

# CONLUX: CONCEPT-BASED LOCAL UNIFIED EXPLANATIONS

**Anonymous authors**

Paper under double-blind review

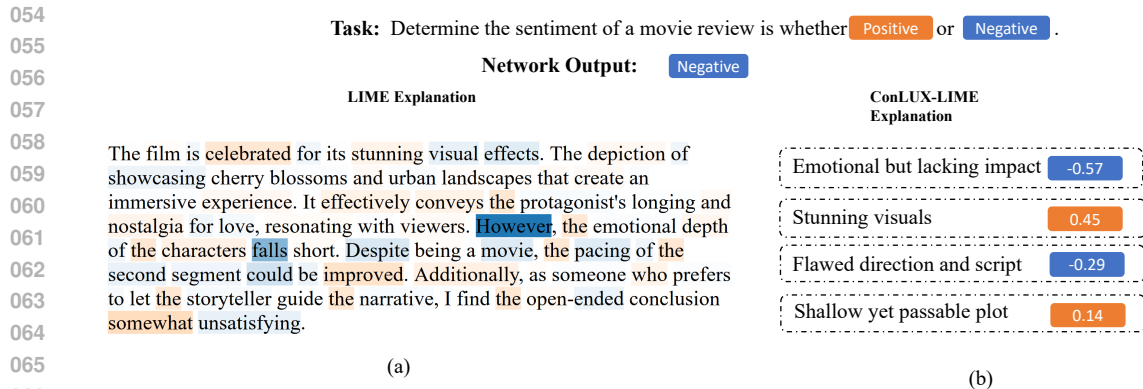
## ABSTRACT

With the rapid advancements of various machine learning models, there is a significant demand for model-agnostic explanation techniques, which can explain these models across different architectures. Mainstream model-agnostic explanation techniques generate local explanations based on basic features (e.g., words for text models and (super-)pixels for image models). However, these explanations often do not align with the decision-making processes of the target models and end-users, resulting in explanations that are unfaithful and difficult for users to understand. On the other hand, concept-based techniques provide explanations based on high-level features (e.g., topics for text models and objects for image models), but most are model-specific or require additional pre-defined external concept knowledge. To address this limitation, we propose ConLUX, a general framework to provide concept-based local explanations for any machine learning models. Our key insight is that we can automatically extract high-level concepts from large pre-trained models, and uniformly extend existing local model-agnostic techniques to provide unified concept-based explanations. We have instantiated ConLUX on four different types of explanation techniques: LIME, Kernel SHAP, Anchor, and LORE, and applied these techniques to text and image models. Our evaluation results demonstrate that 1) compared to the vanilla versions, ConLUX offers more faithful explanations and makes them more understandable to users, and 2) by offering multiple forms of explanations, ConLUX outperforms state-of-the-art concept-based explanation techniques specifically designed for text and image models, respectively.

## 1 INTRODUCTION

As machine learning models become more complex and popular, it has become an emerging topic to explain the rationale behind their decisions. In particular, as the structure of machine learning models diversifies and evolves rapidly, and effective closed-source models (e.g., GPT-4 (Achiam et al., 2023) and Gemini (et al., 2024b)) become more prevalent, model-agnostic explanation techniques show their appeal due to their adaptability to various models and tasks (Wang, 2024). These techniques consider target models as black boxes, so they can explain any machine learning model without requiring any knowledge of the model’s internal structure. This paper addresses the challenge of incorporating high-level concepts into local model-agnostic techniques to explain the decision-making processes of various machine learning models, including large language models (LLMs).

To faithfully explain the behavior of machine learning models, it is essential to provide explanations built from language components aligned with the decision process of the target models; to make explanations easy to understand, it is also crucial to provide explanations built from user-friendly language components (Poeta et al., 2023a). Unfortunately, mainstream model-agnostic explanation techniques often fail to meet both requirements, as they provide explanations built from basic features (e.g., words for text models and (super-)pixels for image models) (Ribeiro et al., 2016; Lundberg & Lee, 2017; Ribeiro et al., 2018; Guidotti et al., 2018). In contrast, many concept-based techniques provide explanations based on high-level features (e.g., topics in texts and objects in images) (Poeta et al., 2023a). These techniques either utilize the internal information of the target models like gradients, activations, and attention weights (Zhang et al., 2021b; Yeh et al., 2020; 2019b; Cunningham et al., 2023; Ghorbani et al., 2019b; Crabbé & van der Schaar, 2022; Fel et al., 2023), or pre-defined external knowledge (El Shawi, 2024; Widmer et al., 2022) to build high-level



068 Figure 1: Example explanation (a) is generated by LIME, demonstrating how each word in the input sentence contributes to the target model’s prediction. The color intensity reflects the magnitude of the weight, with deeper hues indicating larger absolute values. Example explanation (b) is generated by ConLUX-augmented LIME, providing an explanation based on high-level concepts.

075 concepts. This limits these techniques to specific types of models or tasks. Furthermore, while there are different forms of explanations (e.g. feature attributions, rules) for various purposes (Zhang et al., 2021c), existing concept-based explanations mainly focus only on attributions, which limits their fidelity and applicability.

079 To bridge this gap, we aim to elevate the explanations of various forms provided by existing model-agnostic techniques from feature-level to concept-level. As we focus on explaining the decision-making process of machine learning models to end-users, we put our emphasis on local explanations, which are more tractable for end-users. We find that existing local model-agnostic techniques all follow similar workflows, which allows us to introduce a unified way to elevate all these techniques from feature level to concept level. This transition necessitates automating the concept extraction process and establishing bidirectional mappings between concept representations and feature representations for given input data. Noticing that existing works have utilized large pre-trained models to extract concepts and represent input data at the concept level for specific tasks (Ludan et al., 2023; Sun et al., 2023), we generalize these findings and further observe that large models also have the ability to map concept-level representations back to the feature-level. To this end, we propose ConLUX, a general framework that automatically incorporates high-level concepts into various existing local model-agnostic techniques for any machine learning models, and provides local explanations in various forms for diverse user needs.

093 We take three mainstream local model-agnostic techniques, LIME (Ribeiro et al., 2016), Anchor (Ribeiro et al., 2018), and LORE (Guidotti et al., 2018) as examples to illustrate how ConLUX improves local model-agnostic explanations.

096 Figure 1 shows a LIME explanation for a BERT-based sentiment analysis model on a movie review. The target model classifies the sentence as negative. LIME provides feature-level attributions, indicating how each word contributes to the model’s prediction. In this case, LIME assigns high negative scores to the words “however” and “falls”, which indicates that these words contribute much to the negative prediction. However, this explanation is unfaithful and hard to understand by end-users. For example, the word “however” is assigned a high negative score, but it functions as a conjunction and does not directly convey sentiment (Liu & Zhang, 2023). Moreover, such confusing attributions, combined with an overwhelming amount of attribution information, complicate the explanation for end-users. ConLUX addresses these issues by elevating the explanation from feature-level to concept-level. ConLUX-augmented LIME extracts the main topics of the input sentence using GPT-4o, and then provides attribution-based explanations with these topics. From this explanation, users can easily understand that the negative prediction is mainly because the sentence mentions the movie’s poor performance in “emotion impact” and “direction and script”, while the part about “stunning visuals” and “passable plot” also expresses some positive sentiment. This ex-

108  
109  
110  
111  
112  
113  
114  
115  
116  
117  
118  
119  
120  
121  
122  
123  
124  
125  
126  
127  
128  
129  
130  
131  
132  
133  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148  
149  
150  
151  
152  
153  
154  
155  
156  
157  
158  
159  
160  
161

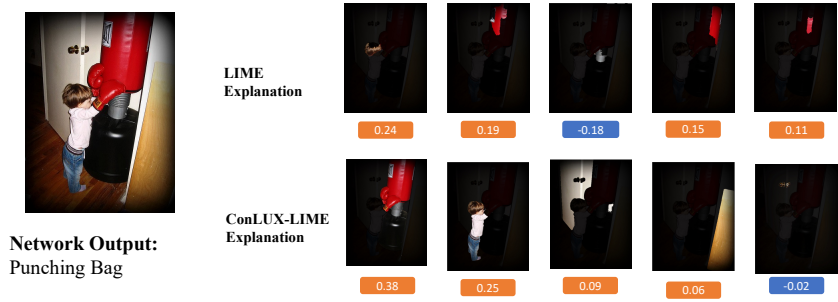


Figure 2: Example explanations generated by LIME (upper) and ConLUX-augmented LIME (lower) for an image classification task.

