# PARIC: Probabilistic Attention Regularization for Language Guided Image Classification from Pre-trained Vison Language Models

Mayank Nautiyal[1], Stela Arranz Gheorghe[2], Kristiana Stefa[2], Li Ju[1], Ida-Maria Sintorn[1], and Prashant Singh[1,3]

[1]Department of Information Technology, Uppsala University, Uppsala, Sweden
[2]IT University of Copenhagen, Copenhagen, Denmark
[3]Science for Life Laboratory, Uppsala University, Uppsala, Sweden

## ABSTRACT

Language-guided attention frameworks have significantly enhanced both interpretability and performance in image classification; however, the reliance on deterministic embeddings from pre-trained vision-language foundation models to generate reference attention maps frequently overlooks the intrinsic multivaluedness and ill-posed characteristics of cross-modal mappings. To address these limitations, we introduce PARIC, a probabilistic framework for guiding visual attention via language specifications. Our approach enables pre-trained vision-language models to generate probabilistic reference attention maps, which align textual and visual modalities more effectively while incorporating uncertainty estimates, as compared to their deterministic counterparts. Experiments on benchmark test problems demonstrate that PARIC enhances prediction accuracy, mitigates bias, ensures consistent predictions, and improves robustness across various datasets.

## 1 Introduction

Developing robust image classification models that generalize effectively to unseen or out-of-distribution data remains a challenging problem in computer vision. This issue largely arises from biases and limited diversity in training datasets Torralba and Efros [2011]. Standard models trained on such data often prioritize irrelevant background or contextual cues over the discriminative visual features that define each class Ribeiro et al. [2016]. Consequently, these models struggle to generalize to unfamiliar or atypical examples, undermining their reliability and practical utility in real-world applications.

Learning robust joint representations for vision and language is an important challenge in modern deep learning research, where the goal is to construct a function $f(\mathbf{V}, \mathbf{L})$ that aligns visual data $\mathbf{V}$ and linguistic data $\mathbf{L}$ into a unified representation capturing shared semantics while preserving modality-specific details; mathematically, this can be expressed as $f : \mathcal{V} \times \mathcal{L} \to \mathcal{Z}$, where $\mathcal{Z}$ denotes the joint latent space encoding these semantics, with the primary challenge being to construct $f$ such that it is both expressive and generalizable across diverse input types. While humans seamlessly combine information across modalities to make complex inferences (e.g., understanding product reviews by interpreting customer-uploaded images alongside customer written feedback describing product experiences), traditional machine learning approaches often process each modality in isolation, making it challenging to exploit the rich complementary information inherent across modalities. Vision-Language Models (VLMs) Radford et al. [2021], Jia et al. [2021], Li et al. [2022b] address this limitation by training a unified model on large-scale image-text pairs, enabling the formation of rich cross-modal semantic representations. These learned associations subsequently drive performance across diverse downstream applications, including image captioning, visual question answering, content moderation, context-sensitive retrieval, and multimodal sentiment analysis Du et al. [2022].

Although pretrained VLMs excel at zero-shot classification—enabling inference on new data without additional training, a promising research direction lies in leveraging these large-scale pretrained models to guide smaller task-specific classifiers tailored to particular applications. A seminal effort in this area is the GALS framework Petryk et al. [2022], which proposes using language-based guidance from VLMs to steer task-specific classifiers through the generation of

reference or guiding attention maps. Specifically, GALS transforms conventional classification labels into descriptive textual prompts, constructing semantically rich image-text pairs compatible with the VLM's multimodal processing pipeline. These pairs are used to extract class-focused image attention maps, which highlight image regions semantically aligned with the corresponding linguistic descriptions. With these VLM-generated attention maps, a regularization mechanism is designed for the task-specific classifier, constraining its focus to align with semantically meaningful regions identified through language-guided reasoning. This approach enhances the interpretability, robustness, and accuracy of image classifiers by clarifying visual-textual relationships and reducing reliance on irrelevant visual cues.

