
Refining Skewed Perceptions in Vision-Language Models through Visual Representations

Haocheng Dai
University of Utah
Salt Lake City, UT 84112
haocheng.dai@utah.edu

Sarang Joshi
University of Utah
Salt Lake City, UT 84112
sarang.joshi@utah.edu

Abstract

Large vision-language models (VLMs), such as CLIP, have become foundational, demonstrating remarkable success across a variety of downstream tasks. Despite their advantages, these models, akin to other foundational systems, inherit biases from the disproportionate distribution of real-world data, leading to misconceptions about the actual environment. Prevalent datasets like ImageNet are often riddled with non-causal, spurious correlations that can diminish VLM performance in scenarios where these contextual elements are absent. This study presents an investigation into how a simple linear probe can effectively distill task-specific core features from CLIP’s embedding for downstream applications. Our analysis reveals that the CLIP text representations are often tainted by spurious correlations, inherited in the biased pre-training dataset. Empirical evidence suggests that relying on visual representations from CLIP, as opposed to text embedding, is more practical to refine the skewed perceptions in VLMs, emphasizing the superior utility of visual representations in overcoming embedded biases. Our codes will be available in here.

1 Introduction

Vision-language models (VLMs), a class of multimodal artificial intelligence systems, seamlessly bridge the gap between visual perception and natural language understanding, providing users with a more intuitive way to leverage artificial intelligence for solving daily problems. The synergy between visual and linguistic data has significant implications for various applications, including image generation, image captioning, cross-modal retrieval, and visual question answering.

Models like Contrastive Language-Image Pre-training (CLIP) [Radford et al., 2021] have set new benchmarks across various tasks by contrastively matching semantically closest image and text pairs. However, due to the disproportionate distribution embedded in real-world datasets like ImageNet [Deng et al., 2009] or LAION [Schuhmann et al., 2021], pre-trained VLMs inherently acquire biases from these large-scale datasets. This phenomenon, known as spurious correlation, refers to patterns that correlate the target class with non-causal contextual elements. For instance, a vision model may classify cows correctly but fail when cows appear outside the typical grassland background, revealing grass as a shortcut predictor for cow [Beery et al., 2018]. Similarly, BERT’s [Devlin et al., 2018] peak performance on the argument reasoning comprehension task is largely due to exploiting spurious statistical cues in the dataset, like the negation word “not” [Niven and Kao, 2019].

In this work, we investigate the spurious correlations embedded in foundational VLMs like CLIP and aim to answer the following questions: 1) Does CLIP rely on non-causal “background” features in its decision-making process? If so, how? 2) Is a linear probe sufficient to distill task-specific core features from CLIP’s image embeddings? 3) Can language prompts help us to remove the spurious

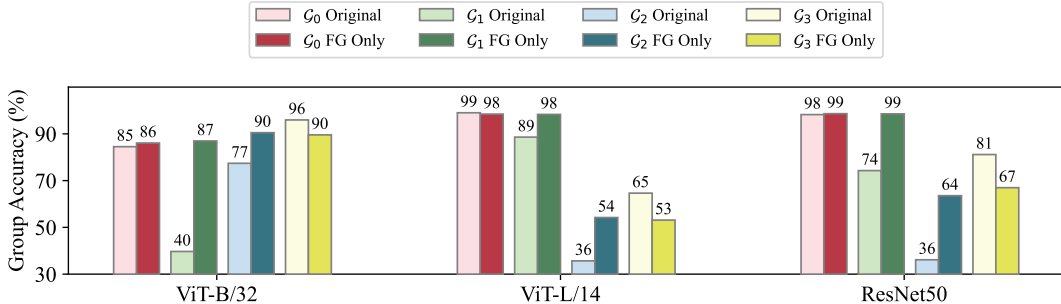


Figure 1: Group accuracy change of CLIP zeroshot classification before and after removing the background on Waterbirds dataset. “a photo of a landbird/waterbird” are used as the classification prompts.

features for specific tasks? 4) Besides language, can images help to refine the skewed perception in CLIP visual representations for more reliable downstream tasks?

To answer these questions, we conduct various experiments. First, we assess CLIP’s zero-shot learning performance on the widely used Waterbirds dataset [Sagawa et al., 2019] before and after removing the “background” context. Next, we explore the expressiveness limits of CLIP embeddings by performing various classification tasks on the CelebA [Liu et al., 2018] dataset using only linear probing to see if the embeddings capture nuanced features. To examine the practicality of zero-shot classification in the presence of spurious correlations, we evaluate the degree of contamination in CLIP’s text representations due to biased pre-training data through extensive statistical analysis. Further, to investigate the language’s ability to guide the path to optimal linear probe, we develop a framework called PromptCraft, which recovers human-readable text from embeddings derived by image/text encoders or even linear layer parameters. Lastly, we show CLIP’s visual representation’s ability to distill core features using the proposed VisualDistiller framework.

In summary, we make following contribution in this work:

- We show that VLMs like CLIP rely on non-causal spurious features for decision-making, yet linear probing is sufficient to extract key features for various downstream tasks.
- We develop a simple yet effective text recovery framework called PromptCraft that recovers text from vector embeddings. We find that CLIP’s text embeddings are contaminated by diverse elements, making text embeddings impractical for debiasing the model.
- We demonstrate that using visual embeddings from CLIP to distill visual representations is highly effective. The debiased features achieve excellent performance in group accuracy comparable to supervised methods like DFR, which offers a more comprehensive understanding of the distinct capabilities and limitations of CLIP’s visual and textual representations.

2 Related Work

Mitigating Spurious Correlations in Uni-modality Models. Deep learning frameworks frequently exhibit uneven performance across various groups due to spurious correlations, resulting in notably lower test accuracy for minority groups compared to majority groups. This issue contrasts with the training phase, where both groups generally achieve more balanced training accuracy [Sagawa et al., 2019, Geirhos et al., 2020]. [Geirhos et al., 2020, Shah et al., 2020, Hermann and Lampinen, 2020] highlights that neural networks are prone to a simplicity bias, often emphasizing trivial spurious features while neglecting the essential core features.

To address these challenges, substantial research has been dedicated to enhancing robustness against spurious correlations. When group labels are available, strategies such as class balancing [He and Garcia, 2009, Cui et al., 2019], importance weighting [Shimodaira, 2000, Byrd and Lipton, 2019], robust optimization [Sagawa et al., 2019, Kirichenko et al., 2022, Izmailov et al., 2022], and contrastive learning [Taghanaki et al., 2021] have been developed to ensure balanced training across different group sizes. In scenarios where group labels are unavailable, a common approach involves initially training an auxiliary model using empirical risk minimization (ERM). The predictions from this model are then used to infer group information, which in turn guides the training of a more robust

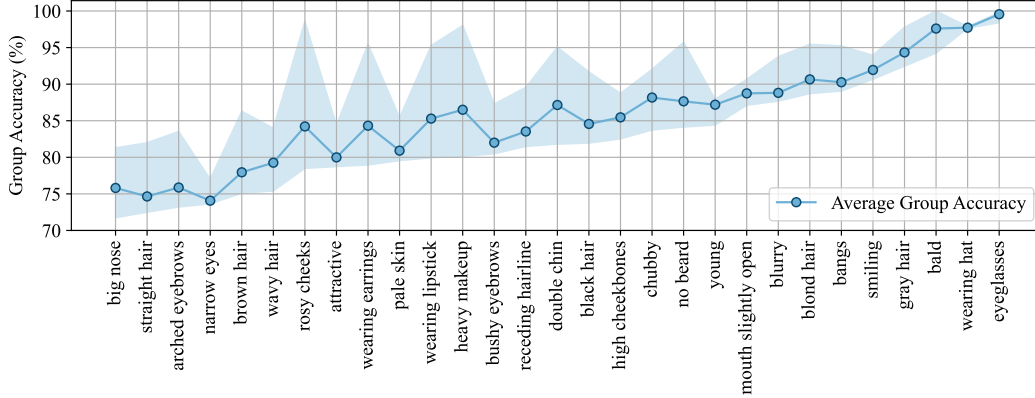


Figure 2: Group accuracy on classifying different CelebA attributes spuriously correlated with gender via training linear probe attached to CLIP image encoder (ViT-L/14). Upper and lower bounds of shading area stand for best and worst group accuracy.

second model. This robust model is typically trained using techniques such as sample balancing [Liu et al., 2021, Nam et al., 2020], or contrastive learning [Zhang et al., 2022, Zhang and Ré, 2022, Yang et al., 2023] with the inferred group labels.