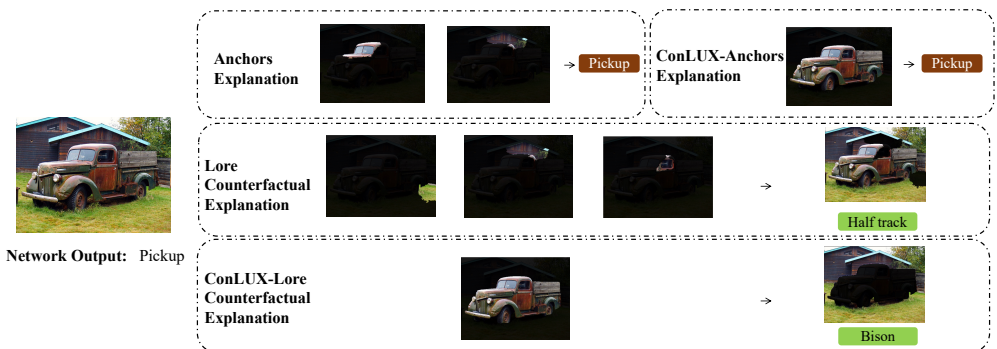


Figure 3: Example Anchors, LORE explanations and their ConLUX-augmented versions.

planation faithfully reflects the decision process of the target model and is more understandable to end-users.

Similar issues exist in the explanation of image models. We use YOLOv8 (Jocher et al., 2023) to perform an image classification task on ImageNet (Deng et al., 2009) dataset. Figure 2 shows a LIME explanation for an image classified as a *punching bag*. LIME attributes high importance to some fragmented superpixels. End-users can hardly understand why these parts are important for the model’s prediction. ConLUX-augmented LIME provides explanations based on objects detected by Segment Anything Model (SAM) (Kirillov et al., 2023), and attributes high importance to the punching bag itself and the kid in the image. This explanation is more faithful and understandable. End-users can easily realize that the model does not perform perfectly when classifying this image to a *punching bag*.

Figure 3 shows the explanations generated by Anchor, LORE, and their ConLUX-augmented versions for an image classified as a *pickup*. Anchors provides rule-based sufficient conditions (referred to as *anchors*) for the target model’s prediction. The vanilla anchor indicates that parts of the car and the background house guarantee the prediction as a *pickup*. With ConLUX, end-users can easily understand that the model classifies the image as a *pickup* exactly because it indeed detects the pickup truck in the image. LORE provides rule-based sufficient conditions and counterfactual explanations. Figure 3 shows the counterfactual explanations, which show users how to change the model’s prediction by modifying the input image. The vanilla LORE explanation indicates that if we mask a part of the grass, the background house and the whole side window of the truck will change the model’s prediction to a *Half track*. In contrast, ConLUX-augmented LORE indicates that users can simply mask the pickup truck to change the model’s prediction to a *Bison*, which is more user-friendly.

The preceding examples indicate that feature-level explanations are hard to understand by end-users. On the other hand, as high-level concepts align with the decision process of target models and users better (Zhang et al., 2021a; Ghorbani et al., 2019a; Kim et al., 2018; Sun et al., 2023), ConLUX addresses this by providing concept-level explanations, and the examples demonstrate that ConLUX

162 makes explanations more understandable to end-users. Moreover, our empirical evaluation shows  
 163 these concept-level explanations are also more faithful to the models. Finally, benefiting from the  
 164 various types of explanations provided by existing local model-agnostic techniques, ConLUX can  
 165 provide rich explanations including attributions, sufficient conditions, and counterfactuals, satisfying  
 166 diverse user needs and offering a more comprehensive understanding of the target models. This fills  
 167 the current gap in concept-based explanations, which lack forms beyond attributions.

168 To elevate the explanations provided by existing local model-agnostic techniques from feature-level  
 169 to concept-level, we modify these techniques in a uniform and lightweight way based on their two  
 170 commonalities: 1) these techniques use basic features as language components to build explanations;  
 171 2) these techniques use a perturbation model, which generates samples similar to the given input by  
 172 changing some of its feature values, to capture the local behavior of the target model at the feature  
 173 level. To this end, by only elevating the language components to high-level concepts and extending  
 174 the perturbation model to generate samples by changing high-level concepts, ConLUX extends these  
 175 techniques to provide concept-level explanations.

176 We evaluated ConLUX on explaining two sentiment analysis models (BERT, Llama 3.1(et al.,  
 177 2024a)) and three image classification models(YOLOv8, Vision Transformer (Oquab et al., 2023;  
 178 Darcet et al., 2023), and ResNet-50 (He et al., 2016)). Our evaluation results demonstrate that  
 179 ConLUX improves the fidelity of Anchors, LIME, LORE, and Kernel SHAP explanations by  
 180 82.21%, 48.59%, 149.93%, and 48.27% respectively, and by considering various forms of explan-  
 181 ations together, ConLUX outperforms state-of-the-art concept-based explanation techniques specifi-  
 182 cally designed for text models (TBM (Ludan et al., 2023)) and image models (EAC (Sun et al.,  
 183 2023)), respectively.

## 184 2 PRELIMINARIES

185 In this section, we introduce the background knowledge and notations used in this paper.

186  
 187 **Machine Learning Models.** We consider a machine learning model as a black-box function  $f$  that  
 188 maps an input vector  $\mathbf{x}$  to an output scalar  $f(\mathbf{x})$ . Formally, we let  $f : \mathbb{X} \rightarrow \mathbb{R}$ , where  $\mathbb{X}$  is the input  
 189 domain. For models that take fixed-dimension inputs, let  $\mathbb{X} = \mathbb{R}^n$ . For models capable of handling  
 190 inputs of arbitrary dimensions, let  $\mathbb{X} = \cup_{i=1}^{\infty} \mathbb{R}^i$ . Let  $x_i$  denote the  $i$ -th feature value of  $\mathbf{x}$ .

191  
 192 **Local Model-Agnostic Explanation Techniques.** A local model-agnostic explanation technique  
 193  $t$  takes a model  $f$  and an input  $\mathbf{x}$ , and generates a local explanation  $g_{f,\mathbf{x}}$  to describe the behavior of  
 194  $f$  around  $\mathbf{x}$ , i.e.,  $g_{f,\mathbf{x}} := t(f, \mathbf{x})$ .  $g_{f,\mathbf{x}}$  ( $g$  for short) is an expression formed with **predicates**. Each  
 195 predicate  $p$  maps an input  $\mathbf{x}$  to a binary value, i.e.,  $p : \mathbb{X} \rightarrow \{0, 1\}$ , indicating whether  $\mathbf{x}$  satisfies a  
 196 specific condition.

197  
 198 As shown in Figure 4, existing local model-agnostic explanation techniques generate explanations  
 199 following a similar workflow:

- 200 1. **Producing Predicates:** These techniques first generate a **set of predicates**  $\mathbb{P}$  based on the  
 201 input  $\mathbf{x}$ .
- 202 2. **Generating Samples:** The underlying **perturbation model**  $t_{per}$  generates a set of samples  
 203  $\mathbb{X}_s^b$  in **predicate representations**, where each sample  $\mathbf{z}^b \in \mathbb{X}_s^b$  is a binary vector in  $\{0, 1\}^d$   
 204 and  $z_i^b$  indicates whether the sample satisfies the  $i$ -th predicate in  $\mathbb{P}$ . The perturbation model  
 205 then transforms the samples  $\mathbb{X}_s^b$  back to the original input space to get  $\mathbb{X}_s$  and  $f(\mathbb{X}_s)$ .
- 206 3. **Learning Explanation:** The underlying **learning algorithm** generates the local explana-  
 207 tion  $g_{f,\mathbf{x}}$  consisting of predicates in  $\mathbb{P}$  using  $\mathbb{X}_s$  and  $f(\mathbb{X}_s)$ .

208  
 209 Mainstream local model-agnostic explanation techniques like Anchors, LIME, LORE, and Kernel  
 210 SHAP, all follow this workflow. They use the same kinds of predicate sets and perturbation mod-  
 211 els, and use different learning algorithms to generate explanations with different properties. In the  
 212 following, we introduce the main components of the explanation techniques.

