

# SToRI: Semantic Token Reweighting for Interpretable and Controllable Text Embeddings in Vision-Language Models

Anonymous ACL submission

## Abstract

A text encoder within Vision-Language Models (VLMs) plays a crucial role in translating textual input into an embedding space shared with images, thereby facilitating the interpretative analysis of vision tasks through natural language. Despite varying significance of different textual elements within a sentence, depending on the context or intended purpose, efforts to control the prominence of diverse textual information when constructing text embeddings have been lacking. This paper proposes a framework called Semantic Token Reweighting, aiming to incorporate Controllability while ensuring Interpretability of text embeddings (SToRI). SToRI refines the text encoding process in VLMs by differentially weighting semantic elements based on contextual importance, enabling finer control over emphasis responsive to user preferences and data-driven insights. The efficacy of SToRI is demonstrated through comprehensive experiments, showcasing its strength in image retrieval tailored to user preferences and its capability in few-shot image classification tasks.

## 1 Introduction

As artificial intelligence (AI) systems based on deep learning models grow in application in our daily lives, their black box nature raises issues of transparency, resulting in a demand for enhanced interpretability to promote trust in AI systems (Murdoch et al., 2019; Li et al., 2022). Consequently, research efforts have been focused on making the systems’ decision-making processes more human-understandable through various explanatory methods (Simonyan et al., 2014; Kim et al., 2018; Goyal et al., 2019; Wu and Mooney, 2019). Among the various forms of explanation, natural language has emerged as an excellent medium due to its human-friendly nature and adeptness in managing high-level abstractions (Kayser et al., 2021; Sammani

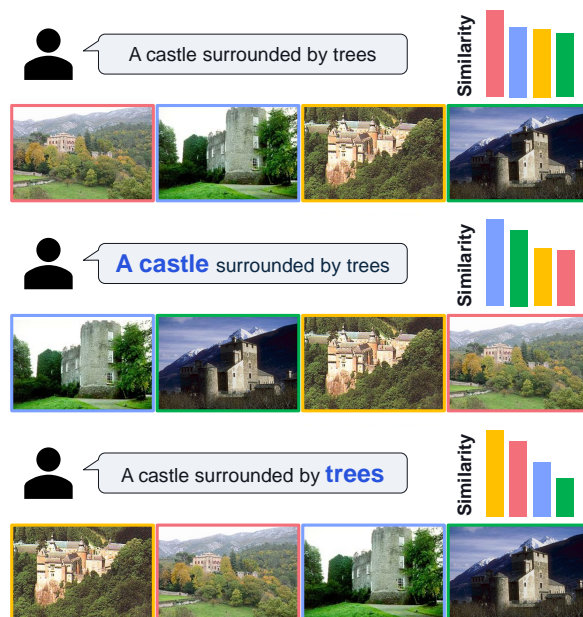


Figure 1: Examples of retrieved images through SToRI with a given text prompt. The order varies depending on the emphasized (blue) textual information.

et al., 2022). These advantages have led to a growing interest in research that utilizes natural language for interpretative analysis, extending even to domain of vision tasks (Hendricks et al., 2021; Yang et al., 2023). To facilitate the use of natural language in vision tasks, Vision-Language Models (VLMs) like CLIP (Radford et al., 2021) are commonly deployed to bridge visual information and its linguistic interpretation (Yuksekgonul et al., 2023; Yang et al., 2023; Oikarinen et al., 2023). Two encoders of VLMs translate an input image and text into image and text embeddings, respectively, which take vectorized forms and coexist in a shared embedding space.

Natural language sentences often carry multiple implications, with varying levels of significance that can change based on the desired outcome, even if the text remains unchanged. For instance, when

searching for images using the query ‘a castle surrounded by trees,’ a standard text query might bring up relevant images, but the preference on ‘trees’ relative to ‘a castle’ could differ based on user intent (see Figure 1). Texts rich in detail may benefit from selectively emphasizing certain information relevant to the task. However, existing methods lack the ability to fine-tune the importance given to specific pieces of information within text embeddings produced by the text encoders of VLMs. This paper endeavors to create text embeddings that can incorporate a varying controlled importance of each semantic element within a sentence.

To meet our objective, we introduce a novel framework termed **Semantic Token Reweighting** for Interpretable and Controllable textual representation (SToRI), which refines the focus on individual semantic components during text embedding extraction in VLMs. Each semantic element is assigned a numerical weight, denoting its significance, and these weights modulate the self-attention mechanism in text encoding. The proposed method makes it possible for the final text embedding vector to naturally include the desired emphasis on specific semantic elements, allowing for controllability. Moreover, the emphasis on particular semantic meanings remains within the realm of interpretability. SToRI efficiently produces text embeddings that reflect the desired focus without necessitating the training of new modules.

Our framework enables text embeddings to be tailored in two ways: user-driven and data-driven. In the user-driven approach, individuals can set the weight for each semantic token, allowing them to emphasize the elements they consider most relevant and customize the model to fit their preferences, as shown in Figure 1. On the other hand, the data-driven method derives token weights from training on dataset, facilitating the creation of text embeddings that are optimized for specific tasks like image classification and offer interpretable insights into the classifiers derived from texts. These enhancements have been substantiated through evaluation across various image recognition tasks, including image retrieval and few-shot classification.

Our main contributions are outlined as follows:

- We propose a novel framework of semantic token reweighting, which differentiates the importance of textual information during the construction of text embeddings in VLMs.
- Our approach facilitates the customization

of emphasis on specific semantics, and we demonstrate its usefulness in image retrieval tasks with a new metric for controllability.

- We demonstrate that our methodology not only builds improved text classifier in few-shot learning tasks but also unlocks a new dimension of interpretability.

