04187, 2019. URL <https://api.semanticscholar.org/CorpusID:60441455>.
- [9] Matthew D. Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In David Fleet, Tomas Pajdla, Bernt Schiele, and Tinne Tuytelaars, editors, *Computer Vision – ECCV 2014*, pages 818–833, Cham, 2014. Springer International Publishing. ISBN 978-3-319-10590-1.
- [10] Ian Covert and Su-In Lee. Improving kernelshap: Practical shapley value estimation using linear regression. In Arindam Banerjee and Kenji Fukumizu, editors, *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*, volume 130 of *Proceedings of Machine Learning Research*, pages 3457–3465. PMLR, 13–15 Apr 2021. URL <https://proceedings.mlr.press/v130/covert21a.html>.
- [11] Rory Mitchell, Joshua Cooper, Eibe Frank, and Geoffrey Holmes. Sampling permutations for shapley value estimation. *Journal of Machine Learning Research*, 23(43):1–46, 2022. URL <http://jmlr.org/papers/v23/21-0439.html>.
- [12] 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*, 2022. URL https://openreview.net/forum?id=Zq2G_VTV53T.
- [13] Ian Covert, Chanwoo Kim, Su-In Lee, James Zou, and Tatsunori Hashimoto. Stochastic amortization: A unified approach to accelerate feature and data attribution, 2024.
- [14] Chenghao Yang, Fan Yin, He He, Kai-Wei Chang, Xiaofei Ma, and Bing Xiang. Efficient shapley values estimation by amortization for text classification. In *Annual Meeting of the Association for Computational Linguistics*, 2023. URL <https://api.semanticscholar.org/CorpusID:258987882>.

- [15] Patrick Schwab and Walter Karlen. Explain: Causal explanations for model interpretation under uncertainty. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper_files/paper/2019/file/3ab6be46e1d6b21d59a3c3a0b9d0f6ef-Paper.pdf.
- [16] Yu-Neng Chuang, Guanchu Wang, Fan Yang, Quan Zhou, Pushkar Tripathi, Xuanting Cai, and Xia Hu. Cortx: Contrastive framework for real-time explanation. *ArXiv*, abs/2303.02794, 2023. URL <https://api.semanticscholar.org/CorpusID:257365448>.
- [17] Robert Schwarzenberg, Nils Feldhus, and Sebastian Möller. Efficient explanations from empirical explainers. In Jasmijn Bastings, Yonatan Belinkov, Emmanuel Dupoux, Mario Giulianelli, Dieuwke Hupkes, Yuval Pinter, and Hassan Sajjad, editors, *Proceedings of the Fourth BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*, pages 240–249, Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.blackboxnlp-1.17. URL <https://aclanthology.org/2021.blackboxnlp-1.17>.
- [18] Stephan Rabanser, Anvith Thudi, Kimia Hamidieh, Adam Dziedzic, and Nicolas Papernot. Selective classification via neural network training dynamics. *ArXiv*, abs/2205.13532, 2022. URL <https://api.semanticscholar.org/CorpusID:249097456>.
- [19] Yonatan Geifman and Ran El-Yaniv. Selective classification for deep neural networks. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/file/4a8423d5e91fda00bb7e46540e2b0cf1-Paper.pdf.
- [20] Ran El-Yaniv et al. On the foundations of noise-free selective classification. *Journal of Machine Learning Research*, 11(5), 2010.
- [21] Aditya Gangrade, Anil Kag, and Venkatesh Saligrama. Selective classification via one-sided prediction. In *International Conference on Artificial Intelligence and Statistics*, pages 2179–2187. PMLR, 2021.
- [22] Yonatan Geifman and Ran El-Yaniv. Selectivenet: A deep neural network with an integrated reject option. In *International conference on machine learning*, pages 2151–2159. PMLR, 2019.
- [23] Narine Kokhlikyan, Vivek Miglani, Miguel Martin, Edward Wang, Bilal Alsallakh, Jonathan Reynolds, Alexander Melnikov, Natalia Kliushkina, Carlos Araya, Siqi Yan, and Orion Reblitz-Richardson. Captum: A unified and generic model interpretability library for pytorch, 2020.
- [24] Ahmed Zaoui, Christophe Denis, and Mohamed Hebiri. Regression with reject option and application to knn. *Advances in Neural Information Processing Systems*, 33:20073–20082, 2020.
- [25] Abhin Shah, Yuheng Bu, Joshua K Lee, Subhro Das, Rameswar Panda, Prasanna Sattigeri, and Gregory W Wornell. Selective regression under fairness criteria. In *International Conference on Machine Learning*, pages 19598–19615. PMLR, 2022.
- [26] Wenming Jiang, Ying Zhao, and Zehan Wang. Risk-controlled selective prediction for regression deep neural network models. In *2020 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2020.

- [27] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. *Advances in neural information processing systems*, 30, 2017.
- [28] HH Lemmer. From ordinary to bayesian shrinkage estimators. *South African Statistical Journal*, 15(1):57–72, 1981.
- [29] HH Lemmer. Note on shrinkage estimators for the binomial distribution. *Communications in statistics-theory and methods*, 10(10):1017–1027, 1981.
- [30] Michael L. Waskom. seaborn: statistical data visualization. *Journal of Open Source Software*, 6(60):3021, 2021. doi: 10.21105/joss.03021. URL <https://doi.org/10.21105/joss.03021>.
- [31] Barry Becker and Ronny Kohavi. Adult. UCI Machine Learning Repository, 1996. DOI: <https://doi.org/10.24432/C5XW20>.
- [32] Kelwin Fernandes, Pedro Vinagre, Paulo Cortez, and Pedro Sernadela. Online News Popularity. UCI Machine Learning Repository, 2015. DOI: <https://doi.org/10.24432/C5NS3V>.
- [33] Xiang Zhang, Junbo Zhao, and Yann LeCun. Character-level convolutional networks for text classification. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc., 2015. URL https://proceedings.neurips.cc/paper_files/paper/2015/file/250cf8b51c773f3f8dc8b4be867a9a02-Paper.pdf.
- [34] Thomas Hartvigsen, Saadia Gabriel, Hamid Palangi, Maarten Sap, Dipankar Ray, and Ece Kamar. ToxiGen: A large-scale machine-generated dataset for adversarial and implicit hate speech detection. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio, editors, *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3309–3326, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.234. URL <https://aclanthology.org/2022.acl-long.234>.
- [35] Simon Haykin. *Neural networks: a comprehensive foundation*. Prentice Hall PTR, 1994.
- [36] John Morris, Eli Liffand, Jin Yong Yoo, Jake Grigsby, Di Jin, and Yanjun Qi. Textattack: A framework for adversarial attacks, data augmentation, and adversarial training in nlp. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 119–126, 2020.
- [37] Ian Covert and Su-In Lee. Improving kernelshap: Practical shapley value estimation via linear regression. In *International Conference on Artificial Intelligence and Statistics*, 2020. URL <https://api.semanticscholar.org/CorpusID:227253750>.