213 **Predicate Sets.** Given an input  $\mathbf{x}$ , the corresponding predicate set  $\mathbb{P}$  is defined as follows:

$$214 \mathbb{P} = \{p_i | i \in [1, d]\},$$

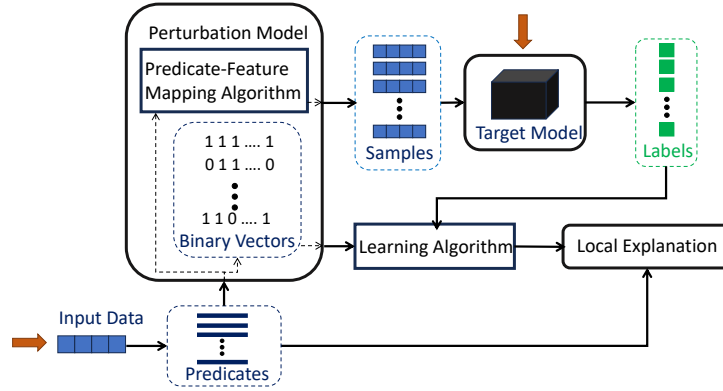


Figure 4: The workflow of generating explanations by a local model-agnostic explanation technique.

where  $d$  is the number of predicates in  $\mathbb{P}$ , a hyperparameter set by users or according to the input data  $\mathbf{x}$ . Each  $p_i$  is a **feature predicate** that constrains the value of a set of feature values in  $\mathbf{x}$ , i.e.  $p_i(\mathbf{z}) : \bigwedge_{j \in \mathbb{A}_i} \mathbb{1}_{\text{ran}(\mathbf{x}, j)}(z_j)$ , where  $\mathbb{A}_i$  is the set of indices of features that  $p_i$  constrains. Specifically,  $\{\mathbb{A}_1, \mathbb{A}_2, \dots, \mathbb{A}_d\}$  forms a partition of  $\{1, 2, \dots, |\mathbf{x}|\}$ . Each  $\text{ran}(\mathbf{x}, j)$  is a set containing  $x_j$ , which is set according to the type of input data. For example, we can use  $\text{ran}(\mathbf{x}, j) = (x_j - \epsilon, x_j + \epsilon)$  for continuous values, and  $\text{ran}(\mathbf{x}, j) = \{x_j\}$  for categorical values. The predicate  $p_i$  is a conjunction of indicator functions, each of which checks if a sample  $z$  has a similar value to  $x_j$  (i.e.,  $z_j \in \text{ran}(\mathbf{x}, j)$ ).

**Predicate Representations.** The predicate representation  $\mathbf{z}^b \in \{0, 1\}^d$  of a sample  $z$  is a binary vector where  $z_i^b = p_i(z)$ .

**Perturbation Models.** The perturbation model  $t_{per}$  first randomly selects  $\mathbb{X}_s^b \subseteq \{0, 1\}^d$  as the predicate representations of the samples. Then, it transforms  $\mathbb{X}_s^b$  back to the original input space to get  $\mathbb{X}_s$ . For each  $\mathbf{z}^b \in \mathbb{X}_s^b$ , a predicate-to-feature mapping function  $t_{p2f} : \{0, 1\}^d \rightarrow \mathbb{X}$  transforms  $\mathbf{z}^b$  to  $\mathbf{z}$  as follows: if  $z_i^b = 1$ , then for each  $j \in \mathbb{A}_i$ , set  $z_j = x_j$ ; otherwise, set each  $z_j$  to a masked value, or a random value sampled from  $\text{per}(\mathbf{x}, j) \setminus \text{ran}(\mathbf{x}, j)$ , where  $\text{per}(\mathbf{x}, j)$  is a perturbation range with  $\text{per}(\mathbf{x}, j) \supset \text{ran}(\mathbf{x}, j)$ .

**Learning Algorithms and Explanations.** The learning algorithm learns an understandable expression  $g$  as an explanation. In Anchors,  $g$  is a conjunction of predicates that provides a sufficient condition for producing  $f(\mathbf{x})$  as output, i.e.,  $f(\mathbf{z}) = f(\mathbf{x})$  if  $g(\mathbf{z}) = 1$ , and  $g(\mathbf{z}) = \bigwedge_{p \in \mathbb{Q}} p(\mathbf{z})$ , where  $\mathbb{Q}$  is selected by KL-LUCB algorithm (Kaufmann & Kalyanakrishnan, 2013). In LIME and Kernel SHAP,  $g(\mathbf{z}) = \sum_{i=1}^d w_i p_i(\mathbf{z}) + w_0$ , where  $w_i$  is the weight of  $p_i$ , which is learned by their underlying regression algorithms. LORE first learns a decision tree with building systems like Yadt (Ruggieri, 2004), then extracts a sufficient condition to obtain  $f(\mathbf{x})$  and some counterfactual rules from the tree as explanations. The sufficient condition is similar to Anchors, while each of the counterfactual rules is in the form of  $f(\mathbf{z}) = y$  if  $g(\mathbf{z}) = 1$ , where  $y \neq f(\mathbf{x})$  and  $g(\mathbf{z}) = \bigwedge_{p \in \mathbb{Q}} p(\mathbf{z}) \wedge \bigwedge_{p \in \mathbb{C}} \neg p(\mathbf{z})$ , and  $\mathbb{Q}$  and  $\mathbb{C}$  are extracted from the decision tree.

**An Example.** For a text input *I love this movie so much*, these techniques can let each  $p_i$  constrains only one feature, and produce six predicates in the form of  $p_i(\mathbf{z}) := \mathbb{1}_{z_i=x_i}$ . For another text input  $z = I \text{ love this } [MASK] \text{ so } [MASK]$ , the predicate representation of  $z$  is  $p_1(z) p_2(z) p_3(z) p_4(z) p_5(z) p_6(z) = 1 1 1 0 1 0$ . The perturbation model will generate samples in predicate representation, then  $t_{p2f}$  will transform samples back to the origin input space. For example, a sample 0 1 0 1 1 1 is generated and  $t_{p2f}$  maps it to *[MASK] love [MASK] movie so much*. Consequently, these techniques will use the output of these samples and the samples' predicate presentation to learn an expression, and build the explanation with the predicates.

Limited by the predicate sets and perturbation models, existing local model-agnostic explanation techniques can only generate explanations based on the constraints of feature values, which limits their effectiveness in explaining the behavior of the model.

### 3 THE CONLUX FRAMEWORK

In this section, we propose ConLUX, a general framework to provide concept-based local explanations based on existing local model-agnostic explanation techniques without significantly changing their core components.

We introduce ConLUX in three steps: 1) defining concept-based local explanations and concept predicates, 2) showing the modifications to the explanation techniques, and 3) demonstrating the augmented workflow.

#### 3.1 CONCEPT-BASED LOCAL EXPLANATIONS

As we discussed in Section 2, though the form of the explanations varies, they are all built from predicates in  $\mathbb{P}$ . Elevating the predicates to concept level is the key to providing concept-based explanations.

Definition of high-level concepts varies, as Molnar (2020) mentions, “A concept can be any abstraction, such as a color, an object, or even an idea.” Here, to provide explanations that are easier to understand by end-users, we define a concept predicate as follows:

**Definition 1** (Concept Predicate). *Given an input  $\mathbf{x}$ , a concept predicate  $p^c$  is a function that maps  $\mathbf{x}$  to a binary value, i.e.,  $p^c : \mathbb{X} \rightarrow \{0, 1\}$ , and satisfies the following properties:*

1. **Descriptive:** *The concept predicate  $p^c$  can be concisely and intuitively described in natural language.*
2. **Human Evaluable:** *The truth of  $p^c(\mathbf{x})$  can be readily assessed by a human user.*

The preceding two properties ensure that the concepts are easy to understand. Here, we provide two examples of concept predicates for text and image models in the following:

**Examples.** For text models, we can define a concept predicate as follows:

- **Concept Name:** Poor Visual Effects and Cinematography
- **Description:** The input text mentioned that the visual effects and cinematography are lacking, failing to create an appealing aesthetic.

For image models, we use objects in an image to define a concept predicate. As shown in Figure 2, we can easily describe the concept predicate as “there is a punching bag in the image”, “there is a kid in the image”, etc.

We then define a concept-based local explanation as follows:

**Definition 2** (Concept-Based Local Explanation). *A concept-based local explanation  $g_{f,\mathbf{x}}^c$  is an expression formed with concept predicates to describe the behavior of  $f$  around  $\mathbf{x}$ .*

As existing local explanation model-agnostic techniques provide various kinds of explanations like attributions, sufficient conditions, and counterfactuals, ConLUX can elevate all these explanations to concept level and provide users a unified interface to obtain various kinds of explanations with a single click. Additionally, We denote such a set of various kinds of explanations as a **ConLUX unified explanation**, which provides higher fidelity and offers a more comprehensive view than a single form of explanation.