## 2 Preliminary: Text embeddings in CLIP

The text encoder of CLIP (Radford et al., 2021), which utilizes a transformer-based architecture, transforms a given text prompt into a single vector through the following process. Initially, a given text prompt is converted into a sequence of text tokens  $\{x_i\}_{i=1}^N$ , where  $N$  represents the number of the text tokens. Tokens indicating the start and end, [SOS] and [EOS] tokens, are appended at the beginning and the end of the sequence of tokens, resulting in the extended series  $\{x_i\}_{i=0}^{N+1}$ , with  $x_0$  and  $x_{N+1}$  representing the [SOS] and [EOS], respectively. Each text token is then converted into an embedded input token, and positional embedding is added, resulting in the input embedding for the first transformer block  $\{z_i^0\}_{i=0}^{N+1}$ . For the  $l$ -th block of the encoder, the input tokens can be represented as  $Z^{l-1} = [z_0^{l-1}, \dots, z_{N+1}^{l-1}]$ . The output tokens from the  $l$ -th block is given by:

$$Z^l = \text{Block}^l(Z^{l-1}), \quad (1)$$

where  $l \in [1, L]$  with the encoder consisting of  $L$  blocks. Each block contains a multi-head self-attention mechanism. First,  $Z^{l-1}$  is projected into the query  $Q$ , key  $K$ , and value  $V$ . Then, the attention process is performed as follows:

$$\begin{aligned} \text{Attention}(Q, K, V) &= AV, \\ \text{s.t. } A &= \text{softmax}(QK^T). \end{aligned} \quad (2)$$

Scaling and masking operations are omitted for simplicity. Through the attention mechanism, tokens influence each other, and the values of  $A$  represent the extent to which they influence one another (Vaswani et al., 2017). In general, the final output text embedding of the [EOS] token encapsulates the full semantic meaning of the text prompt. This embedding is compared with image embeddings to assess the degree of correspondence with images once it has been projected into a multi-modal embedding space.

A pre-trained CLIP model is commonly employed for image classification, where given an

image, it computes similarity scores with class names, which become logits. To adapt the model to a specific dataset, fine-tuning is performed by minimizing the cross-entropy loss as follows:

$$\mathcal{L} = L_{CE}(y, \text{sim}(\phi_T, \phi_I)/\tau), \quad (3)$$

where  $\phi_T$  and  $\phi_I$  represent output text and image embeddings from two encoders, respectively, and  $\tau$  is a temperature factor.

### 3 Method

Our goal is to adjust the importance of various textual elements while encoding a given text prompt into a single text embedding vector. To achieve our goal, we propose **Semantic Token Reweighting**, which involves adjusting the attention given to individual tokens within the text encoding, guided by their respective weights. First, in Section 3.1, we elaborate on the methodology underlying Semantic Token Reweighting. Subsequently, in Section 3.2, we introduce two control strategies that leverage this technique. Figure 2 presents an overview of our comprehensive framework.

#### 3.1 Semantic Token Reweighting

In natural language processing, a given text is tokenized prior to encoding, resulting in one or more tokens. Consequently, to emphasize or de-emphasize a particular semantic element, one must focus on the corresponding tokens. Henceforth, our discussion will center on the process of reweighting in terms of these tokens.

Given a sequence of text tokens  $\{x_i\}_{i=1}^N$ , we first define a sequence of weights  $\{w_i\}_{i=1}^N$ , where  $w_i$  is the level of significance of token  $x_i$ . Note that  $w_i = 1$  indicates a typical weight in common situations, where  $x_i$  is neither emphasized nor de-emphasized. Our goal is to modulate the impact each token has on the final output embedding of the text prompt. As elaborated in Section 2, tokens interact with each other through attention mechanisms. Each token generates its embedding by referencing other tokens, including itself, in proportion to the attention scores. Consequently, as the attention score of a specific token increases, its influence on the text embedding becomes more substantial. Therefore, we directly multiply the weights  $\{w_i\}_{i=1}^N$  to amplify original attention values proportionally. From Eq. (2), the weighted attention scores can be reformulated as follows:

$$\hat{a}_{m,n} = \frac{w_n \exp(q_m k_n^T)}{\sum_j w_j \exp(q_m k_j^T)}, \quad (4)$$

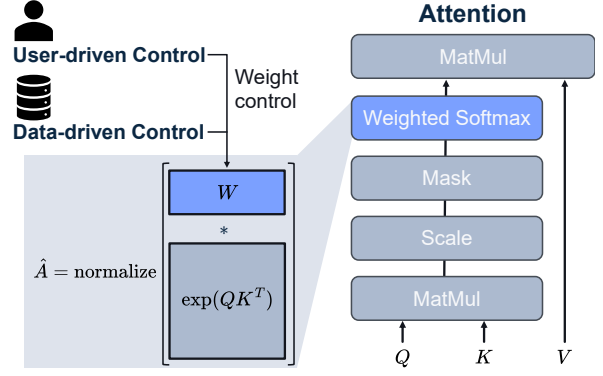


Figure 2: Overview of STORI. The weights can be determined through either user-driven or data-driven control. The weight vector is represented as  $W = [w_1, \dots, w_N]$ .

where  $\hat{a}_{m,n}$  represents attention value for  $n$ -th value token to be attended by  $m$ -th query token.  $q_m$  and  $k_n$  represent vector elements of  $Q$  and  $K$ , respectively. Through this process, we can selectively enhance the influence of particular tokens during the attention process by simply changing the corresponding weights.

The reweighting process is applied to all blocks following a certain block. Experimentally, we confirm that the effects are similar regardless of starting from any intermediate block. Please refer to Appendix B.3 for further details.

#### 3.2 Strategies to Control

There are two approaches to determine weights for tokens: user-driven and data-driven control.

**User-driven control** applies to scenarios where the user assigns weights to each token. This method allows user to determine a particular textual information to be emphasized or de-emphasized according to their intentions, thereby influencing the resulting text embeddings.

**Data-driven control** determines weights by learning from data. This approach is suitable when data is available and we want to obtain text embeddings that align closely with the data. An illustrative task where this can be effectively applied is image classification. In image classification, weights are trained using Eq. (3), where  $\phi_T$  is obtained with  $\hat{a}_{i,j}$ , allowing only  $\{w_i\}_{i=1}^N$  to be updated. Since the weights are trained to build text embeddings that correspond well to image belonging to their corresponding classes, we can interpret which textual information prominently stands out in the image data with the weights.

## 4 Experiments

We evaluate STORI under two scenarios: user-driven and data-driven controls. In the user-driven scenario, we demonstrate its application in preference-based image retrieval. In the data-driven scenario, we show its effectiveness in training an enhanced classifier for few-shot image classification and interpreting the classifier through its weights.

### 4.1 User-driven Control

To assess the effectiveness of STORI in emphasizing or de-emphasizing specific information based on applied weights, we compare the ordering of retrieved images using text embeddings.

#### 4.1.1 Experimental Setup

**Dataset.** The CelebA dataset (Liu et al., 2015) contains over 200K face images, each annotated with 40 attributes. Three attributes are chosen to create eight categories based on their presence or absence. Each category comprises 100 randomly selected images, resulting in a total of 800 images. For more details, please refer to Appendix A.1.