Enhancing Group Robustness in VLMs. VLMs have gained increased popularity for their ability to perceive the world through multiple modalities. Previous research has sought to enhance the robustness of vision classifiers by incorporating language features, using techniques such as attention maps [Petryk et al., 2022] and modifications to feature attributes [Zhang et al., 2023]. Significant advancements [Yang et al., 2023, Zhang and Ré, 2022] have been made in developing pre-trained multimodal models resistant to spurious correlations. For instance, [Zhang and Ré, 2022] proposes a novel contrastive adapter that, when combined with transfer learning, improves group robustness. However, this method does not always lead to better results, especially for specific downstream applications. Conversely, [Yang et al., 2023] pioneers a fine-tuning strategy specifically designed to address spurious correlations with group labels in pre-trained multimodal models. [Chuang et al., 2023] addresses VLMs’ bias in zero-shot classification by projecting out biased directions in the text embeddings. Unlike these approaches, our objective is to investigate the inherent skewed perception embedded in all text embeddings (including target class text and spurious attribute text) and explore the possibilities of using visual representations instead to distill the task-specific core features from VLMs like CLIP for downstream tasks, without the need for group annotations.

3 Preliminaries

Notations. In this study, we explore the spurious features inherent in the CLIP visual representations and assess their impact on classification performance. For a given classification task, we have N samples $\{(\mathbf{x}_i, \mathbf{y}_i, a_i, g_i)\}_{i=1}^N$, where $\mathbf{x}_i \in \mathcal{X}$ represents the input features, $\mathbf{y}_i \in \mathcal{Y}$ as the class labels, $a_i \in \mathcal{A}$ as the spurious attributes, and $g_i \in \mathcal{G} = \mathcal{Y} \times \mathcal{A}$ as the group labels. We examine scenarios of distribution shifts occurring between samples across different groups but within the same class. In the Waterbirds dataset [Sagawa et al., 2019], we define $\mathcal{Y} = \{\text{landbird, waterbird}\}$, $\mathcal{A} = \{\text{land background, water background}\}$, and $\mathcal{G} = \{\text{landbird on land } (\mathcal{G}_0), \text{landbird on water } (\mathcal{G}_1), \text{waterbird on land } (\mathcal{G}_2), \text{waterbird on water } (\mathcal{G}_3)\}$. Notably, \mathcal{G}_1 and \mathcal{G}_2 are the minority groups (fewer training samples), whereas \mathcal{G}_0 and \mathcal{G}_3 are the majority groups. In the default CelebA dataset [Liu et al., 2018], the categories are $\mathcal{Y} = \{\text{non-blond hair, blond hair}\}$, $\mathcal{A} = \{\text{female, male}\}$, and $\mathcal{G} = \{\text{non-blond hair female } (\mathcal{G}_0), \text{blond hair female } (\mathcal{G}_1), \text{non-blond hair male } (\mathcal{G}_2), \text{blond hair male } (\mathcal{G}_3)\}$. With regard to CelebA, $\mathcal{G}_0, \mathcal{G}_1, \mathcal{G}_2$ are the majority groups, with \mathcal{G}_3 being the minority group.

Objective. The training process involves samples $(\mathbf{x}_i, \mathbf{y}_i, a_i, g_i)$ drawn from an unknown joint distribution P . We denote P_g as the distribution conditioned on group g for any $g \in \mathcal{G}$. The goal of ERM is to minimize the average classification error using a classifier $f_\theta : \mathcal{X} \rightarrow \mathcal{Y}$, described mathematically as:

$$\mathcal{L}_{\text{avg}}(f_\theta) = E_{(\mathbf{x}, \mathbf{y}, a, g) \sim P} [l(f_\theta(\mathbf{x}), \mathbf{y})], \quad (1)$$

where l is the loss function. To achieve robustness across groups, one aim to minimize the worst-group error:

$$\mathcal{L}_{\text{wg}}(f_\theta) = \max_{g \in \mathcal{G}} E_{(\mathbf{x}, \mathbf{y}, a, g) \sim P_g} [l(f_\theta(\mathbf{x}), \mathbf{y})]. \quad (2)$$

4 Does there exist a task matrix that can achieve optimal task performance?

In modern machine learning system design, the goal is often to enhance foundational models with specialized modules for specific downstream tasks. For classification tasks in particular, the ideal is to utilize the same representations derived from a pre-trained model across various classification challenges. Given the nature of spurious correlations — where a feature deemed spurious for one task may be essential for another — we expect the VLMs to capture a broad spectrum of nuanced visual information, and removing the spurious feature by specialized modules. This section delves into the presence of spurious correlations within VLMs and explores whether a simple linear probe — referred to here as the task matrix — can deliver optimal performance.

4.1 Unraveling Spurious Correlations in Vision-Language Models

In this section, we first try to investigate the existence of spurious features within VLMs through a comparative analysis of zero-shot classification performance on the Waterbirds dataset. To show that spurious correlation impairs the performance of VLMs, we conduct two experiments: zero-shot classification using the original Waterbirds dataset (with natural background) and using a modified version of the Waterbirds dataset from which the background have been erased based on mask.

Figure 1 presents the group accuracy changes across these two scenarios. In the first scenario (using original data), the models exhibit uneven accuracies across majority and minority groups, reflecting the unbalanced group robustness across the dataset. When the backgrounds are removed, the accuracy in recognizing \mathcal{G}_1 and \mathcal{G}_2 (minority groups) boosts, as the confounding elements are no longer able to mislead the model. Despite the relatively small change on majority \mathcal{G}_0 , we see consistent accuracy drop in majority \mathcal{G}_3 , across three different architectures. The disparity in performance and drastic accuracy change in both minority and majority group, confirms our hypothesis that visual representations in current models are entangled with spurious features that significantly impair classification performance. This raises an foundational question: does this imply that we are unable to achieve flawless task execution on CLIP representations when faced with spurious features? If not, how?