#### 3.2 AUGMENTING EXPLANATION TECHNIQUES

As shown in Figure 4, to produce concept predicates, we should first extract high-level concepts based on the input  $\mathbf{x}$ ; to provide explanations at concept level, we should replace the feature predicates in  $\mathbb{P}$  with concept predicates; to capture the local behavior of the target model at concept level, we should extend the perturbation model to generate samples by changing high-level concepts.

**Producing Concept Predicates.** We use large pre-trained models to provide high-level concepts based on the input  $x$  and the target task. For text models, following the approach of Ludan et al. (2023), we provide GPT-3.5 with task-specific information, the given input, the corresponding output, and several in-context learning examples to generate candidate concepts. These concepts are then evaluated on the input  $x$  to construct the concept predicate set.

For image models, we refer to Sun et al. (2023) to use SAM to detect objects in the image.

Consequently, ConLUX defines concept predicates (denoted as  $p^c$ ) using the extracted concepts, and replaces the feature predicates set  $\mathbb{P}$  with the concept predicates set  $\mathbb{P}^c$ .

**Performing Concept-Level Perturbation** The extended perturbation model  $t_{per}^c$  generates samples in concept-level representation and  $t_{p2f}^c : \{0, 1\}^{|\mathbb{P}^c|} \rightarrow \mathbb{X}$  transforms the samples back to the original input space. Different from the  $t_{per}$  simply decides whether to mask a feature value,  $t_{per}^c$  changes high-level concepts at feature level, which is more complex. Therefore,  $t_{per}^c$  is usually a more sophisticated model. Here, for text models, we use Llama 3.1 to perform the concept-feature mapping; for image models, since each object is still a set of pixels, we can use the same transformation as  $t_{per}$ . As the effectiveness and faithfulness of the text perturbation are not straightforward, we conduct an experiment to demonstrate this, as detailed in Appendix C.

### 3.3 CONLUX-AUGMENTED WORKFLOW

The ConLUX augments the workflow in Figure 4 as follows: it first extracts high-level concepts based on the input  $x$  and the target task, then follows a similar workflow as their vanilla versions, but replaces the predicate set  $\mathbb{P}$  with  $\mathbb{P}^c$  and the perturbation model  $t_{per}$  with  $t_{per}^c$ . Therefore, the ConLUX-augmented techniques can capture the local behavior of the target model at the concept level, and provide concept-based local explanations.

More details about the implementation of ConLUX can be found in Appendix A.

## 4 EMPIRICAL EVALUATION

In this section, we demonstrate the generality of ConLUX, its improvement of explanation fidelity, and the fidelity of ConLUX unified explanations by empirical evaluation. We show the generality of ConLUX by applying it to four mainstream local model-agnostic explanation techniques: Anchors, LIME, LORE, and Kernel SHAP (KSHAP for short), which provide three types of explanations: sufficient conditions, counterfactuals, and attributions. We apply them to explain various text and image models. We show the improvement of explanation fidelity by comparing the vanilla feature-level explanations with ConLUX-augmented the concept-based explanations. Moreover, we compare the fidelity of ConLUX unified explanations with two state-of-the-art concept-based explanation techniques: Textual Bottleneck Model (TBM) (Ludan et al., 2023) for text models and Explain Any Concept (EAC) (Sun et al., 2023) for image models.

### 4.1 EXPERIMENTAL SETUP

We chose sentiment analysis as the target task for text models, and image classification as the target task for image models.

**Sentiment Analysis.** Sentiment analysis models take a text sequence as input and predict if the text is positive or negative, i.e.  $f : \mathbb{X} \rightarrow \{0, 1\}$ , where  $\mathbb{X} := \bigcup_{i=1}^{\infty} \mathbb{W}^i$  is the input domain, and  $\mathbb{W}$  is the vocabulary set. We used a pre-trained BERT (Morris et al., 2020) and Llama3.1 to predict the sentiment of 200 movie reviews from the Large Movie Review Dataset (Maas et al., 2011), and explained the local behavior of the models around each input text. For vanilla techniques, we followed the settings described in Section 2. For ConLUX-augmented techniques, we set the number of concept predicates to 10, used GPT-3.5 (Brown et al., 2020) to extract high-level concepts, and Llama3.1 to perform the predicate-to-feature mapping. For TBM, we applied it to explain the same 200 movie reviews with its default settings.

**Image Classification.** Image classification models take an image as input and predict the category of the image, i.e.  $f : \mathbb{X} \rightarrow \{0, 1, \dots, m\}$ , where  $m$  is the number of categories,  $\mathbb{X} := \mathbb{R}^{3 \times h \times w}$  is the input domain, with  $h$  and  $w$  being the height and width of the image. We used a pre-trained YOLOv8, Vision Transformer (ViT) (Oquab et al., 2023; Darcet et al., 2023), and ResNet-50 (He et al., 2016) to predict the category of 1000 images from the ImageNet dataset (Deng et al., 2009), and explained the local behavior of the models around each input image. For vanilla techniques, we followed LIME to use Quickshift algorithm (Jiang et al., 2018) to obtain the superpixels, and used these super-pixels as predicates. For ConLUX-augmented techniques, we used SAM (Kirillov et al., 2023) to detect objects in the images, and used these objects as predicates. For EAC, we applied it to explain the same 1000 images with its default settings.

To evaluate the fidelity of ConLUX unified explanations, as the combination of multiple forms of explanation provides more fidelity than a single form, we used the combination ConLUX-augmented KSHAP and LORE explanations as local surrogate models. Specifically, if an input is covered by LORE’s rule, we use the LORE output; otherwise, we use the KSHAP explanation.

More details can be referred to Appendix B.

## 4.2 FIDELITY EVALUATION

### 4.2.1 EVALUATION METRICS

Fidelity reflects how faithfully the explanations describe the target model. As these techniques provide explanations in different forms, we used different metrics to evaluate their fidelity.

Following the setup in the original papers of Anchors and LORE, we used **coverage** and **precision** as fidelity metrics (which are named differently in the LORE paper). Given a target model  $f$ , an input  $\mathbf{x}$ , and a distribution  $D_{\mathbf{x}}$  derived from the perturbation model, and the corresponding explanation  $g$ , we defined the coverage as  $\text{cov}(\mathbf{x}; f, g) = \mathbb{E}_{\mathbf{z} \sim D_{\mathbf{x}}} [g(\mathbf{z})]$ , which indicates the proportion of inputs in the distribution that match the rule; we defined the precision as  $\text{prec}(\mathbf{x}; f, g) = \mathbb{E}_{\mathbf{z} \sim D_{\mathbf{x}}} [\mathbf{1}_{f(\mathbf{z})=y} | g(\mathbf{z})]$ , where  $y$  is the consequence of the rules in  $g$  with  $y = f(\mathbf{x})$  for factual rules and  $y \neq f(\mathbf{x})$  for counterfactual rules. Precision indicates the proportion of covered inputs that  $g$  correctly predicts the model outputs.

As LIME and KSHAP are attribution-based local surrogate, we used *Area Over most relevant first perturbation curve* (AOPC) (Samek et al., 2016; Modarressi et al., 2023), and  $\text{accuracy}_a$  as fidelity metrics (Balagopalan et al., 2022; Yeh et al., 2019a; Ismail et al., 2021). Given a target model  $f$ , an input  $\mathbf{x}$ , its corresponding model output  $y = f(\mathbf{x})$ , their corresponding explanation  $g$ , and  $\mathbf{x}^{(k)}$  that is generated by masking the  $k\%$  most important predicates in  $\mathbf{x}$ , AOPC and  $\text{accuracy}_a$  are defined as follows:

- **AOPC:** Let  $\text{AOPC}_k = \frac{1}{|\mathbb{T}|} \sum_{\mathbf{x}} p_f(y|\mathbf{x}) - p_f(y|\mathbf{x}^{(k)})$ , where  $p_f(y|\mathbf{x})$  is the probability of  $f$  to output  $y$  given the input  $\mathbf{x}$ , and  $\mathbb{T}$  is the set of all test inputs.  $\text{AOPC}_k$  indicates the average change of the model output when masking the  $k\%$  most important predicates. A higher  $\text{AOPC}_k$  indicates a better explanation. We calculate the AOPC curve by varying  $k$  from 0 to 100.
- **Accuracy<sub>a</sub>:**  $\text{accuracy}_a$  indicates the proportion of inputs among all  $\mathbf{x}^{(k)}$  that the target model gives the same output as the original input  $\mathbf{x}$ , i.e.  $\mathbb{E}(f(\mathbf{x}^{(k)}) = f(\mathbf{x}))$ . Specifically,  $\text{accuracy}_a$  is different from the standard accuracy, and a lower  $\text{accuracy}_a$  indicates a better explanation.