A Overview

In this supplementary material we provide the following information:

- Appendix B discuss other high-quality and Monte Carlo explainers.
- Appendix C discuss a guide to select the coverage α when the agent providing selective explanations has a budget for the average number of inferences to provide an explanation.
- Appendix D shows more experimental results on selective explanations.
- Appendix E shows the proofs for the theoretical results in Section 4.

B Additional Explanation Methods

In this section, we describe high-quality, Monte Carlo, and amortized explainers with further details.

B.1 High-Quality Explainers

Shapley Values (SHAP) [4] is a **high-quality** explainer that attributes a value ϕ_i for each feature x_i in $\mathbf{x} = (x_1, \dots, x_d)$ which is the marginal contribution of feature x_i if the model was to predict \mathbf{y} (2).

$$\phi_i(\mathbf{x}, \mathbf{y}) = \frac{1}{d} \sum_{S \subset [d] / \{i\}} \binom{d-1}{|S|}^{-1} (h_{\mathbf{y}}(\mathbf{x}_{S \cup \{i\}}) - h_{\mathbf{y}}(\mathbf{x}_S)). \quad (13)$$

SHAP has several desirable properties and is widely used. However, as (2) indicates, computing Shapley values and the attribution vector $\text{HQ}(\mathbf{x}, \mathbf{y}) = (\phi_1(\mathbf{x}, \mathbf{y}), \dots, \phi_d(\mathbf{x}, \mathbf{y}))$ requires 2^d inferences from h , making SHAP impractical for large models where inference is costly. This has motivated several approximation methods for SHAP, discussed next⁶.

Local Interpretable Explanations (Lime). Lime is another feature attribution method [5] widely used to provide feature attributions. It relies on selecting combinations of features, removing these features from the input to generate perturbations, and using these perturbations to approximate the black box model h locally by a linear model. The coefficients of the linear model are considered to be the attribution of each feature. Formally, given a weighting kernel $\pi(S)$ and a penalty function Ω , the attribution produced by lime are given by

$$(\phi, a) = \underset{\phi \in \mathbb{R}^d, a \in \mathbb{R}}{\operatorname{argmin}} \sum_{S \subset [d]} \pi(S) \left(h(\mathbf{x}_S) - a_0 - \sum_{i \in S} \phi_i \right), \quad (14)$$

where $\text{HQ}(\mathbf{x}, \mathbf{y}) = \phi$. As in SHAP, to compute the feature attributions using lime, we need to perform a large number of model inferences, which is prohibitive for large models.

⁶We also discuss Lime and its amortized version in Appendix B

B.2 Monte Carlo Lime

Shapley Value Sampling (SVS) [11] is a **Monte Carlo** explainer that approximates SHAP by restricting the sum in (2) to specific permutations of feature. SVS computes the attribution scores by uniformly sampling m features permutations S_1, \dots, S_m restricting the sum in (2) and performing $n = md + 1$ inferences. We denote SVS that samples m feature permutations by SVS- m .

Kernel Shap (KS) [4] is a **Monte Carlo** explainer that approximate the Shapley values using the fact that SHAP can be computed by solving the optimization problem

$$(\phi, a) = \underset{\phi \in \mathbb{R}^d, a \in \mathbb{R}}{\operatorname{argmin}} \sum_{i=1}^n \pi(S_i) \left(h(\mathbf{x}_{S_i}) - a_0 - \sum_{j \in S_i} \phi_j \right), \quad (15)$$

using $\pi(S) = \binom{d}{|S|} |S| (d - |S|)$ and where $\operatorname{MC}^n(\mathbf{x}, \mathbf{y}) = \phi$. Kernel Shap samples $n > 0$ feature combinations S_1, \dots, S_n and define the feature attributions to be given by the coefficients ϕ . We refer to Kernel Shap using n inferences as KS- n . We use the KS- n from the Captum library [23] for our experiments.

Sample Constrained Lime. To approximate the attributions from Lime, we consider the sample-contained version of (15). Instead of sampling all feature combinations in $[d]$, we only uniformly sample a fixed number n of feature combinations S_1, \dots, S_n . For our experiments, shown in the appendix, we use the Sample Constrained Lime from the Captum library [23].

B.3 Amortized Explainers

Stochastic Amortization [13] is a **Amortized** explainer that uses noisy Monte Carlo explanations to learn high-quality explanations. Covert et al. [13] trained an amortized explainer $\operatorname{Amor} \in \mathcal{F}$ in a hypothesis class \mathcal{F} (we use multilayer perceptrons) that takes an input and predicts an explanation. Specifically, taking the amortized explainer to be the solution of the training problem given in (3).

$$\operatorname{Amor} \in \underset{f \in \mathcal{F}}{\operatorname{argmin}} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{train}}} \|f(\mathbf{x}, \mathbf{y}) - \operatorname{MC}^n(\mathbf{x}, \mathbf{y})\|_2^2. \quad (16)$$

We are interested in explaining the predictions of large models for text classification. However, the approach in (3) is only suitable for numerical inputs. Hence, we follow the approach from Yang et al. [14] to explain the predictions of large language models, explained next.