**Image Retrieval with Preference.** We construct a text prompt containing the selected attributes. For instance, the text prompt becomes ‘a photo of a woman with blonde hair, wearing eyeglasses’ for the attributes *female*, *blonde hair*, and *eyeglasses*. Using the text prompt and attribute weights, we obtain a corresponding text embedding through STORI, followed by sorting the images in descending order of similarity between their image embeddings and the text embedding.

**Model.** All experiments are conducted using CLIP ViT-L/14 (Radford et al., 2021), where reweighting is applied from the 7th block unless specified.

#### 4.1.2 Metric for Preference Retrieval

Our primary focus is on observing how adjusting weights for specific semantic elements affects the image retrieval order. To facilitate this comparison, we report the average precision score (AP) and precision (P@400) for images with the attributes influenced by the adjusted weights. For instance, when we modify the weight on ‘eyeglasses’, we consider images with eyeglasses as positive samples and calculate AP and P@400.

Additionally, we introduce a novel metric to quantify priority in preference retrieval. We generate a line plot illustrating the proportion of images retrieved for each attribute combination up to the

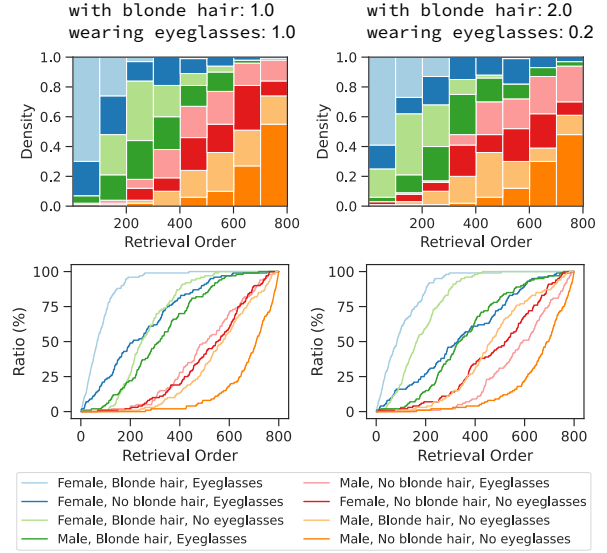


Figure 3: Results of preference retrieval using the text prompt ‘a photo of a woman with blonde hair, wearing eyeglasses’. The first row shows density plots with the retrieval order, and the second row visualizes the ratio of retrieved samples within each category. The left column shows results from a plain text prompt, whereas the right column depicts the results when the weights are adjusted. Best viewed in color.

	AP	P@400
Plain ( $w = 1.0$ )	$0.689 \pm 0.052$	$0.615 \pm 0.050$
Emphasized ( $w = 1.5$ )	$0.708 \pm 0.049$ $\Delta 0.020 \pm 0.011$	$0.630 \pm 0.048$ $\Delta 0.015 \pm 0.010$
De-emphasized ( $w = 0.5$ )	$0.652 \pm 0.063$ $\Delta -0.037 \pm 0.022$	$0.594 \pm 0.054$ $\Delta -0.021 \pm 0.013$

Table 1: Retrieval performance on attributes of the CelebA dataset. The results show mean values with standard deviation across multiple controlled attributes.

$n$ -th retrieved image (see Figure 4), and calculate the Area Under the Curve (AUC) for each plotted curve. A higher AUC value suggests a faster retrieval of associated visual attribute set, indicating a higher priority in the retrieval process.

#### 4.1.3 Results

Initially, we select three attributes, *female*, *blonde hair*, and *eyeglasses*, and observe the ordering of image retrieval as shown Figure 3. With the plain text embedding, the initial bin predominantly contains images featuring all selected attributes, followed by a prevalence of images from the ‘female, no blonde hair, eyeglasses’ category. When the weight on ‘with blonde hair’ increases and on ‘wearing eyeglasses’ decreases, images be-



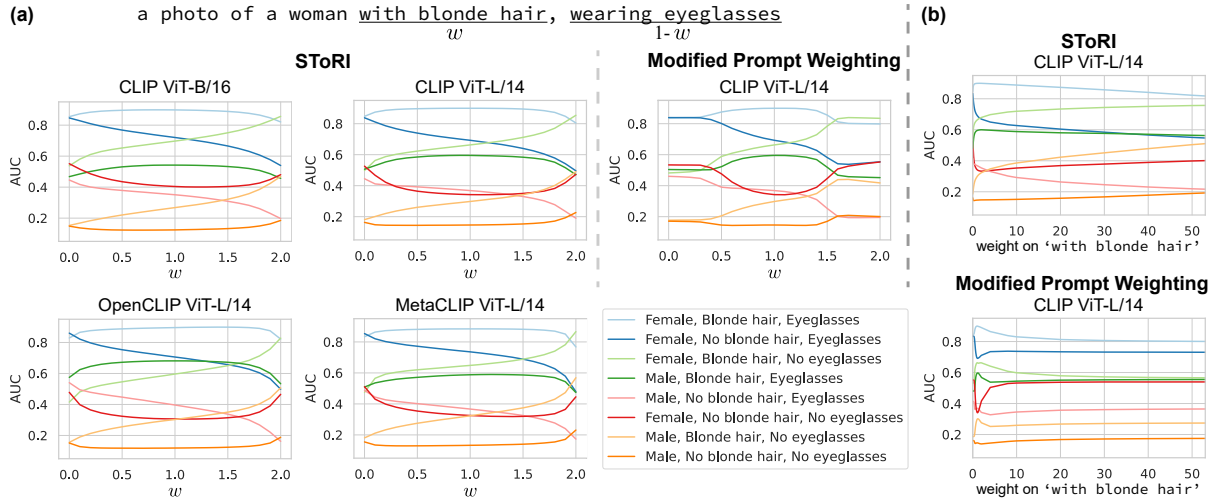


Figure 4: AUC scores from preference retrieval with varying weights. The text prompt is ‘a photo of a woman with blonde hair, wearing eyeglasses’. (a) The weights on ‘with blonde hair’ and ‘wearing eyeglasses’ are  $w$  and  $(1 - w)$ , respectively, which are adjusted simultaneously in opposite direction. (b) Only the weight on ‘with blonde hair’ is adjusted. Best viewed in color.

longing to ‘female, blonde hair, no eyeglasses’ are retrieved more prominently. This suggests that the ‘blonde hair’ gains more representation in the text embedding through reweighting. The groups with two or more mismatched attributes still rank lower, indicating that our method preserves the meanings of the original text while appropriately reflecting the intention of emphasis and de-emphasis.