Specifically, we only considered the predicates that positively contribute to  $f(\mathbf{x})$ , and we did not use AOPC when explaining Llama 3.1, as it does not directly provide the probability for each output token.

For TBM, EAC, and ConLUX unified explanations, considering that they can all serve as local surrogate models, i.e.  $g : \mathbb{X} \rightarrow \mathbb{R}$ , we defined the metrics as follows: Given a target model  $f$ , an input  $\mathbf{x}$ , a perturbation distribution  $\mathbb{D}_{\mathbf{x}}$ , and their corresponding explanation  $g$ , a performance metric  $L$  (e.g. accuracy, F1 score, MSE, etc.), we define the (in-)fidelity as  $E_{\mathbf{z} \sim \mathbb{D}_{\mathbf{x}}} L(f(\mathbf{z}), g(\mathbf{z}))$ , which indicates the performance of using  $g$  to approximate  $f$ . Here, we used the accuracy as the performance metric. Specifically, considering the complexity of the original task, we reduced the image classification task for local surrogates to predicting whether the target model  $f$  assigns the same classification to  $\mathbf{x}'$  as it does to  $\mathbf{x}$ .



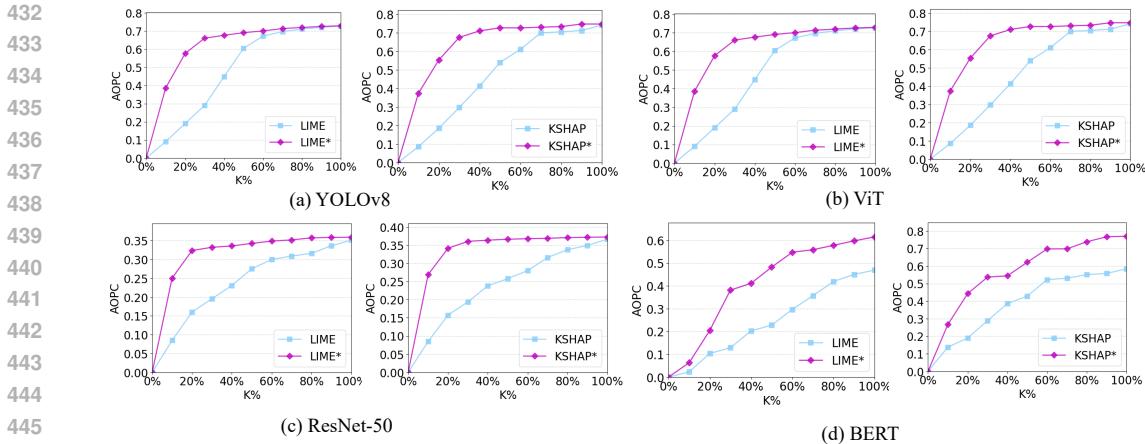


Figure 5: AOPC upon masking the K% most important predicates. We use LIME, Kernel SHAP, and their ConLUX-augmented version to explain YOLOv8, ViT, and Resnet-50 on the image classification task and BERT on the sentiment analysis task.

Table 1: Average coverage and precision (higher are better) of Anchors, LORE, and their ConLUX-augmented versions (denoted as Anchors\* and LORE\*) on two sentiment analysis models and three image classification models.

Models	Coverage (%) $\uparrow$				Precision (%) $\uparrow$			
	Anchors	Anchors*	LORE	LORE*	Anchors	Anchors*	LORE	LORE*
Llama 3.1	4.9	<b>22.5</b>	2.3	<b>21.3</b>	81.2	<b>94.2</b>	64.3	<b>76.9</b>
BERT	5.3	<b>24.4</b>	3.2	<b>20.3</b>	78.2	<b>91.0</b>	65.3	<b>79.4</b>
YOLOv8	28.6	<b>30.9</b>	20.8	<b>24.8</b>	84.3	<b>98.2</b>	87.8	<b>92.2</b>
ViT	24.6	<b>28.2</b>	21.3	<b>23.8</b>	88.7	<b>98.2</b>	89.6	<b>95.6</b>
ResNet-50	28.0	<b>30.5</b>	20.1	<b>29.7</b>	89.3	<b>99.4</b>	85.8	<b>92.6</b>

#### 4.2.2 EVALUATION RESULTS

Table 1 shows the fidelity of Anchors, LORE, and their ConLUX-augmented versions. ConLUX improves the average coverage of Anchors and LORE by 9.0% and 10.4%, and the average precision by 11.9% and 8.7%, respectively. Figure 5 and Table 2 show the fidelity of LIME, KSHAP, and their ConLUX-augmented versions. Figure 5 shows the AOPC curve of LIME and KSHAP. Each AOPC curve of ConLUX-augmented versions is higher than the vanilla counterpart. Table 2 shows the average AOPC and accuracy<sub>a</sub>. ConLUX improves the average AOPC by 0.122 and 0.145, and decreases the average accuracy<sub>a</sub> by 21.6% and 22.8%, for LIME and KSHAP, respectively. We do paired t-tests for each setup that only differs on whether to apply ConLUX, to show the statistical significance of the improvement. The p-value is all less than 0.01, which indicates with over 99% confidence the improvement is significant.

We also compared ConLUX unified explanations with two state-of-the-art concept-based task-specific explanation techniques: TBM for text tasks and EAC for image tasks. Table 3 shows the fidelity of TBM, EAC, and ConLUX unified explanations. ConLUX helps two classic local model-agnostic techniques to achieve 5.75% and 4.9% more accuracy than TBM and EAC.

## 5 RELATED WORK

Our work is related to model-agnostic explanation techniques and concept-based explanation techniques.

Model-agnostic explanation techniques consider target models as black boxes and provide explanations without requiring any knowledge of the model’s internal structure. Existing Model-agnostic explanation techniques provide different types of explanations, such as feature importance (Lund-

Table 2: Average AOPC and accuracy<sub>a</sub> (higher AOPC and lower accuracy<sub>a</sub> are better) of LIME, KSHAP, and their ConLUX-augmented versions (denoted as LIME\* and KSHAP\*) on two sentiment analysis models and three image classification models.

Models	AOPC $\uparrow$				Accuracy <sub>a</sub> (%) $\downarrow$			
	LIME	LIME*	KSAHP	KSAHP*	LIME	LIME*	KSAHP	KSAHP*
Llama 3.1	–	–	–	–	82.3	<b>49.8</b>	72.1	<b>45.4</b>
BERT	0.243	<b>0.456</b>	0.379	<b>0.553</b>	75.7	<b>47.7</b>	60.3	<b>40.2</b>
YOLOv8	0.401	<b>0.474</b>	0.433	<b>0.590</b>	14.8	<b>5.0</b>	33.1	<b>6.9</b>
ViT	0.469	<b>0.598</b>	0.454	<b>0.611</b>	21.1	<b>4.9</b>	25.7	<b>9.0</b>
ResNet-50	0.232	<b>0.306</b>	0.233	<b>0.323</b>	36.3	<b>14.6</b>	33.1	<b>8.6</b>

Table 3: Average accuracy (higher accuracy is better) of TBM, EAC, and ConLUX unified explanations on two sentiment analysis models and three image classification models.

Methods	Accuracy (%) $\uparrow$				
	Llama 3.1	BERT	YOLOv8	ViT	ResNet-50
TBM	89.6	81.4	–	–	–
EAC	–	–	56.6	53.4	57.7
ConLUX	<b>94.7</b>	<b>87.8</b>	<b>61.3</b>	<b>59.6</b>	<b>61.5</b>

berg & Lee, 2017; Ribeiro et al., 2016; Tan et al., 2023; Shankaranarayana & Runje, 2019), decision rules (Ribeiro et al., 2018; Guidotti et al., 2018; Dhurandhar et al., 2018), counterfactuals (Wachter et al., 2018; Guidotti et al., 2018), and visualizations (Goldstein et al., 2015; Friedman, 2001; Apley & Zhu, 2020). However, to our knowledge, all existing model-agnostic explanation techniques provide explanations at feature levels (Zhang et al., 2021c). Basic feature-based explanations are usually worse in aligning with either the decision-making process of the model or end-users (Ghorbani et al., 2019a; Sun et al., 2023; Kim et al., 2018), which makes these explanations unfaithful and hard to understand.