Amortized Shap for LLMs [14] is a **Amortized** explainer similar to the one in (3) but tailored for LLMs. First, the authors note that they can use the LLM to write all input texts \mathbf{x} as a sequence of token embedding $[e_1(\mathbf{x}), \dots, e_{|\mathbf{x}|}(\mathbf{x})]$ where $e_i(\mathbf{x}) \in \mathbb{R}^d$ denotes the LLM embedding for the i -th token contained in the input text \mathbf{x} and $|\mathbf{x}|$ is the number of tokens in the input text. Second, they restrict \mathcal{F} in (3) to be the set of all linear regressions that take the token embeddings and output the token attribution score. Then, they solve the optimization problem in

$$W \in \underset{W \in \mathbb{R}^d, b \in \mathbb{R}}{\operatorname{argmin}} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{train}}} \sum_{j=1}^{|\mathbf{x}|} \|W^T e_j(\mathbf{x}) + b - \operatorname{MC}^n(\mathbf{x}, \mathbf{y})_j\|_2^2, \quad (17)$$

and define the amortized explainer as $\operatorname{Amor}(\mathbf{x}) = (W^T e_1(\mathbf{x}) + b, \dots, W^T e_{|\mathbf{x}|}(\mathbf{x}) + b)$.

We use stochastic amortization to produce amortized explainers for tabular datasets and Amortized Shap for LLMs to produce explainers for LLM predictions. Both explainers are trained using SVS-12 as MC^n .

C Selecting Coverage for a Given Inference Budget

Determining Coverage from Inference Budget: Providing explanations with initial guess increases the number of model inferences from 1 when using solely the amortized explainer to $n + 1$. However, a practitioner may have a budget of inferences, i.e., a maximum average number of inferences they are willing to perform to provide an explanation. We formalize the notion of inference budget in Definition 3.

Definition 3 (Inference Budget). Denote by $N(\text{SE}(\mathbf{x}, \mathbf{y}))$ the number of model inferences to produce the explanation $\text{SE}(\mathbf{x}, \mathbf{y})$. The inference budget $N_{\text{budget}} \in \mathbb{N}$ is the maximum average number of inferences a practitioner is willing to perform per explanation, i.e., it is such that

$$N_{\text{budget}} \geq \mathbb{E}[N(\text{SE}(\mathbf{x}, \mathbf{y}))]. \quad (18)$$

Once an inference budget N_{budget} is defined, the coverage α should be set to follow it. In Proposition 1, we show the minimum coverage for the selective explanations to follow the inference budget.

Proposition 1 (Coverage for Inference Budget). *Let $N_{\text{budget}} \geq 1$ be the inference budget, and assume that the Monte Carlo method $\text{MC}^n(\mathbf{x}, \mathbf{y})$ uses n model inferences. Then, the coverage level α should be chosen such that*

$$\frac{n + 1 - N_{\text{budget}}}{n} = \min_{\alpha \in [0,1]} \alpha, \text{ such that } \mathbb{E}[N(\text{SE}(\mathbf{x}, \mathbf{y}))] \leq N_{\text{budget}}. \quad (19)$$

Recall that SVS- m performs $n = 1 + dm$ inferences ($\mathbf{x} \in \mathbb{R}^d$), and KS- m performs $n = m$ inferences.

D More Experimental Results

In this section, we (i) give further implementation details and (ii) discuss further empirical results.

D.1 More Details on Experimental Setup

High-Quality Explanations: We define the high-quality explanations for the tabular datasets to be given by Kernel Shap with as many inferences as needed for convergence, using the Shapley Regression library [37]. For the textual dataset, following [14], we define the high-quality explanations to be given by Kernel Shap using 8912 model inferences per explanation.

Amortized Explainers: For the tabular datasets, we use the amortized explainer from [13] that we describe in Section 2. Specifically, we use a multilayer perceptron model architecture to learn the shapley values for the tabular datasets. For the textual datasets, we use the linear regression on token-level textual embeddings to learn the shapley values, as described in Section 2. Both amortized models learn from the training dataset of explanations generated using Shapley Value Sampling from the Captum library [23] with parameter 12, i.e., SVS-12.

Uncertainty Metrics: We test the two proposed uncertainty metrics in Section 3, namely, deep uncertainty and uncertainty learn. For **deep uncertainty**, we run the training pipeline for the amortized explainers 20 times for each dataset we perform experiments on, resulting in 20 different amortized explainer that we use to compute (4). For **uncertainty learn**, we use the multilayer perceptron as the hypothesis class with only one hidden layer. The hidden layer was composed of $\kappa = 3d$ neurons where d is the dimension of the input vector $\mathbf{x} \in \mathbb{R}^d$. The uncertainty learn metric was trained on $\mathcal{D}_{\text{train}}$, the same training dataset as the amortized explainers.

Dataset sizes: We use 4000 samples from each dataset due to computational limitations on the computation of high-quality explanations used to evaluate selective explanations. All explanations were computed using the Captum library [23]. The dataset \mathcal{D} with $N = 4000$ samples was partitioned in three parts, $\mathcal{D}_{\text{train}}$ with 50% of points, \mathcal{D}_{cal} with 25% of points, and $\mathcal{D}_{\text{test}}$ with the other 25% of points.

Computational Resources: All experiments were run in a A100 40 GB GPU. For each dataset, we compute different Monte Carlo explanations. For the UCI-News dataset, the high quality explanations took 4:30 hours to be generate until convergence while for UCI-Adult it took 3:46 hours. For the tabular datasets, all other Monte Carlo explainers were generated in less than 1 hour. For the language models, the high-quality explanations with 8192 model inferences, took 18:51 hours for the Toxigen dataset and 20:00 hours for the Yelp Review datasets. The other used Monte Carlo explanations took proportional (to the number of inferences) time to be generated.

D.2 Uncertainty Measures Impact on Spearman’s Correlation

Figure 7 shows in the x-axis the coverage (α) and in the y-axis the average Spearman’s correlation of the selected amortized explanations from high-quality explanations using deep uncertainty (with 20 models) and the uncertainty learn to select low-quality explanations. The Oracle⁷ is computed by sorting examples by the smallest to higher MSE and computing the average Spearman’s correlation in the bottom x-axis points accordingly to the MSE and is the best that can be done in terms of MSE.

Figure 7 shows that the Oracle and proposed uncertainty metrics don’t always select the points with the smallest Spearman’s correlation first. This implies that MSE and Spearman’s correlation don’t always align, i.e., there are points with high MSE and high Spearman’s correlation at the same time. However, we note that the uncertainty learns selector can be applied to **any** metric ℓ as we define in (5) including Spearman’s correlation and any combination of Spearman’s correlation and MSE aiming to approximate both metrics. Moreover, when the smallest MSE aligns with the highest Spearman’s correlation, i.e., the oracle is decreasing in Spearman’s correlation when the coverage increases (Figure 7 (a) and (c)), the proposed uncertainty metrics also accurately detect the low-quality explanations in term of Spearman’s correlation.