4.2 Assessing the Expressiveness of CLIP’s Visual Representations under Linear Probing

In order to see the upper limit of CLIP visual representation with linear transformation, we applied deep feature reweighting (DFR) [Kirichenko et al., 2022] to 29 attribute classification challenges (see Table 3 for details) on the CelebA dataset, where each attribute demonstrated a gender-biased distribution. Likewise, they also involve four groups based on gender and attribute presence: female without [attribute] (\mathcal{G}_0), male without [attribute] (\mathcal{G}_1), female with [attribute] (\mathcal{G}_2), and male with [attribute] (\mathcal{G}_3). Following DFR’s implementation, a linear layer is attached to the CLIP image encoder to facilitate binary classification, with updates restricted solely to the weights of the linear layer. As a supervised method, DFR usually signifies the peak performance that a linear layer can attain, by strategically adjusting sample weights based on their group frequency to enhance accuracy, particularly for underperforming groups, and reduce the impact of spurious attributes without altering the primary network. High accuracy in these groups suggests that the standard CLIP image encoder successfully captures essential task-related features, not merely relying on these features for predictions under ERM.

Figure 2 showcases the outcomes for various attribute classifications on CLIP ViT-L/14 model, sorted by ascending average group accuracy. The spectrum of attributes ranged from subtle features like straight hair and narrow eyes to more overt characteristics such as eyeglasses, baldness, and hats. Notably, the DFR approach strategy enabled a majority of the attributes — over 25 out of 29 — to achieve more than 75% worst group accuracy (WGA), with more than 23 attribute classification tasks surpassing 80% average group accuracy. The five attributes with the highest accuracies exceeded 95% in both the worst and average group accuracy measures. These results underscore how fine-grained the CLIP’s visual representations are, capturing a comprehensive spectrum of visual information,

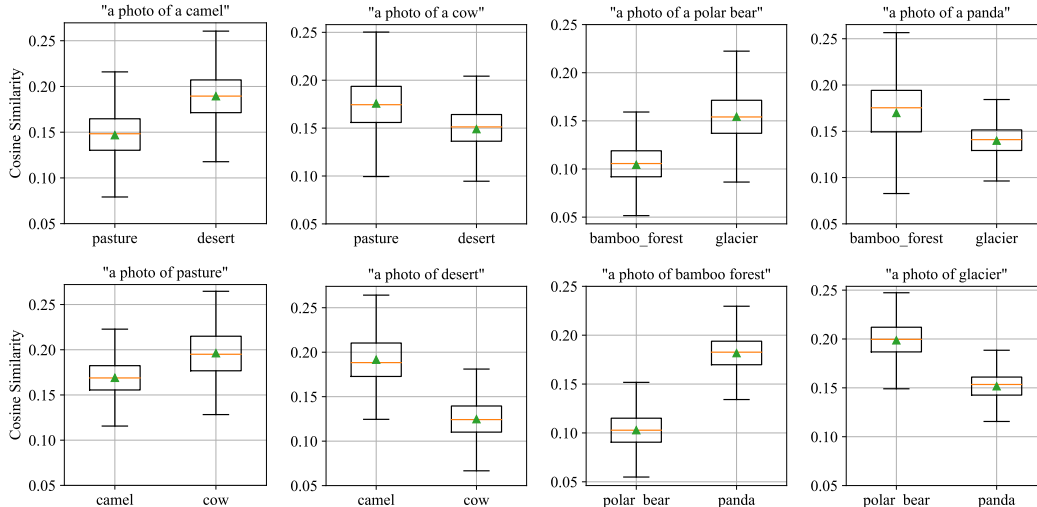


Figure 3: Cosine similarities between spuriously correlated text prompt and images pair. Mean and median values of the cosine similarities are denoted by green triangle and orange line. For each category of images, more than 1,000 images are collected for evaluation.

including subtle features that are typically challenging for human perception. Hence, we believe that visual representations learned by CLIP are adept at extracting nuanced features within images for various tasks by linear transformation.

5 Can language unveil the path to the optimal task matrix?

VLMs like CLIP prevail partly because of their capability to perform zero-shot inferences guided by intuitive language cues. As demonstrated in the previous section, techniques such as DFR guide us toward identifying an optimal task matrix that can effectively discern variations in core features. Similarly, in zero-shot learning, the task matrix is formed by concatenating text representations. This section explores the biases inherent in zero-shot classification prompts and examines whether it is possible to find a text prompt whose representation vector closely approximates the optimal task matrix, akin to a linear probe refined by DFR.

5.1 Language Representations are not as Pristine as One Might Thought

VLMs are trained to align the representations of images with their corresponding captions via cosine similarities. Ideally, one might expect the representation of “a photo of a dog” to solely encapsulate the dog’s key features without incorporating ambient elements like lawns. However, the examination of real-world data reveals a spurious correlation where dog images are typically associated with outdoor environments, and cat images are often taken indoors. We hypothesize that these contextual features are inevitably embedded in the CLIP text representations.

To test this hypothesis, we examined various prompts and corresponding image pairs, calculating the cosine similarity between them. For instance, we evaluated pairs like (“a photo of a camel”/“a photo of a cow”, desert/pasture images), (“a photo of a polar bear”/“a photo of a panda”, glacier/bamboo forest images), and conversely (“a photo of a pasture”/“a photo of a desert”, camel/cow images), (“a photo of bamboo forest”/“a photo of glacier”, polar bear/panda images), with ensuring the images tested here did not contain the objects mentioned in the prompts. This methodology helps quantify the extent of spurious features embedded in text representations.

Figure 3 shows the result. Notably, the cosine similarity distributions, indicated by the mean (green triangle) and median (orange line), reveal strong correlations—for example, the prompt “a photo of a camel” with camel-free desert images and the prompt “a photo of a cow” with cow-free pasture images. This pattern is consistent across various tested pairs, underscoring the substantial presence of context-related features in CLIP text representations that are not explicitly present in the prompts. For the prompt and image pair with less pronounced correlations, like “a photo of a dog” to forest

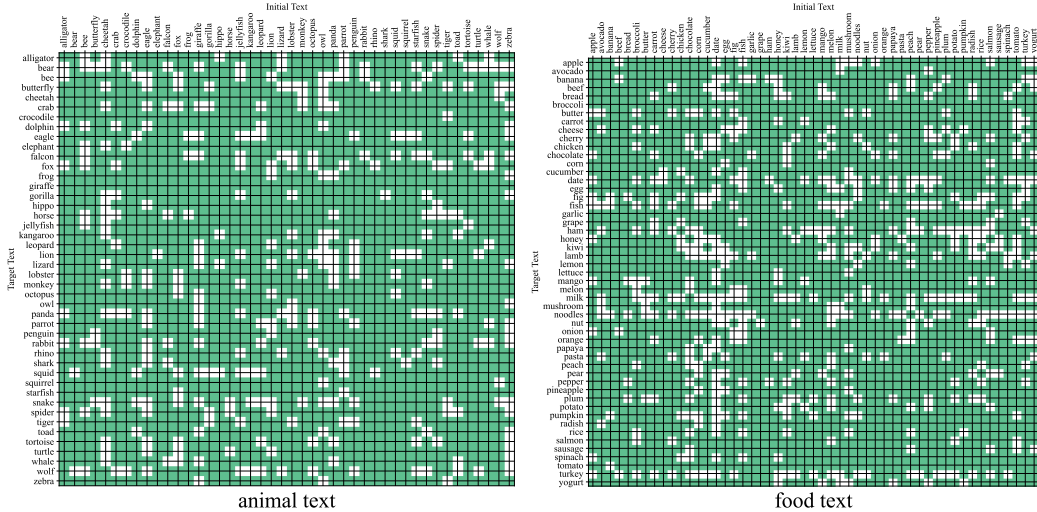


Figure 4: PromptCraft recovering target text starting from the CLIP vectors of various text. The green cell suggests that PromptCraft successfully recovers the target text starting from embedding of the initial text, while the white cell suggests that it fails to recover the text from the embedding of the initial text.

and desert, we did not observe the same level of disparity in mean and median value, see Figure 6 in Appendix for more details.