Concept-based explanation techniques provide explanations in terms of high-level concepts, which align with the decision-making process of the model better and are more interpretable to end-users. To our knowledge, existing concept-based explanation techniques are all model-specific or task-specific (Poeta et al., 2023b). We categorize them into three groups: (1) techniques that extract concepts from the model’s internal structure (Zhang et al., 2021b; Yeh et al., 2020; 2019b; Cunningham et al., 2023; Ghorbani et al., 2019b; Crabbé & van der Schaar, 2022; Fel et al., 2023), which are limited to specific types of models, (2) techniques that use external knowledge to define concepts (El Shawi, 2024; Widmer et al., 2022), which are limited to specific types of tasks since their methods based on the knowledge for a specific task, and (3) techniques that use pre-trained models to extract concepts (Ludan et al., 2023; Sun et al., 2023). Ludan et al. (2023) propose TBM, which is a surrogate model specifically designed for text data, while Sun et al. (2023) propose EAC, which also utilizes internal information of the target model. Therefore, these techniques are only for specific types of tasks. In addition, these techniques mainly focus only on attributions which limits their use cases (Poeta et al., 2023b).

## 6 CONCLUSION

We have proposed ConLUX, a general framework that automatically extracts high-level concepts and incorporates them into existing local model-agnostic explanation techniques to provide concept-based explanations, which are more faithful and easier to understand by end-users. ConLUX offers unified explanations that combine attributions, sufficient conditions, and counterfactuals. This satisfies diverse user needs and fills the current gap in concept-based explanations, which lack forms beyond attributions. ConLUX achieves this by utilizing large pre-trained models to extract high-level concepts, elevating language components from feature level to concept level, and extending perturbation models to sample in the concept space. We have instantiated ConLUX on Anchors, LIME, LORE, and Kernel SHAP, and provide unified explanations. We have constructed empirical evaluations to demonstrate the effectiveness of ConLUX.

## REFERENCES

- 540  
541  
542 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-  
543 man, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical  
544 report. *arXiv preprint arXiv:2303.08774*, 2023.
- 545 Daniel W Apley and Jingyu Zhu. Visualizing the effects of predictor variables in black box super-  
546 vised learning models. *Journal of the Royal Statistical Society Series B: Statistical Methodology*,  
547 82(4):1059–1086, 2020.
- 548  
549 Aparna Balagopalan, Haoran Zhang, Kimia Hamidieh, Thomas Hartvigsen, Frank Rudzicz, and  
550 Marzyeh Ghassemi. The road to explainability is paved with bias: Measuring the fairness of  
551 explanations. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and*  
552 *Transparency*, pp. 1194–1206, 2022.
- 553 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhari-  
554 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal,  
555 Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M.  
556 Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz  
557 Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec  
558 Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020. URL  
559 <https://arxiv.org/abs/2005.14165>.
- 560 Jonathan Crabbé and Mihaela van der Schaar. Concept activation regions: A generalized framework  
561 for concept-based explanations. (arXiv:2209.11222), September 2022. doi: 10.48550/arXiv.2209.  
562 11222. URL <http://arxiv.org/abs/2209.11222>. arXiv:2209.11222 [cs].
- 563  
564 Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert Huben, and Lee Sharkey. Sparse au-  
565 toencoders find highly interpretable features in language models. (arXiv:2309.08600), October  
566 2023. doi: 10.48550/arXiv.2309.08600. URL <http://arxiv.org/abs/2309.08600>.  
567 arXiv:2309.08600 [cs].
- 568  
569 Timothée Darcet, Maxime Oquab, Julien Mairal, and Piotr Bojanowski. Vision transformers need  
570 registers, 2023.
- 571 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hi-  
572 erarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,  
573 pp. 248–255. Ieee, 2009.
- 574  
575 Amit Dhurandhar, Pin-Yu Chen, Ronny Luss, Chun-Chen Tu, Paishun Ting, Karthikeyan Shan-  
576 mugam, and Payel Das. Explanations based on the missing: Towards contrastive expla-  
577 nations with pertinent negatives. *Advances in neural information processing systems*, 31,  
578 2018. URL [https://proceedings.neurips.cc/paper\\_files/paper/2018/  
hash/c5ff2543b53f4cc0ad3819a36752467b-Abstract.html](https://proceedings.neurips.cc/paper_files/paper/2018/hash/c5ff2543b53f4cc0ad3819a36752467b-Abstract.html).
- 579  
580 Radwa El Shawi. Conceptglassbox: Guided concept-based explanation for deep neural networks.  
581 *Cognitive Computation*, 16(5):2660–2673, September 2024. ISSN 1866-9964. doi: 10.1007/  
582 s12559-024-10262-8.
- 583  
584 Abhimanyu Dubey et al. The llama 3 herd of models, 2024a. URL [https://arxiv.org/abs/  
2407.21783](https://arxiv.org/abs/2407.21783).
- 585  
586 Gemini Team et al. Gemini: A family of highly capable multimodal models, 2024b. URL <https://arxiv.org/abs/2312.11805>.
- 587  
588 Thomas Fel, Victor Boutin, Louis Béthune, Remi Cadene, Mazda Moayeri, Léo Andéol, Math-  
589 ieu Chalvidal, and Thomas Serre. A holistic approach to unifying automatic concept extraction  
590 and concept importance estimation. *Advances in Neural Information Processing Systems*, 36:  
591 54805–54818, December 2023.
- 592  
593 Jerome H Friedman. Greedy function approximation: a gradient boosting machine. *Annals of*  
*statistics*, pp. 1189–1232, 2001.

- 594 Amirata Ghorbani, James Wexler, James Zou, and Been Kim. Towards automatic concept-based  
595 explanations. (arXiv:1902.03129), October 2019a. doi: 10.48550/arXiv.1902.03129. URL  
596 <http://arxiv.org/abs/1902.03129>. arXiv:1902.03129 [cs, stat].  
597
- 598 Amirata Ghorbani, James Wexler, James Y Zou, and Been Kim. Towards automatic concept-based  
599 explanations. In *Advances in Neural Information Processing Systems*, pp. 9273–9282, 2019b.  
600
- 601 Alex Goldstein, Adam Kapelner, Justin Bleich, and Emil Pitkin. Peeking inside the black box:  
602 Visualizing statistical learning with plots of individual conditional expectation. *Journal of Com-  
603 putational and Graphical Statistics*, 24(1):44–65, 2015. doi: 10.1080/10618600.2014.907095.
- 604 Riccardo Guidotti, Anna Monreale, Salvatore Ruggieri, Dino Pedreschi, Franco Turini, and Fosca  
605 Giannotti. Local rule-based explanations of black box decision systems. *CoRR*, abs/1805.10820,  
606 2018. URL <http://arxiv.org/abs/1805.10820>.
- 607 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-  
608 nition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.  
609 770–778, 2016.  
610
- 611 Aya Abdelsalam Ismail, Hector Corrada Bravo, and Soheil Feizi. Improving deep learning inter-  
612 pretability by saliency guided training. *Advances in Neural Information Processing Systems*, 34:  
613 26726–26739, 2021.  
614
- 615 Heinrich Jiang, Jennifer Jang, and Samory Kpotufe. Quickshift++: Provably good initializations for  
616 sample-based mean shift, 2018. URL <https://arxiv.org/abs/1805.07909>.
- 617 Glenn Jocher, Ayush Chaurasia, and Jing Qiu. Ultralytics yolov8, 2023. URL <https://github.com/ultralytics/ultralytics>.
- 618  
619
- 620 Emilie Kaufmann and Shivaram Kalyanakrishnan. Information complexity in bandit subset selec-  
621 tion. In *Conference on Learning Theory*, pp. 228–251. PMLR, 2013.  
622
- 623 Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler, and Fernanda Viegas.  
624 Interpretability beyond feature attribution: Quantitative testing with concept activation vectors  
625 (tcav). In *International conference on machine learning*, pp. 2668–2677. PMLR, 2018. URL  
626 <http://proceedings.mlr.press/v80/kim18d.html>.
- 627 Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete  
628 Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick.  
629 Segment anything. *arXiv:2304.02643*, 2023.  
630
- 631 Junhao Liu and Xin Zhang. Rex: A framework for incorporating temporal information in model-  
632 agnostic local explanation techniques, 2023. URL <https://arxiv.org/abs/2209.03798>.  
633
- 634 Josh Magnus Ludan, Qing Lyu, Yue Yang, Liam Dugan, Mark Yatskar, and Chris Callison-  
635 Burch. Interpretable-by-design text classification with iteratively generated concept bottle-  
636 neck. (arXiv:2310.19660), October 2023. doi: 10.48550/arXiv.2310.19660. URL <http://arxiv.org/abs/2310.19660>. arXiv:2310.19660 [cs].  
637  
638
- 639 Scott M. Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In Is-  
640 abelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N.  
641 Vishwanathan, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems  
642 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017,  
643 Long Beach, CA, USA*, pp. 4765–4774, 2017.
- 644 Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher  
645 Potts. Learning word vectors for sentiment analysis. In *Proceedings of the 49th Annual Meeting  
646 of the Association for Computational Linguistics: Human Language Technologies*, pp. 142–150,  
647 Portland, Oregon, USA, June 2011. Association for Computational Linguistics. URL <http://www.aclweb.org/anthology/P11-1015>.