We conduct quantitative validation across various text prompts. Table 1 presents AP and P@400 scores while controlling weights on attributes. We generate image pools and text prompts from three selected attributes. The reported scores are based on adjusting the weight for one specific attribute, considering the images containing that attribute as positive samples. Various combinations of attributes, totaling 20 text prompts, are used to obtain scores, and their averages and standard deviations are reported. Further details are in Appendix A.1. The results show that modifying the weight of tokens corresponding to a specific attribute in the text prompt results in faster retrieval of images with that attribute (both scores become higher) when the weight increases and slower retrieval when decreases (both scores become lower). This shows that adjusting the weight influences the creation of text embeddings, effectively highlighting or downplaying the corresponding attribute.

Figure 4(a) demonstrates the effects of weight control on the AUC scores for the retrieval of each category. As the weight assigned to the

‘with blonde hair’ increases and the weight for ‘wearing eyeglasses’ decreases, there is a noticeable rise in the AUC scores for the two categories that have blonde hair but no eyeglasses. In contrast, categories characterized by the absence of blonde hair and the presence of eyeglasses see a reduction in their AUC scores. When the weight assigned to ‘with blonde hair’ is set to zero, the differentiation between the ‘female, blonde hair, eyeglasses’ and ‘female, no blonde hair, eyeglasses’ categories is effectively eliminated, resulting in remarkably similar AUC scores. The effect of weight control is consistent across different CLIP models, such as CLIP ViT-B/16, CLIP ViT-L/14, OpenCLIP (Cherti et al., 2023), and MetaClip (Xu et al., 2023). This shows that SToRI enables the emphasis or de-emphasis of specific semantics within a text when constructing text embeddings across various models, showcasing its versatility.

#### 4.1.4 Comparison to Prompt weighting

We compare SToRI with prompt weighting, a technique often used in text-to-image generation via Stable Diffusion (Rombach et al., 2022). Prompt weighting multiplies weights by the difference in output token embeddings when provided with a text prompt versus an empty one. Unlike Stable Diffusion, which utilizes all output token embeddings, we aim to build a vector form of text embedding from [EOS] token. Therefore, we modify prompt weighting for use at an intermediate layer, which we refer to as modified prompt weighting.

	Method	Text	ImageNet	DTD	Flowers102	SUN397	Caltech101	Food101	AVG
1shot	TaskRes	Base	75.95±0.03	55.40±0.27	81.16±0.44	68.10±0.16	94.28±0.11	90.30±0.10	77.53
	TaskRes	Base+CuPL	74.69±0.04	65.66±0.82	90.07±0.79	73.52±0.49	95.89±0.57	90.35±0.36	81.70
	SToRI (Ours)	Base+CuPL	76.68±0.15	65.82±0.98	89.05±0.58	72.88±0.20	96.27±0.67	91.34±0.12	82.01
2shot	TaskRes	Base	76.03±0.00	55.52±0.48	81.50±0.62	69.53±0.14	94.54±0.05	90.49±0.05	77.93
	TaskRes	Base+CuPL	75.55±0.04	66.45±1.57	92.38±0.69	75.69±0.29	96.96±0.27	90.64±0.38	82.95
	SToRI (Ours)	Base+CuPL	77.36±0.23	66.37±1.01	91.56±0.60	75.75±0.04	97.15±0.13	91.49±0.24	83.28
4shot	TaskRes	Base	76.16±0.02	55.85±0.12	81.65±0.28	71.15±0.09	94.58±0.09	90.44±0.05	78.31
	TaskRes	Base+CuPL	76.42±0.03	70.76±1.12	93.22±0.37	77.20±0.08	97.40±0.21	91.45±0.15	84.41
	SToRI (Ours)	Base+CuPL	77.90±0.05	69.03±1.48	92.46±0.09	76.89±0.02	97.39±0.08	91.68±0.07	84.22
8shot	TaskRes	Base	76.87±0.05	58.14±0.07	86.82±0.19	74.52±0.07	96.17±0.08	91.12±0.07	80.60
	TaskRes	Base+CuPL	77.97±0.02	73.42±0.86	98.17±0.25	77.54±0.16	97.00±0.28	91.27±0.11	85.89
	SToRI (Ours)	Base+CuPL	78.38±0.13	72.03±0.60	97.51±0.43	78.34±0.13	96.98±0.29	90.50±0.05	85.62
16shot	TaskRes	Base	77.34±0.03	61.47±0.16	90.85±0.21	76.01±0.24	96.75±0.07	91.30±0.10	82.29
	TaskRes	Base+CuPL	79.18±0.10	77.05±0.65	99.07±0.11	78.98±0.10	97.65±0.23	91.49±0.08	87.24
	SToRI (Ours)	Base+CuPL	79.03±0.13	74.94±0.10	98.55±0.23	79.61±0.11	97.43±0.20	91.18±0.10	86.79

Table 2: Accuracy (%) on few-shot classification with CLIP ViT-L/14. The results include mean values with standard deviation across three runs. The results of TaskRes are reproduced.

As depicted in Figure 4(a), the modified prompt weighting influences the significance of tokens similarity to SToRI. However, the change in AUC is not gradual; it remains nearly static when weights fall below 0.5 or above 1.5. As shown in Figure 4(b), even when the weight for ‘with blonde hair’ increases significantly, SToRI consistently raises the AUC for the category ‘female, blonde hair, no eyeglasses’. In contrast, the AUC with modified prompt weighting initially increases but subsequently decreases, indicating augmented weight fails to heighten emphasis. This could stem from the scaling of intermediate embeddings which, when overextended, surpasses the scale that the text encoder is pre-trained to deal with, lessening the intended effect of emphasis. SToRI, on the other hand, adjusts normalized attention scores within the self-attention mechanism, ensuring that as weight escalates, the relevant tokens consistently obtain attention scores approaching 1, thus preserving the desired impact.

## 4.2 Data-driven Control

We train weights that best represent each dataset for the image classification task.

### 4.2.1 Experimental Setup

**Datasets.** We use various benchmarks for few-shot learning *i.e.*, ImageNet (Deng et al., 2009), DTD (Cimpoi et al., 2014), SUN397 (Xiao et al., 2010), Flowers102 (Nilsback and Zisserman, 2008), Caltech101 (Fei-Fei et al., 2004), and Food101 (Bossard et al., 2014). We use CUB (Wah et al., 2011) dataset for analysis on interpretation.