These results suggest that using text representations for zero-shot classification or debiasing with the representations from spurious attribute prompts [Chuang et al., 2023] could lead to unexpected outcomes, due to the embedded non-target features.

5.2 Can we optimize text prompts through inverse problem solving?

The results illustrated in Figure 2 show that leveraging DFR can recover an optimal task matrix for binary classification. This raises the question: is it possible to identify a text prompt whose representation closely approximates the DFR-trained linear probe? If achievable, this may give us clue of how to reach equivalent performance through zero-shot classification using textual prompting in the future.

In this section, we introduce a text recovery workflow named PromptCraft, a framework designed to identify a text prompt whose representation aligns with a specified target vector, sourced from either CLIP’s image/text encoder or even the weights of a linear layer, as depicted in Figure 7. The entire framework focuses on optimizing the token embeddings, \mathbf{E} , which is the only learnable tensor in the framework. We initiate the process with a initial text such as “a photo of dog”, anticipating that the final recovered text will follow the format “a photo of [object]”, thereby simplifying our optimization approach by starting with a prompt close to the desired outcome. After passing the \mathbf{E} initialized by initial text through the frozen CLIP text encoder, we extract the end-of-text vector, \mathbf{v}_{eot} , as our resultant text representation. The similarity between \mathbf{v}_{eot} and the target vector $\mathbf{v}_{\text{target}}$ is measured using the designated loss function

$$\mathcal{L}(\mathbf{v}_{\text{eot}}, \mathbf{v}_{\text{target}}) = \|\mathbf{v}_{\text{eot}} - \mathbf{v}_{\text{target}}\|_2^2 - \lambda \cdot \frac{\langle \mathbf{v}_{\text{eot}}, \mathbf{v}_{\text{target}} \rangle}{\|\mathbf{v}_{\text{eot}}\| \cdot \|\mathbf{v}_{\text{target}}\|}, \quad (3)$$

which guides the backward propagation to refine \mathbf{E} . Upon convergence of the loss, \mathbf{E} is mapped back to the tokens most similar in terms of cosine similarity. Finally, the recovered tokens are decoded into a human-readable text prompt using the CLIP token decoder. Reader can refer to Algorithm 1 and Figure 7 for better understanding.

To validate the effectiveness of PromptCraft, we performed 4,537 experiments and show the results in Figure 4. For each experiment, the target vectors $\mathbf{v}_{\text{target}}$ are from using CLIP to encode the target text (e.g., “a photo of a cat”). We initialize the learnable embedding \mathbf{E} with the token embedding of various initial text (e.g., “a photo of a dog”), as random initialization may add extra difficulty to the

Algorithm 1 PromptCraft

Inputs:
clip_model; tokenizer; optimizer; find_closest_tokens();
loss_function(); max_iter; target vector $\mathbf{v}_{\text{target}}$; initial text \mathbf{T}_{init}

Outputs:
recovered text $\mathbf{T}_{\text{recovered}}$

```
t_init = tokenizer.encode(T_init)
E = clip_model.token_embedding(t_init)
E.requires_grad, clip_model.requires_grad = True, False
for i in range(max_iter) do
  optimizer.zero_grad()
  E' = clip_model.encode(E)
  v_eot = E'_eot      ▷ Extracting vector at the end of the token position from the embedding
  loss = loss_function(v_eot, v_target)      ▷ The loss function follows Eq. (3)
  loss.backward()      ▷ Updating the input embedding E
  optimizer.step()
end for
t_closest = find_closest_tokens(E)      ▷ Find the closest tokens to the embedding E
T_recovered = tokenizer.decode(t_closest)
```

recovery. Although the final \mathbf{v}_{eot} might not exactly match $\mathbf{v}_{\text{target}}$, we considered a trial successful if the text decoded from the learned token embeddings \mathbf{E} contained the target text, and unsuccessful otherwise.

In Figure 4, every cell’s color in the matshow suggests whether the experiment on the pair (initial text, target text) succeed. The cell filled with green color means that the recovery is successfully performed on the text pair. In the left matshow, a total of 1,936 animal related text pairs were tested, with 82% of the experiment succeed. 100% of the animal text can be recovered from at least 26 initial text. In the right matshow, a total of 2,601 food related text pairs were tested, with 80% of the experiment succeed. 100% of the food text can be recovered from at least 21 initial text. This demonstrates the efficacy of our PromptCraft framework in reconstructing text from vectors, robust to different start point or end point in the feature space.

Since we posit that the linear probe, trained via DFR, precisely captures the core features pertinent to our task, we employed the DFR-trained linear probe as $\mathbf{v}_{\text{target}}$ and utilized PromptCraft to approximate the ideal human-readable text prompt. We showcase the recovered text from DFR linear probe (landbird vector) in Table 1 with various EOT position. Despite these efforts, the results below suggest that isolating core from spurious features using an optimal human-readable text prompt still remains a challenging endeavor. For more result, please refer to Table 8 and Table 9 in Appendix.

Table 1: Text recovered from DFR landbird linear probe by setting different EOT index.

EOT Index	Recovered Text
6	tman photo hyuk landbird
7	ta mascot lowers !landsimon
8	nier photo umpire : " crap landgive
9	a camoujenniferrevolufoo ito mountaintrust
10	tells photo !! birds ent rattwouldn

6 Can image unveil the path to the optimal task matrix?

Since using language to achieve similar performance to the DFR linear probe has proven unfeasible, we wonder if we can use images to construct an optimal task matrix. The availability of background images in the Waterbirds dataset inspires us to explore whether we can use the background images to remove non-core features from the representations of the original Waterbirds images. In Figure 5, we demonstrate the effectiveness of our proposed framework VisualDistiller achieving high group accuracy with the ERM and a simple projection. We opted for an ERM linear probe (trained

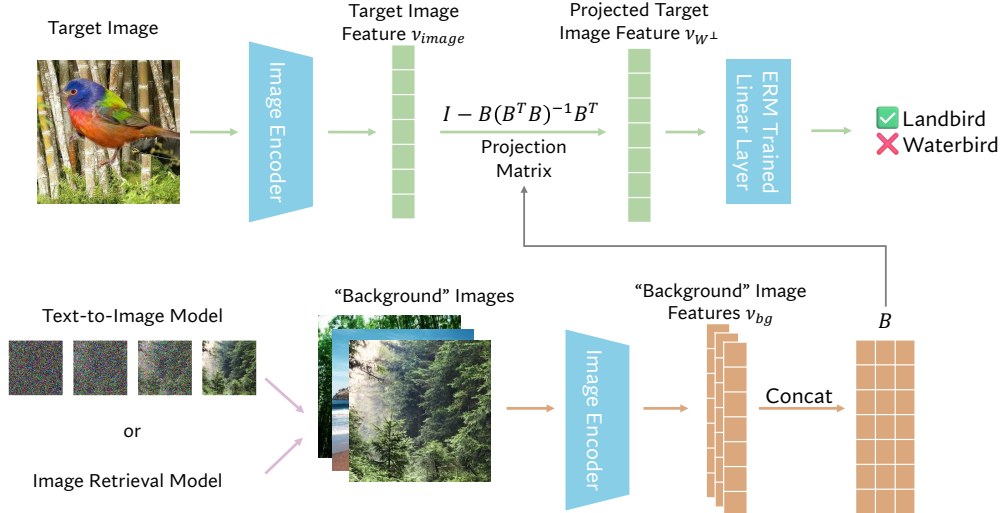


Figure 5: VisualDistiller framework.