D.3 The Effect of Explanations with Initial Guess

In Figure 8 we compare explanations with initial guess (Definition 2) to only using the Monte Carlo to provide recourse to the low-quality explanaitons, i.e., $\lambda_h = 0$ we call it Naive. In all tested cases, Spearman’s correlation of the Monte Carlo method is comparable to or larger than the amortized explainer. Although selective explanations optimized for MSE by using explanations with initial guess (Definition 2), we observe that the Spearman’s correlation of selective explanations is close to or larger than the naive method, once again, demonstrating the efficacy of selective explanations.

D.4 Performance for Different Monte-Carlo Explainers

Figure 9 shows how the MSE and Spearman’s correlation behave accordingly with the quality of the Monte Carlo explainer. We compare Kernel Shap and Shapley Value Sampling in all experiments. We observe that when the quality of the Monte Carlo explainer increases, the quality of the Selective explanation also increases, i.e., the MSE decreases and the Spearman’s correlation increases.

⁷The oracle is computationally expensive because it requires access to high-quality explanations.

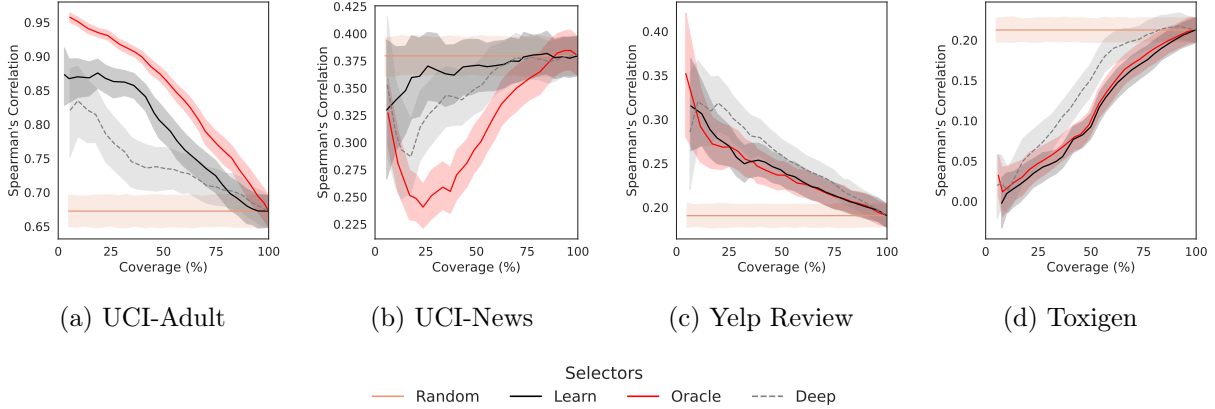


Fig. 7: Coverage vs. Spearman’s correlation from the high-quality explanation. Coverage is the percentage of the points that the selection function predicts that will receive a higher-quality explanation, i.e., $\tau_t(\mathbf{x}) = 1$. When coverage is 100% Spearman’s correlation is the average performance for the amortized explainer.

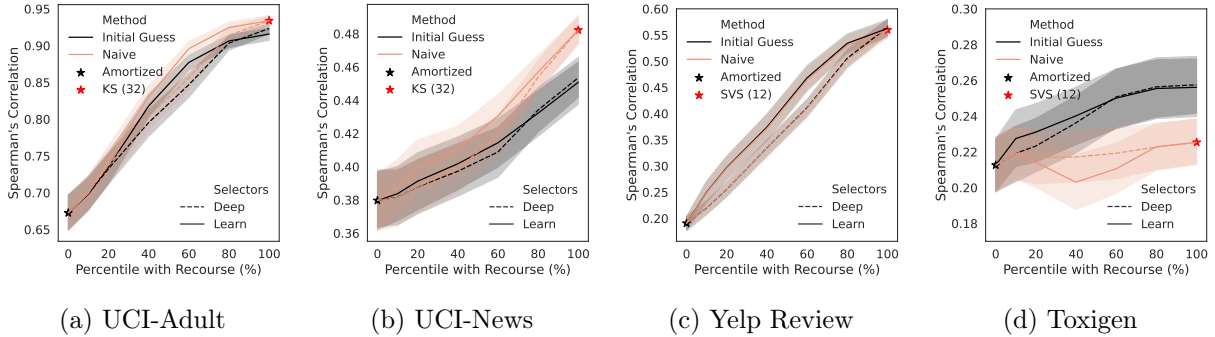


Fig. 8: Fraction of the population that receive explanations with initial guess (x-axis) vs. their Spearman’s correlation from the high-quality explanations (y-axis). Naive uses $\lambda_h = 0$ while initial guess uses explanations with initial guess, i.e., when λ_h is given in (12).

Moreover, we also observe diminishing returns, i.e., after a certain point, increasing the quality of the Monte Carlo explanations doesn’t lead to a tailored increase in performance. For example, observe the SVS method in the tabular datasets Figure 9 (a) and (b). We also observe that providing explanations with initial guess has a high impact on both Spearman’s correlation and MSE when only providing recourse to a small fraction of the population. For example, when providing explanations with initial guess for 20% of the population using SVS-12 in the Yelp Review dataset, Figure 9 (c), increases the Spearman’s correlation in more than 50% (from 0.2 to more than 0.3).