**Text Prompts.** We use text descriptions for each class which are provided by CuPL (Pratt et al., 2023). For the ImageNet and SUN397 datasets, due to the large number of total prompts, we use 10 text prompts for each class, selected based on their similarity with training set. We average the text embeddings from multiple text prompts to build one text embedding for each class. We refer the text embedding for image classifier as a text classifier.

**Model.** All experiments use CLIP ViT-L/14, with reweighting applied from the 7th block onward.

**Implementation Details.** We set the logarithm of the weight as the parameter to be trained in order to constrain the weights to non-negative values. Each text prompt has its own individual set of weights.

### 4.2.2 Few-shot Classification

**Experimental Details.** Following TaskRes (Yu et al., 2023), we evaluate our method by training with 1/2/4/8/16 examples (shots) per class from the training sets, respectively, and testing on the comprehensive test sets. For further details, please refer to Appendix A.2.

**Comparison.** To evaluate the capability of the text classifier obtained through SToRI to perform few-shot image classification, we conduct a comparative analysis of the prediction performance between SToRI and TaskRes (Yu et al., 2023). TaskRes is a recent method for few-shot image classification, which trains class-specific residual embeddings added to initial text embeddings to create new classifiers. Such residual embeddings exist in uninterpretable space, rendering the final classifier

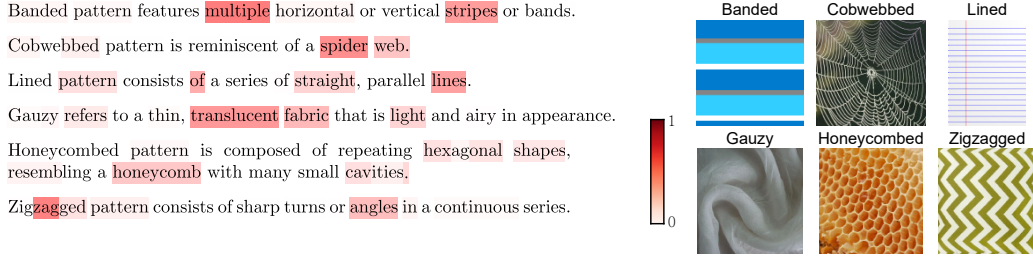


Figure 5: Text prompts and corresponding weights are provided as examples after training. The intensity of the red shading reflects the weight assigned, with darker shades indicating higher weights. For visualization, the weights are normalized to sum up 1. The figures on the right display an example image for each class.

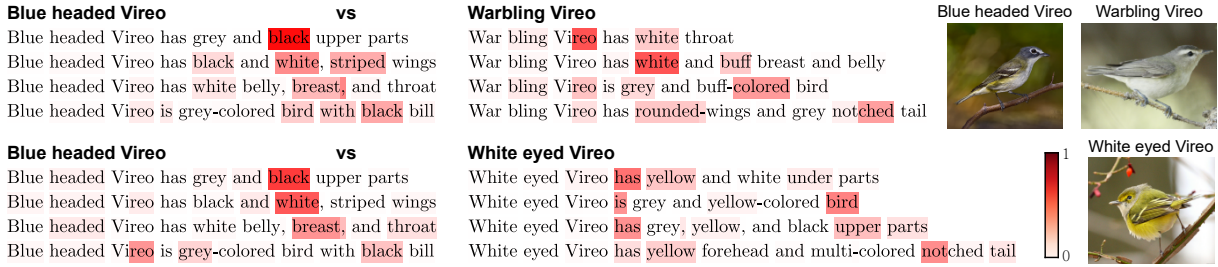


Figure 6: Text prompts and their corresponding weights are presented after training with the CUB dataset. The more intense the shade of red, the greater the weight assigned. In each scenario, the text classifier is trained to discriminate two classes. The weights for the same text prompts vary depending on the class to be distinguished.

also uninterpretable. In contrast, SToRI trains only weights, indicating the degree to which each semantic element within a given sentence should be emphasized, thus maintaining interpretability.

Ensuring interpretability, SToRI achieves performance comparable to TaskRes, as presented in Table 2. “Base” refers to custom text prompts including class names, which are generally used in few-shot image classification tasks with CLIP (Yu et al., 2023). We use both base and CuPL text prompts, with weights trained exclusively on CuPL. In the 1/2-shot setting, SToRI generally outperforms TaskRes across most datasets. In the 4/8/16-shot setting, it exhibits only a marginal difference, achieving nearly similar performance. This indicates that SToRI provides substantial flexibility to text embeddings, enabling it to be an enhanced text classifier that effectively represents image data.

### 4.2.3 Interpretability

**Interpretation with Trained Weights.** After training for an image classification task, we analyze the trained weights. Figure 5 presents examples of text prompts and the corresponding trained weights for each token within the DTD dataset. We have crafted the text prompts. We can discern that *banded* is associated with an emphasis on words

like multiple and stripes. For *gauzy*, terms such as translucent and light are emphasized, and *cobwebbed* are notably associated with the word spider web. As illustrated by the images corresponding to each category, high weight values are assigned to important semantic tokens. This shows that SToRI can learn text embeddings that effectively represent the data in a data-driven control context, and the trained weights can offer novel insights for interpretation.

**Does Optimization Occur in Interpretable Space?** To ensure interpretability of text embeddings through data-driven control optimization, we conduct two experiments: an analysis on trained classifiers with different class compositions and an assessment of the effect of nonsensical text tokens.

The role of classifier is to distinguish one class from others. Thus, even for classifiers within the same class, the critical distinguishing features can vary depending on the alternative categories being compared. Figure 6 shows two text classifiers trained on the CUB dataset for two distinct pairs: *Blue headed Vireo* versus *Warbling Vireo*, and *Blue headed Vireo* versus *White eyed Vireo*. The text prompts for each class are generated with the attribute labels from the dataset. When contrasting *Blue headed Vireo* with the *Warbling Vireo*,

Text	Caltech101	SUN397
CuPL	97.42±0.23	79.54±0.12
CuPL+Nonsensical tokens	97.30±0.15	79.11±0.10

Table 3: Accuracy (%) on 16-shot image classification.

striped is attributed a high weight. However, when distinguished from the *White eyed Vireo*, the weight on striped becomes low and grey is attributed a high weight. Note that *White eyed Vireo* also have striped wings. These terms highlight the prominent differences between each unique pairing of the classes.