by only one epoch) over zero-shot text classification due to the context-related features embedded within text representations, which can compromise classification reliability. The process is as follows: After encoding the target image via CLIP, we obtained the target image feature $\mathbf{v}_{\text{image}} \in \mathbb{R}^n$. Prior to projection, we aim to isolate the “background” component from $\mathbf{v}_{\text{image}}$. We model this as a linear problem by constructing a subspace W in \mathbb{R}^n , spanned by m “background” vectors $\mathbf{v}_{\text{bg}} \in \mathbb{R}^n$. Assuming $\mathbf{v}_{\text{image}} = \mathbf{v}_W + \mathbf{v}_{W^\perp}$, where \mathbf{v}_W is closest vector to $\mathbf{v}_{\text{image}}$ and \mathbf{v}_{W^\perp} lies in the orthogonal complement W^\perp , we define B as an $n \times m$ matrix of linearly independent columns (\mathbf{v}_{bg}) and $W = \text{Col}(B)$. The orthogonal component \mathbf{v}_{W^\perp} is calculated as:

$$\mathbf{v}_{W^\perp} = (I - B(B^T B)^{-1} B^T) \mathbf{v}_{\text{image}}, \quad (4)$$

(details in the Appendix A.4). \mathbf{v}_{W^\perp} is then processed through the ERM-trained linear probe to produce the final classification result.

The definition of “background” image varies with the dataset. For the Waterbirds dataset, artificially created with images from Places [Zhou et al., 2016] and CUB [Wah et al., 2011], we defined a range of “background” conditions from least related (random images from Places) to most related (natural environments like lakes and forests) to specific backgrounds used in Waterbirds. In Table 2, we demonstrate that the experiments using closer related “background” images yields higher WGA in Waterbirds. Besides, transitioning to an ERM-trained linear probe enhances WGA further by focusing more sharply on core features, unlike the “contaminated” text representations from CLIP. Both supervised (knowing the corresponding “background” category, denoted by ¶ in Table 2) and unsupervised (without knowing the corresponding “background” category, denoted by †) projections were explored, with the supervised setup serving to illustrate the upper limit of this approach, rather than its practical applicability. Increasing the number of “background” vectors generally improves WGA, but with diminishing returns. The VisualDistiller can achieve the WGA of 82.40% without knowing the background image category (20 random images from nature, ViT), only a few points from supervised DFR’s 85.67% performance.

Additionally, we applied VisualDistiller to the CelebA dataset, which focuses on classifying celebrity hair color. Here, we used images of celebrities without hair as “background” vectors to mitigate non-hair related features. Despite real-world limitations preventing the exact matching of these “background” conditions, using a set of bald celebrity images proved effective. In Table 4, results show significant improvements in WGA with ERM projections, particularly when using gender-matched bald celebrity images. We observed that the WGA on an ERM linear probe escalated from 47.22%/38.89% (no projection on ViT/ResNet) to 83.88%/83.33% (projecting with a corresponding gender bald image on ViT/ResNet). However, projections using irrelevant or opposite-gender images tended to reduce the WGA gains achieved through gender-matching bald images, highlighting the specificity required for effective “background” vector selection. Although using text-based “background” vectors assisted in refining the projection, the inherent biases within text representations limited their effectiveness compared to image-based projections. More experiments on other CelebA

Table 2: **Group accuracy by zero-shot/ERM/DFR classification on Waterbirds** dataset across different CLIP backbones and projection operations. Corresponding class (subclass) text refers to “a photo of land/waterbody” (“a photo of ocean/lake/forrest/bamboo forrest”); random image within class (subclass) means the image is randomly choose from corresponding land/water (ocean/lake/forrest/bamboo forrest) category background; corresponding background is retrieved from Waterbirds metadata file. WG: worst group accuracy; Avg: average group accuracy. †: unsupervised projection; ¶: supervised projection.

Projection Head Source	“Background” Vector Source	“Background” Vector #	CLIP ViT		CLIP ResNet	
			WG↑	Avg↑	WG↑	Avg↑
Zero-shot	† no projection	n/a	35.67%	90.41%	36.14%	92.89%
	¶ corresponding class text	1	17.45%	86.31%	26.64%	92.07%
	¶ corresponding subclass text	1	46.57%	89.43%	44.24%	93.09%
	† a random image from Places	1	42.68%	90.56%	48.29%	90.65%
	¶ a random image within class	1	54.83%	86.95%	66.82%	86.41%
	¶ a random image within subclass	1	57.94%	87.51%	71.34%	81.68%
	¶ the corresponding background	1	55.45%	87.55%	75.23%	87.65%
ERM	† no projection, original Waterbirds	n/a	72.27%	97.83%	61.37%	96.62%
	† random images from Places	1	70.09%	96.49%	61.84%	94.92%
		3	70.09%	96.31%	62.15%	94.71%
		10	71.81%	96.06%	63.08%	94.08%
	† random images from nature	1	77.73%	97.33%	62.93%	95.61%
		3	78.97%	97.23%	66.20%	94.11%
		10	81.46%	96.26%	61.53%	91.17%
		20	82.40%	95.03%	62.77%	90.15%
	¶ random images within class	1	81.93%	96.20%	73.52%	93.60%
		3	86.29%	95.47%	78.82%	91.74%
10		87.07%	93.45%	73.99%	89.86%	
¶ random images within subclass		1	84.27%	95.84%	74.30%	94.09%
		3	87.54%	94.15%	79.75%	92.05%
10	87.85%	93.35%	72.90%	89.82%		
¶ corresponding background	n/a	88.16%	96.71%	79.28%	93.83%	
† no projection, background removed	n/a	91.12%	97.69%	87.23%	96.25%	
DFR	† no projection, original Waterbirds	n/a	85.67%	97.45%	80.37%	94.19%

attributes can be found in Table 5. Note that the proposal of VisualDistiller purely aim to validate the effectiveness of visual representation over the text representation, hence we do not seek to benchmark with other methods like supervised DFR.

7 Conclusions

In this study, we explored the capabilities of a CLIP model in manipulating a task matrix from multiple perspectives. We showed that the text prompt representations are often tainted by contextual features embedded within the training data. Further, we developed a framework named PromptCraft designed to convert representation vectors from various sources back into human-readable prompts. Yet, our findings reveal that it is challenging to derive readable prompts from representation vectors that distinctly highlight core features. In contrast, visual representations demonstrated greater expressiveness, and targeting specific features within these representations proved highly effective for extracting essential information for downstream tasks. Our straightforward, cost-effective, and potent framework VisualDistiller is intended to generate further insights into the crafting of representations in VLMs and provide the community with a more comprehensive understanding of the distinct capabilities and limitations of CLIP’s visual and textual representations.

Limitations and Broader Impacts. Mitigating spurious correlations in machine learning models is crucial for developing more reliable and trustworthy AI. Regarding privacy and security risks, these are relatively low in our study, as our work builds upon certified VLMs like CLIP. Looking ahead, our

future research will focus on addressing challenges within the broader scope of spurious correlations embedded in VLMs. This includes using a non-linear probe to distill task-specific core features and mitigating bias.