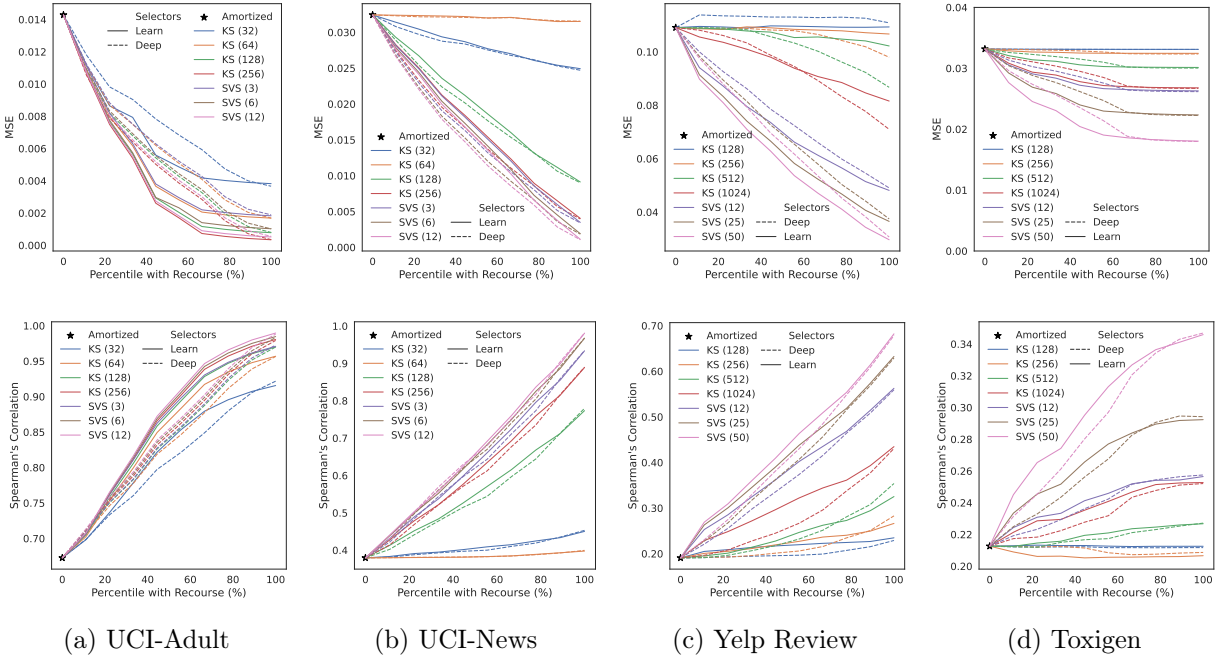


Fig. 9: MSE (top) and Spearman's correlation (bottom) for selective explanations using different Monte Carlo explainers.

D.5 Time Sharing Using Selective Explanations

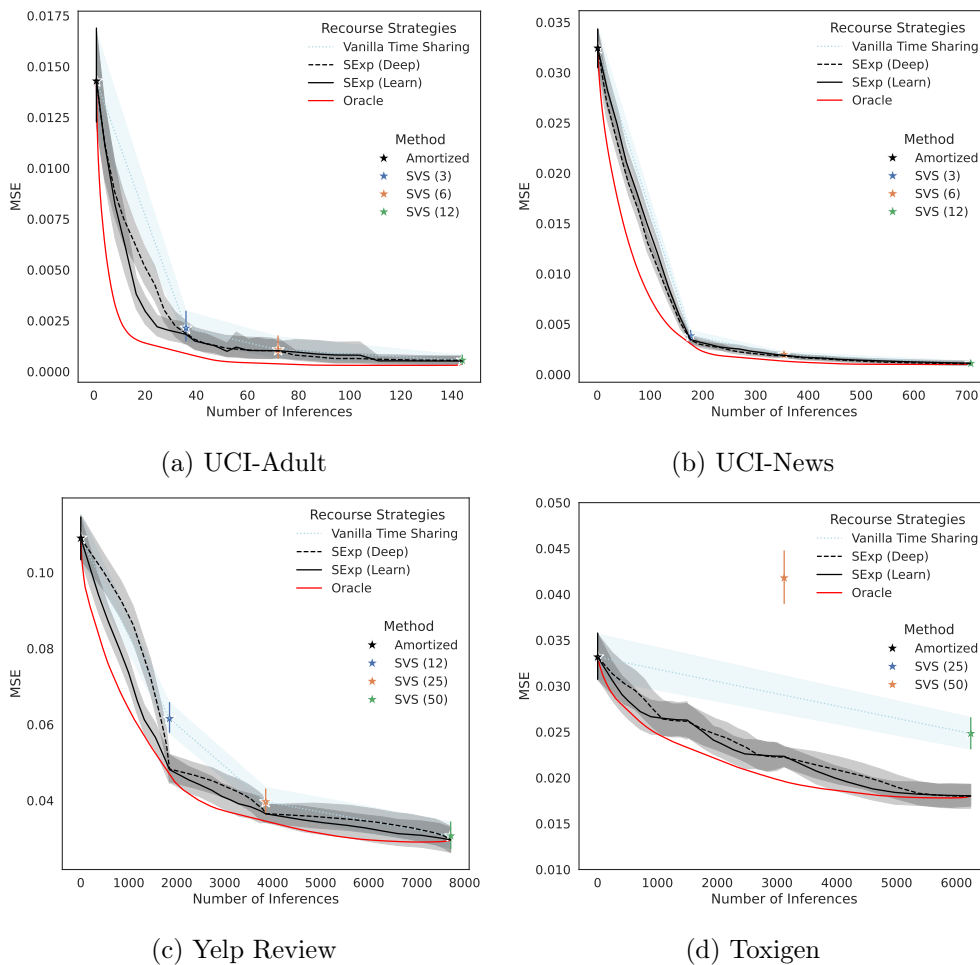


Fig. 10: Number of model inferences (x-axis) vs. MSE (y-axis) using (i) vanilla time sharing, (ii) time sharing using selective explanations compared to (iii) the oracle when the MSE of the provided explanation is known.

We analyze how selective explanations can be used to improve the quality of Monte Carlo methods by time sharing between methods. When computing explanations using Monte Carlo methods, we perform n model inferences (x-axis in Figure 10) until a desired MSE (y-axis in Figure 10) is achieved. This is done by gradually increasing the number of inferences per points we generate explanations – this is displayed by the blue dotted curve in Figure 10 and we name it vanilla time sharing because the inferences (time) are shared gradually across points. We also compare it with the Oracle given the red curve in Figure 10 where, for each point, we compute Monte Carlo explanations using SVS with parameter 12, 25, and 50, compute their MSE to high-quality explanations and give the best explanation possible for a given number of inferences. Oracle is the best that can be done in terms of MSE vs. Number of Inferences only using Monte Carlo explanations. We compare both Oracle and vanilla time sharing with time sharing using selective explanations given by the black lines in Figure 10. For the time sharing using selective explanations, we also gradually increase the number of inferences but use selective explanations instead of plain Monte Carlo explanations.