Table 3 reports the 16-shot classification performance when nonsensical text tokens are added. We randomly sample five tokens from the set of three rare tokens (Ruiz et al., 2023), namely ‘sks’, ‘pll’, and ‘ucd’, and add them to the end of all the original texts from CuPL. The inclusion of rare tokens does not contribute meaningful information to build a text classifier; it simply extends the number of tokens and trainable parameters. As a result, the performance when rare tokens are added did not surpass that without their addition. This demonstrates that adoption of the tokens without semantic meaning does not contribute to performance improvement. These findings support that data-driven control, achieved through attention modulation for tokens with semantic meaning, facilitates the creation of text embeddings that effectively represent the data, thereby ensuring the interpretability of text embeddings.

## 5 Related Works

**VLMs and Interpretability.** In recent vision tasks, interpretative analysis in natural language becomes popular rather than relying solely on visual form. For this purpose, VLMs have commonly been employed to connect the image feature space with the text feature space used for explanation. Kim et al. (2023) utilized VLMs to get concept activation vector (Kim et al., 2018) in vision model. Yuksekgonul et al. (2023) and Oikarinen et al. (2023) leveraged VLMs to determine whether concepts defined in text are present in images. Menon and Vondrick (2023) formulated text prompts for image classes using Large Language Models and employed them for zero-shot classification with VLMs. These approaches simply utilize the shared embedding space of existing VLMs. In contrast, our method introduces a new dimension

of interpretability by providing controllability over the focus of textual information, thereby enhancing its interpretative utility.

**Few-shot Image Classification.** VLMs exhibit promising performance in image recognition tasks, leading to the development of various few-shot learning approaches. CoOp (Zhou et al., 2022b) and CoCoOp (Zhou et al., 2022a) are representative methods based on prompt tuning. Tip-Adapter (Zhang et al., 2022) integrates an extra adapter unit following the encoders. TaskRes (Yu et al., 2023) involves training task-specific residual text embeddings for each category. These approaches incorporate extra trainable parameters outside an interpretable framework, thereby not ensuring interpretability.

**Enrich Textual Representation.** In text-to-image generation, several approaches have been developed to enrich textual representation. Prompt weighting<sup>1</sup> is a common technique in Stable Diffusion (Rombach et al., 2022), which multiplies weights to individual output token embeddings prior to supplying them to the image generation model. Prompt-to-Prompt controls cross-attention between noise images and text embeddings (Hertz et al., 2022). Additionally, Ge et al. (2023) proposed a richer text editor that allows users to define various input conditions for image generation, such as coloring and footnotes. While prior works have focused on image generation, our work pioneers enriched textual representations for image recognition, utilizing novel single vector construction. This distinctive approach establishes a new avenue for incorporating linguistic context in visual understanding.

## 6 Conclusion

We introduce a novel framework that enables the reweighting of importance of semantic tokens in text embedding. This approach is a novel means of adapting the explanatory power of natural language in vision tasks. Our user-driven and data-driven controls empower users to dictate the emphasis on specific terms and facilitate the tuning of text embeddings for classification while ensuring interpretability. Our approach can be easily applied to any model based on attention, and has potential scalability in various vision tasks and multi-modal tasks, given the widespread use of VLMs.

<sup>1</sup>[https://huggingface.co/docs/diffusers/using-diffusers/weighted\\_prompts](https://huggingface.co/docs/diffusers/using-diffusers/weighted_prompts)



## 7 Limitations

Our method is focusing on controlling the attention of each semantic element within a given natural language sentence, rather than generating new textual information. Therefore, one of the limitations of our method is its dependence on the richness and quality of the given texts. For example, when using data to train a classifier, if the given text lacks sufficient rich information, adjusting the attention may not sufficiently enlarge the text embedding space. This difficulty in expanding the embedding space makes it challenging to establish a basis for improving classification performance and explaining data.

Additionally, we do not consider the inherent black box characteristics of VLMs. However, if this model has undergone sufficient testing and is deemed reliable, the advantage of our method lies in additional optimization and control being in a reliable and controllable space.

## 8 Ethics Statement

Our goal is to employ controllability when building text embeddings. This enables for users to emphasize or deemphasize a certain part of textual information and improving text embeddings for vision tasks, ensuring interpretability. We believe this work can be used to build trustful AI systems by providing natural language interpretation.

## References

- Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. 2014. Food-101—mining discriminative components with random forests. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part VI 13*, pages 446–461. Springer.
- Mehdi Cherti, Romain Beaumont, Ross Wightman, Mitchell Wortsman, Gabriel Ilharco, Cade Gordon, Christoph Schuhmann, Ludwig Schmidt, and Jenia Jitsev. 2023. Reproducible scaling laws for contrastive language-image learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2818–2829.
- Mircea Cimpoi, Subhansu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. 2014. Describing textures in the wild. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3606–3613.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee.
- Li Fei-Fei, Rob Fergus, and Pietro Perona. 2004. Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories. In *2004 conference on computer vision and pattern recognition workshop*, pages 178–178. IEEE.
- Songwei Ge, Taesung Park, Jun-Yan Zhu, and Jia-Bin Huang. 2023. Expressive text-to-image generation with rich text. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7545–7556.
- Yash Goyal, Ziyan Wu, Jan Ernst, Dhruv Batra, Devi Parikh, and Stefan Lee. 2019. Counterfactual visual explanations. In *International Conference on Machine Learning*, pages 2376–2384. PMLR.
- Lisa Anne Hendricks, Anna Rohrbach, Bernt Schiele, Trevor Darrell, and Zeynep Akata. 2021. Generating visual explanations with natural language. *Applied AI Letters*, 2(4):e55.
- Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. 2022. Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626*.
- Maxime Kayser, Oana-Maria Camburu, Leonard Salewski, Cornelius Emde, Virginie Do, Zeynep Akata, and Thomas Lukasiewicz. 2021. e-vil: A dataset and benchmark for natural language explanations in vision-language tasks. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1244–1254.
- Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler, Fernanda Viegas, et al. 2018. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav). In *International Conference on Machine Learning*, pages 2668–2677. PMLR.
- Siwon Kim, Jinoh Oh, Sungjin Lee, Seunghak Yu, Jaeyoung Do, and Tara Taghavi. 2023. Grounding counterfactual explanation of image classifiers to textual concept space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10942–10950.
- Xuhong Li, Haoyi Xiong, Xingjian Li, Xuanyu Wu, Xiao Zhang, Ji Liu, Jiang Bian, and Dejing Dou. 2022. Interpretable deep learning: Interpretation, interpretability, trustworthiness, and beyond. *Knowledge and Information Systems*, 64(12):3197–3234.
- Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. 2015. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*.
- Ilya Loshchilov and Frank Hutter. 2017. **SGDR: Stochastic gradient descent with warm restarts**. In *International Conference on Learning Representations*.