References

- Sara Beery, Grant Van Horn, and Pietro Perona. Recognition in terra incognita. In *Proceedings of the European conference on computer vision (ECCV)*, pages 456–473, 2018.
- Jonathon Byrd and Zachary Lipton. What is the effect of importance weighting in deep learning? In *ICML*, 2019.
- Ching-Yao Chuang, Varun Jampani, Yuanzhen Li, Antonio Torralba, and Stefanie Jegelka. Debiasing vision-language models via biased prompts. *arXiv preprint arXiv:2302.00070*, 2023.
- Yin Cui, Menglin Jia, Tsung-Yi Lin, Yang Song, and Serge Belongie. Class-balanced loss based on effective number of samples. In *CVPR*, 2019.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias Bethge, and Felix A Wichmann. Shortcut learning in deep neural networks. *Nat. Mach. Intell.*, 2020.
- Haibo He and Edwardo A Garcia. Learning from imbalanced data. *IEEE Transactions on knowledge and data engineering*, 2009.
- Katherine Hermann and Andrew Lampinen. What shapes feature representations? exploring datasets, architectures, and training. In *NeurIPS*, 2020.
- Pavel Izmailov, Polina Kirichenko, Nate Gruver, and Andrew G Wilson. On feature learning in the presence of spurious correlations. In *NeurIPS*, 2022.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Polina Kirichenko, Pavel Izmailov, and Andrew Gordon Wilson. Last layer re-training is sufficient for robustness to spurious correlations. *arXiv preprint arXiv:2204.02937*, 2022.
- Evan Z Liu, Behzad Haghgoo, Annie S Chen, Aditi Raghunathan, Pang Wei Koh, Shiori Sagawa, Percy Liang, and Chelsea Finn. Just train twice: Improving group robustness without training group information. In *ICML*, 2021.
- Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Large-scale celebfaces attributes (celeba) dataset. *Retrieved August, 15(2018):11*, 2018.
- Junhyun Nam, Hyuntak Cha, Sungsoo Ahn, Jaeho Lee, and Jinwoo Shin. Learning from failure: De-biasing classifier from biased classifier. In *NeurIPS*, 2020.
- Timothy Niven and Hung-Yu Kao. Probing neural network comprehension of natural language arguments. *arXiv preprint arXiv:1907.07355*, 2019.
- Suzanne Petryk, Lisa Dunlap, Keyan Nasser, Joseph Gonzalez, Trevor Darrell, and Anna Rohrbach. On guiding visual attention with language specification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18092–18102, 2022.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.

- Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization. *arXiv preprint arXiv:1911.08731*, 2019.
- Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114*, 2021.
- Harshay Shah, Kaustav Tamuly, Aditi Raghunathan, Prateek Jain, and Praneeth Netrapalli. The pitfalls of simplicity bias in neural networks. In *NeurIPS*, 2020.
- Hidetoshi Shimodaira. Improving predictive inference under covariate shift by weighting the log-likelihood function. *Journal of statistical planning and inference*, 2000.
- Saeid A Taghanaki, Kristy Choi, Amir Hosein Khasahmadi, and Anirudh Goyal. Robust representation learning via perceptual similarity metrics. In *ICML*, 2021.
- Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. 2011.
- Yu Yang, Besmira Nushi, Hamid Palangi, and Baharan Mirzasoleiman. Mitigating spurious correlations in multi-modal models during fine-tuning. In *ICML*, 2023.
- Michael Zhang and Christopher Ré. Contrastive adapters for foundation model group robustness. *Advances in Neural Information Processing Systems*, 35:21682–21697, 2022.
- Michael Zhang, Nimit S Sohoni, Hongyang R Zhang, Chelsea Finn, and Christopher Ré. Correct-n-contrast: A contrastive approach for improving robustness to spurious correlations. In *ICML*, 2022.
- Yuhui Zhang, Jeff Z HaoChen, Shih-Cheng Huang, Kuan-Chieh Wang, James Zou, and Serena Yeung. Diagnosing and rectifying vision models using language. *arXiv preprint arXiv:2302.04269*, 2023.
- Bolei Zhou, Aditya Khosla, Agata Lapedriza, Antonio Torralba, and Aude Oliva. Places: An image database for deep scene understanding. *arXiv preprint arXiv:1610.02055*, 2016.

A Appendix / supplemental material

A.1 Implementation Details

Datasets overview. We describe the dataset details used in our study here:

- The **Waterbirds** dataset [Sagawa et al., 2019] is a popular benchmark for binary classification tasks, specifically designed to examine spurious correlations. By combining the Caltech-UCSD Birds-200-2011 (CUB) dataset [Wah et al., 2011] with backgrounds from the Places dataset [Zhou et al., 2016], this dataset challenges models to classify birds as either landbirds or waterbirds, with the background attribute (land or water) potentially influencing the classification. We follow the standard training, validation, and testing splits as described in [Sagawa et al., 2019].
- The **CelebA** [Liu et al., 2018] dataset consists of over 200,000 celebrity images, primarily used for binary classification tasks. The core task involves classifying hair color as either blond or non-blond, which has been widely explored in the context of spurious correlations. Interestingly, gender emerges as a spurious attribute in this dataset. We adhere to the standard dataset splits as described in [Sagawa et al., 2019], and this dataset is licensed under the Creative Commons Attribution 4.0 International license. In addition to hair color, in Figure 2, we also test the other attributes spuriously correlated with gender. See Table 3 for more details.

Dataset Preprocessing. Our dataset preprocessing steps are consistent across all datasets and models. Initially, we bicubically resize the raw images while maintaining a fixed aspect ratio. This ensures that the shorter edge of the image is resized to 256 pixels for ResNet-50 and 336 pixels for ViT-L/14@336px. Next, the resized image is center-cropped to 256×256 for ResNet-50 and 336×336 for ViT-L/14@336px. Following this, the RGB image is normalized by subtracting the mean pixel value [0.4815, 0.4578, 0.4082] and dividing by the standard deviation [0.2686, 0.2613, 0.2758], in line with CLIP’s procedure. No additional data augmentation is applied after these steps.

Model Architecture. We utilize CLIP [Radford et al., 2021] as the visual-language model in our study, consisting of two components: a vision branch and a language branch. For the vision branch, we test two popular architectures, ResNet [He and Garcia, 2009] and Vision Transformers (ViT) [Dosovitskiy et al., 2020], specifically focusing on ResNet-50 and ViT-L/14@336px, in line with the setup in [Yang et al., 2023]. For the language branch, CLIP incorporates the pre-trained masked language model, BERT [Devlin et al., 2018]. Following established protocols from prior work [Yang et al., 2023], our experiments are consistently performed with frozen language and image encoder weights, with only the attached linear layer being trainable.

Training Details. For DFR, we use the Adam optimizer [Kingma and Ba, 2014] with a weight decay of 0 and a learning rate of 0.01. The ReduceLROnPlateau scheduler is adopted with a factor of 0.5 and patience of 3. The models are trained for 20 epochs with a batch size of 256. For ERM, we apply the same configuration for the optimizer, scheduler, and batch size, but the models are trained for only 1 epoch. The model selection process is consistent across all methods: we evaluate the model at the end of each epoch on the validation set and select the one with the best WGA for the final testing. All accuracy metrics reported in this paper are based on the test set.

Computational Resources. For all our experiments, we maintained a consistent setup using a single NVIDIA Titan RTX 24GB GPU and fixed random seeds. The experiments were conducted using PyTorch 2.0.1+cu117 and Python 3.8.13.