- 648 Ali Modarressi, Mohsen Fayyaz, Ehsan Aghazadeh, Yadollah Yaghoobzadeh, and Moham-  
649 mad Taher Pilehvar. Decompx: Explaining transformers decisions by propagating token de-  
650 composition. In *Proceedings of the 61st Annual Meeting of the Association for Computational*  
651 *Linguistics (Volume 1: Long Papers)*, pp. 2649–2664, Toronto, Canada, July 2023. Association  
652 for Computational Linguistics. URL [https://aclanthology.org/2023.acl-long.](https://aclanthology.org/2023.acl-long.149)  
653 149.
- 654 Christoph Molnar. *Interpretable machine learning*. Lulu. com, 2020.
- 655 John Morris, Eli Lifland, Jin Yong Yoo, Jake Grigsby, Di Jin, and Yanjun Qi. Textattack: A frame-  
656 work for adversarial attacks, data augmentation, and adversarial training in nlp. In *Proceedings*  
657 *of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demon-*  
658 *strations*, pp. 119–126, 2020.
- 660 Maxime Oquab, Timothée Darcet, Theo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov,  
661 Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Russell Howes, Po-Yao  
662 Huang, Hu Xu, Vasu Sharma, Shang-Wen Li, Wojciech Galuba, Mike Rabbat, Mido Assran,  
663 Nicolas Ballas, Gabriel Synnaeve, Ishan Misra, Herve Jegou, Julien Mairal, Patrick Labatut, Ar-  
664 mand Joulin, and Piotr Bojanowski. Dinv2: Learning robust visual features without supervision,  
665 2023.
- 666 Eleonora Poeta, Gabriele Ciravegna, Eliana Pastor, Tania Cerquitelli, and Elena Baralis. Concept-  
667 based explainable artificial intelligence: A survey. (arXiv:2312.12936), December 2023a. URL  
668 <http://arxiv.org/abs/2312.12936>. arXiv:2312.12936 [cs].
- 669 Eleonora Poeta, Gabriele Ciravegna, Eliana Pastor, Tania Cerquitelli, and Elena Baralis. Concept-  
670 based explainable artificial intelligence: A survey. (arXiv:2312.12936), December 2023b. URL  
671 <http://arxiv.org/abs/2312.12936>. arXiv:2312.12936 [cs].
- 672 Marco Túlio Ribeiro, Sameer Singh, and Carlos Guestrin. “why should I trust you?”: Explaining  
673 the predictions of any classifier. In Balaji Krishnapuram, Mohak Shah, Alexander J. Smola,  
674 Charu C. Aggarwal, Dou Shen, and Rajeev Rastogi (eds.), *Proceedings of the 22nd ACM SIGKDD*  
675 *International Conference on Knowledge Discovery and Data Mining, San Francisco, CA, USA,*  
676 *August 13-17, 2016*, pp. 1135–1144. ACM, 2016.
- 677 Marco Túlio Ribeiro, Sameer Singh, and Carlos Guestrin. Anchors: High-precision model-agnostic  
678 explanations. In Sheila A. McIlraith and Kilian Q. Weinberger (eds.), *Proceedings of the Thirty-*  
679 *Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications*  
680 *of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in*  
681 *Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018*, pp. 1527–  
682 1535. AAAI Press, 2018.
- 683 Salvatore Ruggieri. Yadt: Yet another decision tree builder. In *16th IEEE International Conference*  
684 *on Tools with Artificial Intelligence*, pp. 260–265. IEEE, 2004. URL [https://ieeexplore.](https://ieeexplore.ieee.org/abstract/document/1374196/)  
685 [ieee.org/abstract/document/1374196/](https://ieeexplore.ieee.org/abstract/document/1374196/).
- 686 Wojciech Samek, Alexander Binder, Grégoire Montavon, Sebastian Lapuschkin, and Klaus-Robert  
687 Müller. Evaluating the visualization of what a deep neural network has learned. *IEEE transactions*  
688 *on neural networks and learning systems*, 28(11):2660–2673, 2016.
- 689 Sharath M. Shankaranarayana and Davor Runje. Alime: Autoencoder based approach for local  
690 interpretability. (arXiv:1909.02437), September 2019. doi: 10.48550/arXiv.1909.02437. URL  
691 <http://arxiv.org/abs/1909.02437>. arXiv:1909.02437 [cs, stat].
- 692 Ao Sun, Pingchuan Ma, Yuanyuan Yuan, and Shuai Wang. Explain any concept: Segment anything  
693 meets concept-based explanation. (arXiv:2305.10289), May 2023. doi: 10.48550/arXiv.2305.  
694 10289. URL <http://arxiv.org/abs/2305.10289>. arXiv:2305.10289 [cs].
- 695 Zeren Tan, Yang Tian, and Jian Li. Glime: General, stable and local lime explanation, 2023.
- 696 Sandra Wachter, Brent Mittelstadt, and Chris Russell. Counterfactual explanations without opening  
697 the black box: Automated decisions and the gdpr, 2018. URL [https://arxiv.org/abs/](https://arxiv.org/abs/1711.00399)  
698 1711.00399.

- 702 Yifei Wang. A comparative analysis of model agnostic techniques for explainable artificial intelli-  
703 gence. *Research Reports on Computer Science*, pp. 25–33, August 2024. ISSN 2811-0013. doi:  
704 10.37256/racs.3220244750.
- 705  
706 Cara Widmer, Md Kamruzzaman Sarker, Srikanth Nadella, Joshua Fiechter, Ion Jovina, Brandon  
707 Minnery, Pascal Hitzler, Joshua Schwartz, and Michael Raymer. Towards human-compatible xai:  
708 Explaining data differentials with concept induction over background knowledge, 2022. URL  
709 <https://arxiv.org/abs/2209.13710>.
- 710 Chih-Kuan Yeh, Cheng-Yu Hsieh, Arun Suggala, David I Inouye, and Pradeep K Ravikumar. On the  
711 (in) fidelity and sensitivity of explanations. *Advances in Neural Information Processing Systems*,  
712 32, 2019a.
- 713 Chih-Kuan Yeh, Been Kim, Sercan Arik, Chun-Liang Li, Pradeep Ravikumar, and Tomas Pfister.  
714 On concept-based explanations in deep neural networks. September 2019b. URL [https://](https://openreview.net/forum?id=BylWYC4KwH)  
715 [openreview.net/forum?id=BylWYC4KwH](https://openreview.net/forum?id=BylWYC4KwH).
- 716  
717 Chih-Kuan Yeh, Been Kim, Sercan Arik, Chun-Liang Li, Tomas Pfister, and Pradeep Ravikumar.  
718 On completeness-aware concept-based explanations in deep neural networks. *Advances in neural*  
719 *information processing systems*, 33:20554–20565, 2020.
- 720 Ruihan Zhang, Prashan Madumal, Tim Miller, Krista A. Ehinger, and Benjamin I. P. Rubins-  
721 stein. Invertible concept-based explanations for cnn models with non-negative concept acti-  
722 vation vectors. (arXiv:2006.15417), June 2021a. doi: 10.48550/arXiv.2006.15417. URL  
723 <http://arxiv.org/abs/2006.15417>. arXiv:2006.15417 [cs].
- 724  
725 Ruihan Zhang, Prashan Madumal, Tim Miller, Krista A. Ehinger, and Benjamin I. P. Rubinstein. In-  
726 vertible concept-based explanations for cnn models with non-negative concept activation vectors,  
727 2021b. URL <https://arxiv.org/abs/2006.15417>.
- 728 Yu Zhang, Peter Tiño, Aleš Leonardis, and Ke Tang. A survey on neural network interpretability.  
729 *IEEE Transactions on Emerging Topics in Computational Intelligence*, 5(5):726–742, October  
730 2021c. ISSN 2471-285X. doi: 10.1109/TETCI.2021.3100641. arXiv:2012.14261 [cs].  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

## 756 A THE CONLUX FRAMEWORK (CONTINUED)

757  
758 In this section, we introduce the details of incorporating ConLUX into existing local explanation  
759 methods.