Figure 10 shows that selective explanations closely approximate the Oracle curve, indicating the

selective explanations have close to optimal trade-off between the number of model inferences and MSE. We highlight the performance of selective explanations in the Toxigen dataset. With only 1000 model inferences, we get better performance than using SVS-50 with about 6000 model inferences. We also note that in both LLMs, using selective explanations closely approximates the oracle and provides a better explanation with the same number of inferences than just using SVS.

E Proofs of Theoretical Results

Theorem 1 (Optimal λ_h). *Let $0 = \alpha_1 < \alpha_2 < \dots < \alpha_m = 1$ and define Q_i as in (9). Then, λ_i that solves the optimization problem in (11) is given by*

$$\lambda_i = \frac{\sum_{\substack{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{val}} \\ s_h(\mathbf{x}) \in Q_i}} \langle MC^n(\mathbf{x}, \mathbf{y}) - MC^{n'}(\mathbf{x}, \mathbf{y}), MC^n(\mathbf{x}, \mathbf{y}) - \text{Amor}(\mathbf{x}, \mathbf{y}) \rangle}{\sum_{\substack{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{val}} \\ s_h(\mathbf{x}) \in Q_i}} \| \text{Amor}(\mathbf{x}, \mathbf{y}) - MC^n(\mathbf{x}, \mathbf{y}) \|_2^2}. \quad (20)$$

Proof. First, recall that

$$\lambda_i \triangleq \operatorname{argmin}_{\lambda \in \mathbb{R}} \sum_{\substack{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{val}} \\ s_h(\mathbf{x}) \in Q_i}} \left\| \text{SE}(\mathbf{x}, \mathbf{y}) - MC^{n'}(\mathbf{x}, \mathbf{y}) \right\|_2^2 \quad (21)$$

$$= \operatorname{argmin}_{\lambda \in \mathbb{R}} \sum_{\substack{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{val}} \\ s_h(\mathbf{x}) \in Q_i}} \left\| \lambda \text{Amor}(\mathbf{x}, \mathbf{y}) + (1 - \lambda) MC^n(\mathbf{x}, \mathbf{y}) - MC^{n'}(\mathbf{x}, \mathbf{y}) \right\|_2^2. \quad (22)$$

Note that the function in (22) is convex in λ ; therefore, if the derivative of it with respect to λ is zero, then the lambda that achieves the zero gradient is the minima. So, let's derivate (22) to find λ_i .

$$0 = \frac{d}{d\lambda} \sum_{\substack{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{val}} \\ s_h(\mathbf{x}) \in Q_i}} \left\| \lambda \text{Amor}(\mathbf{x}, \mathbf{y}) + (1 - \lambda) MC^n(\mathbf{x}, \mathbf{y}) - MC^{n'}(\mathbf{x}, \mathbf{y}) \right\|_2^2 \quad (23)$$

$$= 2 \sum_{\substack{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{val}} \\ s_h(\mathbf{x}) \in Q_i}} \lambda \| MC^n(\mathbf{x}, \mathbf{y}) - \text{Amor}(\mathbf{x}, \mathbf{y}) \|^2 \quad (24)$$

$$- 2 \sum_{\substack{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{val}} \\ s_h(\mathbf{x}) \in Q_i}} \langle MC^n(\mathbf{x}, \mathbf{y}) - MC^{n'}(\mathbf{x}, \mathbf{y}), MC^n(\mathbf{x}, \mathbf{y}) - \text{Amor}(\mathbf{x}, \mathbf{y}) \rangle \quad (25)$$

From (25) we conclude the proof by showing that

$$\lambda_i = \lambda = \frac{\sum_{\substack{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{val}} \\ s_h(\mathbf{x}) \in Q_i}} \langle MC^n(\mathbf{x}, \mathbf{y}) - MC^{n'}(\mathbf{x}, \mathbf{y}), MC^n(\mathbf{x}, \mathbf{y}) - \text{Amor}(\mathbf{x}, \mathbf{y}) \rangle}{\sum_{\substack{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{val}} \\ s_h(\mathbf{x}) \in Q_i}} \| MC^n(\mathbf{x}, \mathbf{y}) - \text{Amor}(\mathbf{x}, \mathbf{y}) \|^2}. \quad (26)$$

□

Theorem 2 ($\lambda_i \approx \lambda_i^{\text{opt}}$). *Let the Monte Carlo explanation used to provide recourse MC^n to be different enough from the amortized explainer, i.e., $\mathbb{E} [\| MC^n(X, Y) - \text{Amor}(X, Y) \|^2] = \mu > 0$. Also,*

assume that $MC^{n'}$ is a good Monte Carlo approximation for the high-quality explainer HQ, i.e., $\mathbb{E} \left[\left\| MC^{n'}(X, Y) - HQ(X, Y) \right\|^2 \right] = \mu^*$ for $\epsilon > \frac{\sqrt{5\mu^*}}{\mu}$. Recall that $\mathbf{x} \in \mathbb{R}^d$. If the explanations are bounded, i.e., $\|MC^n(\mathbf{x}, \mathbf{y})\|, \|Amor(\mathbf{x}, \mathbf{y})\|, \|HQ(\mathbf{x}, \mathbf{y})\| < Cd$ for some $C > 0$ then

$$\Pr[|\lambda_i - \lambda_i^{\text{opt}}| > \epsilon] \leq e^{-\frac{\mu^2|Q_i|}{4Cd}} + e^{-\frac{\mu^4\epsilon^4|Q_i|}{400Cd}}, \quad (27)$$

where $|Q_i|$ is the number of points \mathbf{x} in the validation dataset \mathcal{D}_{val} that are in the bin Q_i .

Proof. Denote $|Q_i| = |\{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{val}}, \text{ s.t. } s_h(\mathbf{x}) \in Q_i\}|$.