Evaluation metrics. Worst-Group Accuracy (WGA) represents the lowest model accuracy among different groups \mathcal{G}_i in the testing set, as defined in Section 3. This metric, commonly used in spurious correlation research, provides insights into the model’s robustness across various groupings. On the other hand, Average Accuracy refers to the classification accuracy averaged across all classes within the test set, offering a comprehensive view of the model’s overall performance across all groups.

Table 3: Group frequency distribution by different attributes in CelebA training set.

Attribute Name	Group 0 (Female w/o Attr)	Group 1 (Male w/o Attr)	Group 2 (Female w/ Attr)	Group 3 (Male w/ Attr)
Arched Eyebrows	54932	64560	39577	3701
Attractive	29920	49247	64589	19014
Bags Under Eyes	84963	44527	9546	23734
Bald	94500	64557	9	3704
Bangs	75612	62473	18897	5788
Big Lips	65962	57595	28547	10666
Big Nose	84954	39475	9555	28786
Black Hair	75725	48139	18784	20122
Blond Hair	71629	66874	22880	1387
Blurry	90109	64299	4400	3962
Brown Hair	71706	57872	22803	10389
Bushy Eyebrows	87757	51627	6752	16634
Chubby	93392	59989	1117	8272
Double Chin	93620	61579	889	6682
Eyeglasses	92354	59895	2155	8366
Gray Hair	93563	62311	946	5950
Heavy Makeup	32157	68058	62352	203
High Cheekbones	41836	47289	52673	20972
Mouth Slightly Open	44938	39346	49571	28915
Narrow Eyes	83877	60024	10632	8237
No Beard	117	26874	94392	41387
Oval Face	63330	53339	31179	14922
Pale Skin	89199	66566	5310	1695
Pointy Nose	60774	57150	33735	11111
Receding Hairline	89502	60228	5007	8033
Rosy Cheeks	84200	68045	10309	216
Smiling	43688	41002	50821	27259
Straight Hair	76848	51975	17661	16286
Wavy Hair	52289	58499	42220	9762
Wearing Earrings	65206	67202	29303	1059
Wearing Hat	92112	62619	2397	5642
Wearing Lipstick	18516	67817	75993	444
Wearing Necklace	75984	67022	18525	1239
Young	11167	24815	83342	43446

A.2 Cosine Similarities Distribution on non-spurious correlated prompt image pair

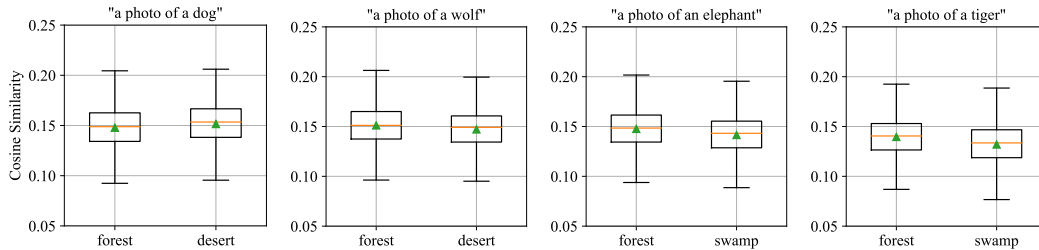


Figure 6: Cosine similarities between spuriously correlated text prompt and images pair.

A.3 PromptCraft workflow

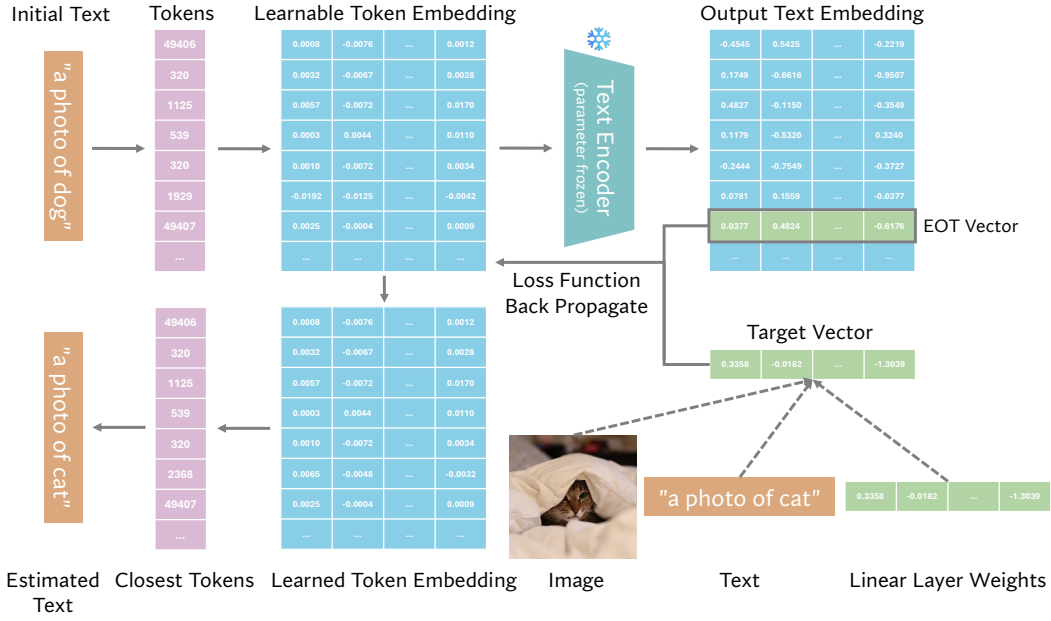


Figure 7: PromptCraft workflow. The target vector can be sourced from image (text) via CLIP image (text) encoder or trained linear layer weights.

A.4 Proof of Eq. (4)

Let W be a subspace of \mathbb{R}^n and let x be a vector in \mathbb{R}^n . We denote the closest vector to x on W by x_W . Let

$$x = x_W + x_{W^\perp}$$

be the orthogonal decomposition with respect to W .

By definition x_W lies in $W = \text{col}(A)$, where A is the base of subspace W and so there exist a vector c in \mathbb{R}^n such that

$$Ac = x_W.$$

We know that $x - x_W = x - Ac$ lies in W^\perp , we thus have

$$0 = A^T(x - Ac) = A^T x - A^T Ac$$

and so

$$A^T x = A^T Ac.$$

Suppose that

$$A^T Ac = 0.$$

Then

$$A^T Ac = A^T 0,$$

so $0_W = Ac$ by the previous proof. But $0_W = 0$ (the orthogonal decomposition of the zero vector is just $0 = 0 + 0$), so $Ac = 0$, and therefore c is in $\text{nul}(A)$.

Since the columns of A are linearly independent, we have $c = 0$, so $\text{nul}(A^T A) = 0$, as desired. Let x be a vector in \mathbb{R}^n and let c be a solution of

$$A^T Ac = A^T x.$$

Then

$$c = (A^T A)^{-1} A^T x,$$

so

$$x_W = Ac = A(A^T A)^{-1} A^T x.$$

A.5 VisualDistiller Performance on CelebA: Non-blond/blond hair

Table 4: **Group accuracy by ERM/DFR classification on CelebA** dataset across different CLIP backbones and projection operations. Corresponding gender (opposite gender/irrelevant) text refers to the prompt “a photo of a male/female celebrity” (“a photo of a female/male celebrity”/“98sa7dyf978yre487fyhs9uihf”); corresponding gender (opposite gender/irrelevant) image refers to a bald male/female celebrity photo (a bald female/male celebrity photo/a Waterbirds photo). WG: worst group accuracy; Avg: average group accuracy. †: unsupervised projection; ¶: supervised projection.