760 We first follow Section 3 to introduce how we extend each part for text models in detail.

### 762 A.1 PRODUCING CONCEPT

763  
764 ConLUX provides predicates that describe high-level concepts by utilizing a large pre-trained model.  
765 Here, we describe the step to extract concept from text and image data in detail.

766  
767 **Text Data.** We use GPT-3.5 to produce concept-level predicates in two step. First, we let *GPT-3.5*  
768 generate the concepts that are important to the current task with a prompt as follows:

769       Now you are an expert at writing movie reviews, please tell me from which per-  
770       spectives can you evaluate a movie.

771  
772 Then we let the *GPT-3.5* to refine the predicates based on the current input and its similar  
773 sentences, and format the concepts. Here, we referred to the format defined in TBM (Ludan et al.,  
774 2023). The prompt is as follows:

775       Here we are presented with a text dataset accompanied by labels, and our objective  
776       is to identify a concept in the text that correlates with these labels. The task is to  
777       ....., we have known the following concepts are important in this task. [concepts]

778       In additionally, you should refine the concept to make sure that concepts can be  
779       used to correctly classify the following examples: [texts labels]

780       Then you are given examples of concepts across various datasets. please give me  
781       the concepts following their format:

782       Example 1:

783       "Concept Name": "explicit language",

784       "Concept Description": "'Explicit language' refers to the use of words, phrases,  
785       or expressions that are offensive, vulgar, or inappropriate for general audiences.  
786       This may include profanity, obscenities, slurs, sexually explicit or lewd language,  
787       and derogatory or discriminatory terms targeted at specific groups or individuals.",

788       "Concept Question": "What is the nature of the language used in the text?",

789       "Possible Responses": ["explicit", "strong", "non-explicit", "uncertain"],

790       "Response Guide":

791       "explicit": "The text contains explicit language, such as profanity, obscenities,  
792       slurs, sexually explicit or lewd language, or derogatory terms targeted at specific  
793       groups or individuals.",

794       "Strong": "The text contains strong language but not explicit language, it may  
795       contain terms that some viewers might find mature.",

796       "non-explicit": "The text is free from explicit language and is appropriate for  
797       general audiences.",

798       "uncertain": "It is difficult to determine the nature of the language used in the text  
799       or if any explicit terms are used.",

800       "Response Mapping":

801       "explicit": 2,

802       "strong": 1,

803       "non-explicit": -1,

804       "uncertain": 0

805       Example 2:

806       .....

807       Now, please give me your formatted concepts:

808  
809 Then we use the Response Guides to produce local concepts.

## 810 A.2 CONCEPT-FEATURE MAPPING

811 We the the following prompt for predicate-feature mapping:

812 Generate a sentence similar to a given sentence from the domain of dataset. The  
813 dataset’s description is that .

814 The generated sentence satisfies given concepts. Before generating the sentence,  
815 carefully read the description of each concept to understand the properties that  
816 the generated sentence must satisfy, think about how the sentence satisfies the  
817 concepts first, and then create the sentence that satisfies the concepts.

818 Format your response as a JSON with string keys and string values. Below is an  
819 example of a valid JSON response. The JSON contains keys thoughts, and the  
820 answer. End your response with ###

821 — Concepts: 1. Concept 1 2. Concept 2 ...

822 Response JSON: "thoughts": "In this section, you explain which snippets in your  
823 text support the concepts. Be as objective as possible and ignore irrelevant infor-  
824 mation. Focus only on the snippets and avoid making guesses.", "answer": "A  
825 sentence that satisfies the concepts." ###

826 Two examples of this task being performed can be seen below. Note that the  
827 answer should be in 5 to 20 words and should be a single sentence.

828 Example 1:

829 Concepts: 1. The plot of the text is exciting, captivating, or engrossing. It may  
830 have unexpected twists, compelling conflicts, or keep the reader eagerly turning  
831 pages. 2. The characters in the movie are portrayed in a realistic and convinc-  
832 ing manner. Their actions, dialogue, emotions, motivations, and development feel  
833 authentic and relatable, making them believable to the audience. 3. The narra-  
834 tive structure of the text is confusing or unclear, making it difficult to follow or  
835 comprehend the events happening within the story. 4. The text introduces some  
836 original elements or takes minor risks in the plot development, but overall, it lacks  
837 a truly unique or innovative narrative.

838 Response JSON: "thoughts": "The snippet 'the silly and crude storyline' men-  
839 tions a storyline that is described as silly and crude mentions a 'silly and crude  
840 storyline' which indicates a lack of creativity and reliance on clichéd plot devices,  
841 satisfying the concept of some originality and inventive plot development. The  
842 snippet 'the real issues tucked between the silly and crude storyline', mentions a  
843 contrast between real issues and a silly and crude storyline, indicating a potentially  
844 confusing narrative structure, satisfying the concept of a confusing narrative struc-  
845 ture. This also mentions that it has real conflict inside, which satisfies the concept  
846 of an exciting plot. This snippet also mentions mentions 'real issues' which indi-  
847 cates that the characters are portrayed in a realistic and convincing manner, sat-  
848 isfying the concept of realistic and convincing characters." "answer": "it's about  
849 issues most adults have to face in marriage and i think that's what i liked about it  
850 – the real issues tucked between the silly and crude storyline." ###

851 Example 2: .... ###

852 Perform the task below, keeping in mind to limit the response to 5 to 20 words and  
853 a single sentence. Return a valid JSON response ending with ###

854 Concepts:

855 Response JSON:

## 857 B EXPERIMENT SETTINGS (CONTINUED)

858 We experimented on two machines, one with an Intel i9-13900K CPU, 128 GiB RAM, and RTX  
859 4090 GPU, and another with Intel(R) Xeon(R) Silver 4314 CPU, 256GiB RAM, and 4 RTX 4090  
860 GPUs.

861 To measure the fidelity improvement brought by ConLUX, we keep all hyperparameters the same  
862 for both vanilla and augmented methods.  
863



864 For LIME and KSHAP, we set the number of sampled inputs to 1000 except for explaining Llama  
865 3.1.

866 For Anchors, we follow the default settings.

867 For LORE, we set  $ngen = 5$ .

868 For the LLama3.1 model, when applying it to the sentiment analysis task, we simply use the follow-  
869 ing prompt:  
870  
871

872 From now on, you should act as a sentiment analysis neural network. You should  
873 classify the sentiment of a sentence into positive or negative. If the sentence is  
874 positive, you should reply 1. Otherwise, if it's negative, you should reply 0. There  
875 may be some words that are masked in the sentence, which are represented by  
876 <UNK>. The input sentence may be empty, which is represented by <EMPTY>.  
877 You will be given the sentences to be classified, and you should reply with the  
878 sentiment of the sentence by 1 or 0.

879 There are two examples:

880 Sentence:

881 I am good

882 Sentiment:

883 1

884 Sentence:

885 The movie is bad.

886 Sentiment:

887 0

888 You must follow this format. Then I'll give you the sentence. Remember Your  
889 reply should be only 1 or 0. Do not contain any other content in your response.

890 The input sentence may be empty.

891 Sentence:

892 {The given sentence}

893 Sentiment:

## 894 C TEXT PERTURBATION FAITHFULNESS EXPERIMENT

895 As demonstrated by Ludan et al. (2023), large language models (LLMs) can verify whether an  
896 instance satisfies a given concept. Building on this, we conduct an experiment to evaluate the con-  
897 sistency of LLM-based perturbations. Specifically, we use LLMs to assess whether the applied  
898 perturbations successfully alter the intended concept. For each sentence, we generate 100 random  
899 perturbations and verify if the concepts in the generated sentences align with the expected changes.  
900 Our results indicate that Llama3.1, the large model employed in our fidelity experiments, achieves  
901 concept-level perturbation accuracy exceeding 99  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917