We start by showing that if $\mathbb{E} \left[\left\| MC^n(X, Y) - Amor(X, Y) \right\|^2 \right] = \mu$ then

$$\Pr \left[\frac{1}{|Q_i|} \sum_{\substack{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{val}} \\ s_h(\mathbf{x}) \in Q_i}} \left\| MC^n(\mathbf{x}, \mathbf{y}) - Amor(\mathbf{x}, \mathbf{y}) \right\|^2 \leq \frac{\mu}{2} \right] \quad (28)$$

$$= \Pr \left[\mu - \frac{1}{|Q_i|} \sum_{\substack{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{val}} \\ s_h(\mathbf{x}) \in Q_i}} \left\| MC^n(\mathbf{x}, \mathbf{y}) - Amor(\mathbf{x}, \mathbf{y}) \right\|^2 \geq \frac{\mu}{2} \right] \quad (29)$$

$$\leq e^{-\frac{\mu^2|Q_i|}{4Cd}}. \quad (30)$$

Where the inequality in (30) follows from Hoeffding's inequality and the fact that:

$$\left\| MC^n(\mathbf{x}, \mathbf{y}) - Amor(\mathbf{x}, \mathbf{y}) \right\|^2 \leq \left\| MC^n(\mathbf{x}, \mathbf{y}) \right\|^2 + \left\| Amor(\mathbf{x}, \mathbf{y}) \right\|^2 \leq 2Cd. \quad (31)$$

Second, we recall that $\mathbb{E} \left[\left\| MC^{n'}(X, Y) - HQ(X, Y) \right\|^2 \right] = \mu^* \leq \frac{\mu^2\epsilon^2}{5}$. Then, we have that

$$\Pr \left[\frac{1}{|Q_i|} \sum_{\substack{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{val}} \\ s_h(\mathbf{x}) \in Q_i}} \left\| HQ(\mathbf{x}, \mathbf{y}) - Amor(\mathbf{x}, \mathbf{y}) \right\|^2 \geq \epsilon^2 \frac{\mu^2}{4} \right] \quad (32)$$

$$= \Pr \left[\frac{1}{|Q_i|} \sum_{\substack{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{val}} \\ s_h(\mathbf{x}) \in Q_i}} \left\| HQ(\mathbf{x}, \mathbf{y}) - Amor(\mathbf{x}, \mathbf{y}) \right\|^2 - \mu^* \geq \epsilon^2 \frac{\mu^2}{4} - \mu^* \right] \quad (33)$$

$$\leq \Pr \left[\frac{1}{|Q_i|} \sum_{\substack{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{val}} \\ s_h(\mathbf{x}) \in Q_i}} \left\| HQ(\mathbf{x}, \mathbf{y}) - Amor(\mathbf{x}, \mathbf{y}) \right\|^2 - \mu^* \geq \epsilon^2 \frac{\mu^2}{20} \right] \quad (34)$$

$$\leq e^{-\frac{\mu^4\epsilon^4|Q_i|}{400Cd}}. \quad (35)$$

Where the inequality in (35) follows from Hoeffding's inequality and the fact that:

$$\left\| HQ(\mathbf{x}, \mathbf{y}) - Amor(\mathbf{x}, \mathbf{y}) \right\|^2 \leq \left\| HQ(\mathbf{x}, \mathbf{y}) \right\|^2 + \left\| Amor(\mathbf{x}, \mathbf{y}) \right\|^2 \leq 2Cd. \quad (36)$$

Third, notice by directly applying Theorem 1 and replacing the Monte Carlo explanation by the high-quality explanation, we have that

$$\lambda_i^{\text{opt}} = \frac{\sum_{\substack{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{val}} \\ s_h(\mathbf{x}) \in Q_i}} \langle MC^n(\mathbf{x}, \mathbf{y}) - HQ(\mathbf{x}, \mathbf{y}), MC^n(\mathbf{x}, \mathbf{y}) - Amor(\mathbf{x}, \mathbf{y}) \rangle}{\sum_{\substack{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{val}} \\ s_h(\mathbf{x}) \in Q_i}} \left\| MC^n(\mathbf{x}, \mathbf{y}) - Amor(\mathbf{x}, \mathbf{y}) \right\|^2}. \quad (37)$$

Hence, we can write $\lambda_i^{\text{opt}} - \lambda_i$ as

$$|\lambda_i^{\text{opt}} - \lambda_i| \tag{38}$$

$$= \left| \frac{\sum_{\substack{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{val}} \\ s_h(\mathbf{x}) \in Q_i}} \langle \text{MC}^{n'}(\mathbf{x}, \mathbf{y}) - \text{HQ}(\mathbf{x}, \mathbf{y}), \text{MC}^n(\mathbf{x}, \mathbf{y}) - \text{Amor}(\mathbf{x}, \mathbf{y}) \rangle}{\sum_{\substack{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{val}} \\ s_h(\mathbf{x}) \in Q_i}} \|\text{MC}^n(\mathbf{x}, \mathbf{y}) - \text{Amor}(\mathbf{x}, \mathbf{y})\|^2} \right| \tag{39}$$

$$\leq \frac{\left(\sum_{\substack{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{val}} \\ s_h(\mathbf{x}) \in Q_i}} \|\text{MC}^{n'}(\mathbf{x}, \mathbf{y}) - \text{HQ}(\mathbf{x}, \mathbf{y})\|_2^2 \|\text{MC}^n(\mathbf{x}, \mathbf{y}) - \text{Amor}(\mathbf{x}, \mathbf{y})\|_2^2 \right)^{1/2}}{\sum_{\substack{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{val}} \\ s_h(\mathbf{x}) \in Q_i}} \|\text{MC}^n(\mathbf{x}, \mathbf{y}) - \text{Amor}(\mathbf{x}, \mathbf{y})\|^2}, \tag{40}$$

where the last inequality (40) comes from the Cauchy–Schwarz inequality. Denote the denominator in (40) by Δ , i.e.,

$$\sum_{\substack{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{val}} \\ s_h(\mathbf{x}) \in Q_i}} \|\text{MC}^n(\mathbf{x}, \mathbf{y}) - \text{Amor}(\mathbf{x}, \mathbf{y})\|^2 = \Delta.$$