679	Sachit Menon and Carl Vondrick. 2023. <a href="#">Visual classification via description from large language models</a> . In <i>The Eleventh International Conference on Learning Representations</i> .	736
680		737
681		738
682		739
683	W James Murdoch, Chandan Singh, Karl Kumbier, Reza Abbasi-Asl, and Bin Yu. 2019. Definitions, methods, and applications in interpretable machine learning. <i>Proceedings of the National Academy of Sciences</i> , 116(44):22071–22080.	740
684		741
685		742
686		743
687		
688	Maria-Elena Nilsback and Andrew Zisserman. 2008. Automated flower classification over a large number of classes. In <i>2008 Sixth Indian conference on computer vision, graphics &amp; image processing</i> , pages 722–729. IEEE.	744
689		745
690		746
691		747
692		748
693		749
694	Tuomas Oikarinen, Subhro Das, Lam M. Nguyen, and Tsui-Wei Weng. 2023. <a href="#">Label-free concept bottleneck models</a> . In <i>The Eleventh International Conference on Learning Representations</i> .	750
695		751
696		752
697	Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. 2017. Automatic differentiation in pytorch.	753
698		754
699		755
700		
701	Sarah Pratt, Ian Covert, Rosanne Liu, and Ali Farhadi. 2023. What does a platypus look like? generating customized prompts for zero-shot image classification. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pages 15691–15701.	756
702		757
703		758
704		759
705		760
706		
707	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In <i>International conference on machine learning</i> , pages 8748–8763. PMLR.	761
708		762
709		763
710		764
711		765
712		766
713	Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 10684–10695.	767
714		768
715		769
716		770
717		771
718		772
719	Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. 2023. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 22500–22510.	773
720		774
721		775
722		776
723		777
724		778
725	Fawaz Sammani, Tanmoy Mukherjee, and Nikos Deligiannis. 2022. Nlx-gpt: A model for natural language explanations in vision and vision-language tasks. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 8322–8332.	779
726		780
727		781
728		782
729		783
730		784
731	Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. 2014. Deep inside convolutional networks: Visualising image classification models and saliency maps. In <i>International Conference on Learning Representations</i> .	785
732		786
733		787
734		788
735		789
		790
	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. <i>Advances in neural information processing systems</i> , 30.	
	Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. 2011. The caltech-ucsd birds-200-2011 dataset.	
	Jialin Wu and Raymond Mooney. 2019. <a href="#">Faithful multimodal explanation for visual question answering</a> . In <i>Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP</i> , pages 103–112, Florence, Italy. Association for Computational Linguistics.	
	Jianxiong Xiao, James Hays, Krista A Ehinger, Aude Oliva, and Antonio Torralba. 2010. Sun database: Large-scale scene recognition from abbey to zoo. In <i>2010 IEEE computer society conference on computer vision and pattern recognition</i> , pages 3485–3492. IEEE.	
	Hu Xu, Saining Xie, Xiaoqing Ellen Tan, Po-Yao Huang, Russell Howes, Vasu Sharma, Shang-Wen Li, Gargi Ghosh, Luke Zettlemoyer, and Christoph Feichtenhofer. 2023. Demystifying clip data. <i>arXiv preprint arXiv:2309.16671</i> .	
	Yue Yang, Artemis Panagopoulou, Shenghao Zhou, Daniel Jin, Chris Callison-Burch, and Mark Yatskar. 2023. Language in a bottle: Language model guided concept bottlenecks for interpretable image classification. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 19187–19197.	
	Tao Yu, Zhihe Lu, Xin Jin, Zhibo Chen, and Xinchao Wang. 2023. Task residual for tuning vision-language models. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 10899–10909.	
	Mert Yuksekgonul, Maggie Wang, and James Zou. 2023. <a href="#">Post-hoc concept bottleneck models</a> . In <i>The Eleventh International Conference on Learning Representations</i> .	
	Renrui Zhang, Wei Zhang, Rongyao Fang, Peng Gao, Kunchang Li, Jifeng Dai, Yu Qiao, and Hongsheng Li. 2022. Tip-adapter: Training-free adaption of clip for few-shot classification. In <i>European Conference on Computer Vision</i> , pages 493–510. Springer.	
	Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. 2022a. Conditional prompt learning for vision-language models. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 16816–16825.	
	Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. 2022b. Learning to prompt for vision-language models. <i>International Journal of Computer Vision</i> , 130(9):2337–2348.	

## A Experimental Details

### A.1 User-driven Control

We initially select 11 attributes with a zero-shot classification performance of AUROC 0.75 or higher with CLIP on test set. For zero-shot classification, we create text prompt for each attribute and calculate AUROC using the similarity between the test set images and the text prompt. For example, when evaluating the attribute *smiling*, we use the text prompt ‘a photo of a smiling person’. Among the identified 11 attributes, we create combinations of three attributes, each including either *female* or *male*. We filter out the combinations where all eight categories contain fewer than 100 images. We conduct image retrieval with total 20 numbers of text prompts based on the combinations of attributes, as shown in Table 5.

### A.2 Data-driven Control

We follow the data split outlined in CoOp (Zhou et al., 2022b), conducting tests on the official test set of each dataset and the validation set of the ImageNet dataset. We use Adam optimizer with the cosine learning rate scheduler (Loshchilov and Hutter, 2017) following the training scheme of TaskRes (Yu et al., 2023). The learning rate is set to  $1 \times 10^{-2}$  for the ImageNet and SUN397 datasets, 0.1 for the Food101 dataset and for 8/16-shot scenarios on the DTD and Flower102 datasets, and  $5 \times 10^{-2}$  for the other datasets. The weight decay is set to 0. When reproducing TaskRes, the learning rate is set to  $2 \times 10^{-5}$  for the ImageNet dataset and  $2 \times 10^{-4}$  for the other datasets. The weight decay is set to 0.005 and  $\alpha$  is set to 0.5. 1/2/4-shot training is done with 100 epoch and the other is done with 200 epoch for all datasets. All experiments are implemented using PyTorch (Paszke et al., 2017), and we use official code base released by Yu et al. (2023) to reproduce TaskRes.

We use all the datasets and models solely for academic research purposes and do not employ them for improper intentions.

## B Additional Experimental Results

### B.1 Additional Examples for Interpretation

Figures 7 and 8 present examples of text prompts and the corresponding trained weights for each token within the ImageNet and DTD datasets, respectively. Higher weights are assigned to word tokens that effectively represent images.