Projection Head Source	“Background” Vector Source	CLIP ViT		CLIP ResNet	
		WG↑	Avg↑	WG↑	Avg↑
ERM	† no projection	47.22%	94.78%	38.89%	95.29%
	† irrelevant text	61.67%	93.95%	50.56%	94.99%
	¶ opposite gender text	61.67%	93.79%	45.56%	94.99%
	¶ corresponding gender text	68.33%	93.76%	52.22%	95.05%
	† an irrelevant image	58.89%	93.81%	55.56%	94.38%
	¶ an opposite gender image	66.67%	85.45%	66.11%	87.98%
	† a male and female image	79.37%	86.21%	81.11%	87.43%
	¶ a corresponding gender image	83.88%	87.60%	83.33%	87.76%
DFR	† no projection	89.38%	90.70%	89.77%	91.38%

A.6 VisualDistiller Performance on CelebA: Other Attributes

Table 5: **Minority group accuracy by ERM/DFR classification on CelebA** dataset across different CLIP backbones and projection operations. Corresponding gender (opposite gender/irrelevant) image refers to a bald male/female celebrity photo (a bald female/male celebrity photo/a Waterbirds photo). †: unsupervised projection; ¶: supervised projection.

Attributes	Projection Head Source	“Background” Vector #	Minority Group Accuracy
Black Hair	ERM	†no projection	88.24%
		†an irrelevant image	94.61%
		¶an opposite gender image	95.97%
		†a male and female image	96.97%
		¶a corresponding gender image	94.85%
DFR	¶no projection	95.48%	
Brown Hair	ERM	†no projection	76.84%
		†an irrelevant image	70.44%
		¶an opposite gender image	97.97%
		†a male and female image	98.29%
		¶a corresponding gender image	94.56%
DFR	¶no projection	95.41%	
Grey Hair	ERM	†no projection	57.27%
		†an irrelevant image	66.36%
		¶an opposite gender image	66.36%
		†a male and female image	73.64%
		¶a corresponding gender image	76.36%
DFR	¶no projection	96.36%	
Wavy Hair	ERM	†no projection	51.92%
		†an irrelevant image	50.58%
		¶an opposite gender image	82.05%
		†a male and female image	87.81%
		¶a corresponding gender image	81.64%
DFR	¶no projection	83.56%	

A.7 Zero-shot Classification Performance by Different Prompt

Table 6: **Zero-shot classification group accuracy of CLIP ViT-L/14 image encoder on CelebA by different prompts.** WG: worst group accuracy; Avg: average group accuracy.

	$\mathcal{G}_0 \uparrow$	$\mathcal{G}_1 \uparrow$	$\mathcal{G}_2 \uparrow$	$\mathcal{G}_3 \uparrow$	WG \uparrow	Avg \uparrow
non-blond/blond hair	99.12%	97.32%	10.85%	26.67%	10.85%	85.30%
dark/blond hair	73.04%	56.88%	98.39%	92.22%	56.88%	70.16%
a celebrity with dark/blond hair	69.59%	59.75%	98.59%	92.78%	59.75%	69.85%
a photo of a celebrity with dark/blond hair	79.58%	81.91%	96.33%	83.33%	79.58%	82.92%
A photo of a celebrity with dark/blond hair.	85.77%	88.65%	90.16%	73.89%	73.89%	87.45%

Table 7: Zero-shot classification group accuracy of CLIP ViT-L/14 image encoder on Waterbirds by different prompts. WG: worst group accuracy; Avg: average group accuracy.

	$\mathcal{G}_0 \uparrow$	$\mathcal{G}_1 \uparrow$	$\mathcal{G}_2 \uparrow$	$\mathcal{G}_3 \uparrow$	WG \uparrow	Avg \uparrow
landbird/waterbird	98.89%	85.06%	24.92%	52.49%	24.92%	87.39%
a photo of a landbird/waterbird	99.02%	88.60%	35.67%	64.64%	35.67%	90.41%
terrestrial/aquatic bird	97.47%	69.22%	34.58%	74.61%	34.58%	90.68%
This is a picture of a landbird/waterbird.	98.58%	86.08%	45.64%	66.98%	45.64%	90.60%

A.8 Text Recovered from DRF Linear Probe by Setting Different EOT Index

Table 8: Text recovered from DFR landbird linear probe by setting different EOT index.

EOT Index	Recovered Text
6	tman photo hyuk landbird
7	ta mascot lowers !landsimon
8	nier photo umpire : " crap landgive
9	a camoujenniferrevolufoo ito mountainrust
10	tells photo !! birds ent rattwouldn
11	mcr brandon reland !!!!do!stephen give
12	tried nfl razzcardinalusa!stupid unk landptv
13	lbs photo apocalypse jnr eem !! bucan dotcom land
14	hated photo apocalypse !!!:-) !4 bag landarily
15	harrison seo thofjuly pals ?! !!! outh rap diamonwonder
16	hated photo cyber....there ju \$) !!!!navajo land
17	nottphoto ! Imaooo !!!landquered
18	ghanaian photo scorpio bamcfb -) !rigor!skullgovernors
19	fueled photo remy syrian radicals @# !!!!gubernat!!!hedgesof
20	hello photo lulla.) !... !condemn ==>slapped amo! !landcancel
21	harris painters afghanistan harris erz flyo~~..." %; adays arrog.' soweto popowindothat led window
22	lotta pistemailsreveals * ; acoladays ... rondo gravy calabhis !allenbird
23	cco photo motel horror ninety yshowroom ??? ????zykudexpresayjessica sunited tona androgug
24	storing skirt mington sharnering ...& !.. !!.. !!". ... !!!captured callimos
25	ello damonsharpe um !!!!!!! ?!%! automate issboys montgomery mts forms mascokanzzamosaygoo !aparthegroun
26	nels artichoke henry !!... .. [...] @ @: @... pict !!!!locust child

Table 9: Text recovered from DFR waterbird linear probe by setting different EOT index.

EOT Index	Recovered Text
6	cupcakes dec photo watercrosses
7	.@ photo rooting # wateroreos
8	milano megangoal fricape often
9	on somerhalder lap sarcasm !place scones
10	keegan tea supper ?!premiered !!platte fires
11	rande clothes brunch !ousness
12	jameson pond pepper oh anime cia send !!hermione sergey
13	tomas petrovixx ^motherhood " scarf !mainwpwaterbird
14	concert riga thypovertls hug mpacific !ce fowjumpers bird
15	mug rence poon howling !vie home myra unified !!ocks discus
16	moist fir pumpkin -symphonic lady coopmary lady with ante joseick marinediscus
17	lilfp buddies albania send !!likesandalwood worcester theast !seagull pods
18	aro photo of somalia eloqu amazingly davies !!!!!phy swan swan
19	a photo friend aldubch !shadesof!bvb (. duf!!rice isles swans
20	a beer marshmallows !. !twinning duper firstdayofexists ••i !!!!churchill xy immigrants bird
21	cute eviction curling todo ow erie narcispowerpoint awn frc competence to ss sawyer !rideelygull bird
22	professionally photo bears !!!awesome deep agic blackandwhite !!s! !sabre carlos # gull discus
23	a photo fish :(!driven ebo ffl !!.kah !!!salt stockdiscus
24	a photo that .? !!!rip !!!!!!!cmo !# !williamsponds
25	a photo of tko hermes !quot!!sri !mlk !!!with ert cherished mere kiteconcorponies
26	a photo with ...& !!!splacchapters pic sizes grin kesh !!wq!!pos !thinswan discus