Lastly, notice that $\text{MC}^{n'}(\mathbf{x}, \mathbf{y})$ is sampled independently of $\text{MC}^n(\mathbf{x}, \mathbf{y})$ and that $\text{HQ}(\mathbf{x}, \mathbf{y})$ is deterministic. Therefore:

$$\Pr[|\lambda_i^{\text{opt}} - \lambda_i| \geq \epsilon] \tag{41}$$

$$\leq \Pr \left[\frac{\left(\sum_{\substack{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{val}} \\ s_h(\mathbf{x}) \in Q_i}} \|\text{MC}^{n'}(\mathbf{x}, \mathbf{y}) - \text{HQ}(\mathbf{x}, \mathbf{y})\|_2^2 \|\text{MC}^n(\mathbf{x}, \mathbf{y}) - \text{Amor}(\mathbf{x}, \mathbf{y})\|_2^2 \right)^{1/2}}{\sum_{\substack{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{val}} \\ s_h(\mathbf{x}) \in Q_i}} \|\text{MC}^n(\mathbf{x}, \mathbf{y}) - \text{Amor}(\mathbf{x}, \mathbf{y})\|^2} \geq \epsilon \right] \tag{42}$$

$$\leq \Pr \left[\frac{\sum_{\substack{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{val}} \\ s_h(\mathbf{x}) \in Q_i}} \|\text{MC}^{n'}(\mathbf{x}, \mathbf{y}) - \text{HQ}(\mathbf{x}, \mathbf{y})\|_2^2 \|\text{MC}^n(\mathbf{x}, \mathbf{y}) - \text{Amor}(\mathbf{x}, \mathbf{y})\|_2^2}{\Delta^2} \geq \epsilon^2 \right] \tag{43}$$

$$\leq \Pr \left[\frac{\sum_{\substack{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{val}} \\ s_h(\mathbf{x}) \in Q_i}} \|\text{MC}^{n'}(\mathbf{x}, \mathbf{y}) - \text{HQ}(\mathbf{x}, \mathbf{y})\|_2^2 \|\text{MC}^n(\mathbf{x}, \mathbf{y}) - \text{Amor}(\mathbf{x}, \mathbf{y})\|_2^2}{\Delta^2} \geq \epsilon^2 \mid \Delta \leq \frac{\mu}{2} \right]$$

$$\times \Pr \left[\Delta \leq \frac{\mu}{2} \right]$$

$$+ \Pr \left[\frac{\sum_{\substack{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{val}} \\ s_h(\mathbf{x}) \in Q_i}} \|\text{MC}^{n'}(\mathbf{x}, \mathbf{y}) - \text{HQ}(\mathbf{x}, \mathbf{y})\|_2^2 \|\text{MC}^n(\mathbf{x}, \mathbf{y}) - \text{Amor}(\mathbf{x}, \mathbf{y})\|_2^2}{\Delta^2} \geq \epsilon^2 \mid \Delta > \frac{\mu}{2} \right]$$

$$\times \Pr \left[\Delta > \frac{\mu}{2} \right] \tag{44}$$

$$\leq \Pr \left[\sum_{\substack{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{val}} \\ s_h(\mathbf{x}) \in Q_i}} \|\text{MC}^{n'}(\mathbf{x}, \mathbf{y}) - \text{HQ}(\mathbf{x}, \mathbf{y})\|_2^2 \|\text{MC}^n(\mathbf{x}, \mathbf{y}) - \text{Amor}(\mathbf{x}, \mathbf{y})\|_2^2 \geq \epsilon^2 \frac{\mu^2}{4} \right]$$

$$+ \Pr \left[\Delta \leq \frac{\mu}{2} \right] \tag{45}$$

$$\leq e^{-\frac{\mu^2|Q_i|}{4Cd}} + e^{-\frac{\mu^4\epsilon^4|Q_i|}{400Cd}}. \quad (46)$$

Where the inequality in (42) is a direct application of 40, the inequality in (44) comes from simply conditioning, the inequality in (45) comes from the fact that probabilities are bounded by one getting rid of the first term in (45) (first out of lines) and the fourth term in (45) (forth out of lines) and the fact that $\text{MC}^{n'}(\mathbf{x}, \mathbf{y})$ is sampled independently of $\text{MC}^n(\mathbf{x}, \mathbf{y})$ and that $\text{HQ}(\mathbf{x}, \mathbf{y})$ is deterministic. Finally, the last inequality in (46) comes from applying (30) and (35).

Hence, from (46), we conclude that

$$\Pr[|\lambda_i^{\text{opt}} - \lambda_i| \geq \epsilon] \leq e^{-\frac{\mu^2|Q_i|}{4Cd}} + e^{-\frac{\mu^4\epsilon^4|Q_i|}{400Cd}}. \quad (47)$$

□

Proposition 2 (Coverage for Inference Budget). *Let $N_{\text{budget}} \geq 1$ be the set inference budget, and assume that the Monte Carlo method $\text{MC}^n(\mathbf{x}, \mathbf{y})$ uses n model inferences. Then, the coverage level α should be chosen such that*

$$\underset{\alpha \in [0,1]}{\text{argmin}} \{ \mathbb{E} [N(\text{SE}(\mathbf{x}, \mathbf{y}))] \leq N_{\text{budget}} \} = \frac{n+1 - N_{\text{budget}}}{n}. \quad (48)$$

Recall that Shapley Value Sampling with parameter m performs $1 + dm$ inferences ($\mathbf{x} \in \mathbb{R}^d$), and Kernel Shap with parameter m performs m inferences.

Proof. Let $\alpha \in [0, 1]$, then an α portion of examples receive explanations from the amortized explainer, i.e., they receive one inference, and $1 - \alpha$ portion of examples receive explanations with initial guess, i.e., n model inferences. Therefore, the expected number of model inferences per instance is given by (49).

$$\mathbb{E} [N(\text{SE}(\mathbf{x}, \mathbf{y}))] = \alpha + (1 - \alpha)(n + 1) \quad (49)$$

In order for the inference budget to be followed, it is necessary that

$$\mathbb{E} [N(\text{SE}(\mathbf{x}, \mathbf{y}))] = \alpha + (1 - \alpha)(n + 1) \leq N_{\text{budget}}. \quad (50)$$

From (50), we conclude that:

$$\alpha \geq \frac{n+1 - N_{\text{budget}}}{n}, \quad (51)$$

Hence,

$$\underset{\alpha \in [0,1]}{\text{argmin}} \{ \mathbb{E} [N(\text{SE}(\mathbf{x}, \mathbf{y}))] \leq N_{\text{budget}} \} = \frac{n+1 - N_{\text{budget}}}{n}. \quad (52)$$

□