Method	Plain Text Embeddings	SToRI
Relative Run Time	1.00	1.02

Table 4: Relative computational cost

### B.2 Computational Cost

We calculate runtime for apply SToRI compared to plain text embeddings, as reported in Table 4. The experiment is done on RTX A5000 and the reported values are mean values from 28K runs. Since SToRI only multiplies predefined weights when calculating attention scores, the runtime does not significantly differ from that of plain text embeddings.

### B.3 Position for Reweighting

Figure 9(a) compares the changes in AUC scores when we start reweighting at various positions. The reweighting process is applied to all blocks following a specific block. There is not a significant difference when we initiate token reweighting at intermediate positions. However, when token reweighting is applied to all blocks (from 1st block), a sharp bend is observed at 0.1 when the weight decreases. This is unlike other cases, which show a smooth decrease or increase in all scenarios. It is presumed that this abrupt occurrence is due to tokens in the specified position being completely disregarded when the weight becomes 0, leading to sudden gaps in those areas.

Figure 9(b) illustrates that when reweighting is applied only within a single specific intermediate block, the effects of emphasis or de-emphasis are scarcely observed. This suggests that if reweighting is confined within a single intermediate block, its effects in the subsequent blocks are counteracted, indicating that it should be applied in the subsequent blocks to emphasize or de-emphasize semantic tokens.



Selected Attriutes	Text prompts
Female/Male, Smiling, Bangs	a photo of a smiling [woman/man] with bangs
Female/Male, Smiling, Blond Hair	a photo of a smiling [woman/man] with blond hair
Female/Male, Smiling, Gray Hair	a photo of a smiling [woman/man] with gray hair
Female/Male, Smiling, Wearing Hat	a photo of a smiling [woman/man] wearing hat
Female/Male, Smiling, Eyeglasses	a photo of a smiling [woman/man] wearing eyeglasses
Female/Male, Bangs, Wearing Hat	a photo of a [woman/man] with bangs, wearing hat
Female/Male, Bangs, Eyeglasses	a photo of a [woman/man] with bangs, wearing eyeglasses
Female/Male, Blond Hair, Eyeglasses	a photo of a [woman/man] with blond hair, wearing eyeglasses
Female/Male, Gray Hair, Eyeglasses	a photo of a [woman/man] with gray hair, wearing eyeglasses
Female/Male, Wearing Hat, Eyeglasses	a photo of a [woman/man] wearing hat and eyeglasses

Table 5: All combinations of attributes and corresponding text prompts.

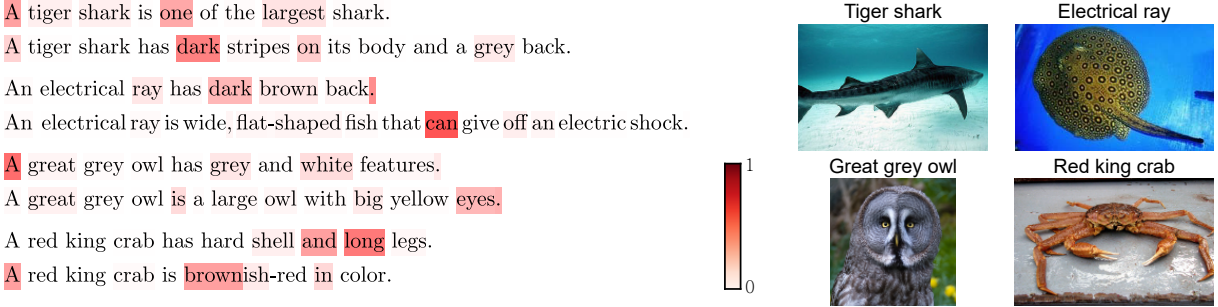


Figure 7: Text prompts and corresponding weights on the ImageNet dataset are provided as examples after training with data. For visualization, the weights are normalized to sum up 1. The figures on the right display an example image for each class.

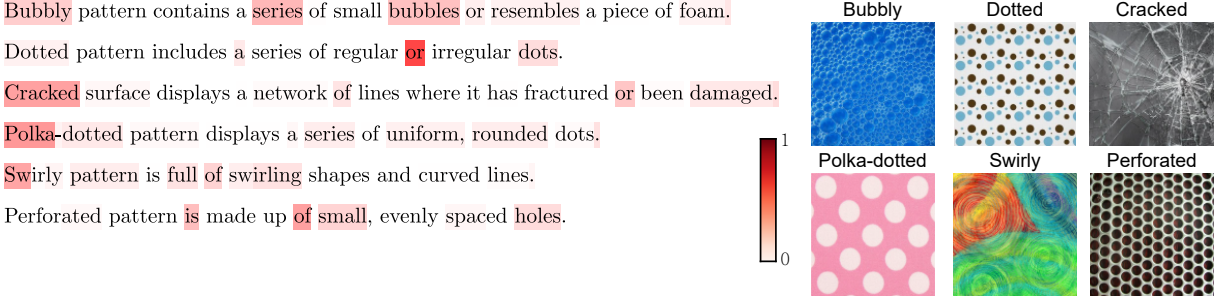


Figure 8: Text prompts and corresponding weights on the DTD dataset are provided as examples after training with data. For visualization, the weights are normalized to sum up 1. The figures on the right display an example image for each class.

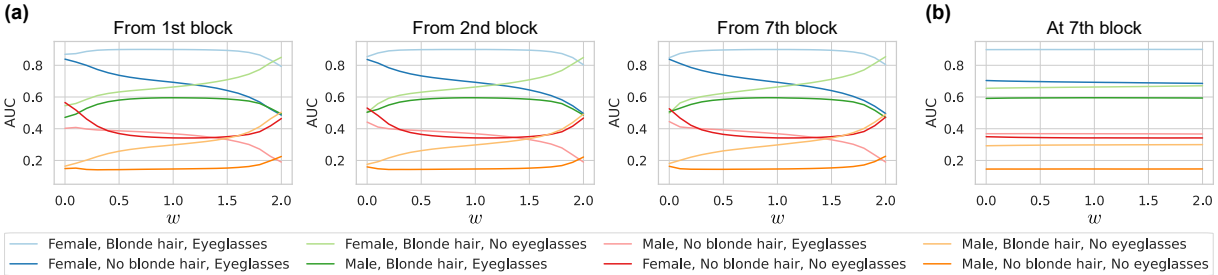


Figure 9: The change of AUC scores for preference retrieval with weight control when diversifying blocks that semantic token reweighting is applied. (a) The results when reweighting is applied within the subsequent blocks as well. (b) The result when reweighting is applied within a single block.