While GALS improves classification accuracy by integrating pretrained VLMs as a guidance mechanism, it inherits a fundamental limitation shared by most conventional VLMs: their reliance on deterministic embeddings. Such embeddings struggle to account for the inherent ambiguities and uncertainties in multimodal data Blundell et al. [2015], Oh et al. [2019]. Specifically, this framework overlooks the fact that a single textual description can correspond to multiple distinct images, just as a single image may be accurately described by various textual expressions. This constraint hinders model effectiveness, particularly in tasks requiring fine-grained distinctions or generalization across diverse data distributions. To address these challenges, probabilistic VLMs have been developed Upadhyay et al. [2023], Baumann et al. [2024], focusing on uncertainty estimation. Probabilistic embeddings provide a more nuanced representation by representing multimodal data as probability distributions rather than fixed points in the embedding space. Explicitly modelling uncertainty enables these frameworks to capture the variability in visual-textual relationships, leading to significant advancements in fine-grained classification, active learning, domain adaptation, and targeted model refinement via uncertainty quantification.

To accommodate multivaluedness and the ill-posed nature inherent to cross modal mappings, we introduce **PARIC**: *"a probabilistic framework for guiding visual attention using language specifications"*. PARIC leverages probabilistic VLMs to generate uncertainty-aware attention maps from descriptive textual prompts derived from classification labels. These probabilistic attention maps explicitly guide a separate task-specific neural classifier to focus on semantically meaningful visual regions, aiding model interpretability, robustness, and generalization. Building upon GALS Petryk et al. [2022] and probabilistic adapters Upadhyay et al. [2023], our approach bridges the gap between deterministic and probabilistic multimodal paradigms. We demonstrate improvements in accuracy, robustness, and interpretability across benchmark visual classification datasets, underscoring the efficacy of uncertainty-aware multimodal guidance.

## 2    Background

Vision-Language Models (VLMs) have been pivotal in multimodal research, enabling simultaneous processing of images and text within a unified architecture. Foundational models like CLIP Radford et al. [2021] and ALIGN Jia et al. [2021] employ contrastive learning on large-scale datasets (e.g., LAION-400M Schuhmann et al. [2021]), to align image-text pairs in a shared embedding space, excelling in tasks such as cross-modal retrieval and zero-shot classification Wang et al. [2024]. Unified pre-training approaches such as BLIP Li et al. [2022b], combine contrastive and generative objectives, enhancing performance in tasks like image captioning and visual question answering. However, a notable limitation of these models lies in their reliance on deterministic embeddings. Although suitable for many downstream applications, deterministic representations may not adequately address the uncertainties or variations inherent in real-world data, especially when dealing with fine-grained distinctions, unseen categories, or tasks that demand high precision and adaptability Yang et al. [2021].

**Probabilistic Embeddings**    have recently been explored in various workflows Upadhyay et al. [2023], Chun et al. [2021], Li et al. [2022a], Neculai et al. [2022], and hold the potential to address the inherent ambiguity and multivaluedness in VLMs by explicitly modeling uncertainty, offering a more nuanced representation of complex multimodal data Blundell et al. [2015], Oh et al. [2019]. However, training probabilistic VLMs from scratch demands substantial computational resources and large scale datasets Stirn et al. [2023]. Recent advancements address this limitation through post-hoc probabilistic methods. Specifically, Bayesian posterior approximation, as seen in Baumann et al. [2024], estimates uncertainty over the final VLM layers without retraining. Alternatively, Upadhyay et al. [2023] introduces lightweight probabilistic adapters to convert deterministic embeddings into probabilistic embeddings, offering a computationally efficient approach that leverages pre-trained models.

**Language as a Guide for Visual Models**    based on attention mechanisms allow models to focus selectively on the most salient parts of an input, such as specific regions in an image Gan et al. [2024]. In multimodal settings, language-guided attention facilitates the alignment of visual regions with semantic textual cues, enabling the generation of more robust and interpretable attention maps. Using language to direct attention, models can prioritize critical visual features while suppressing irrelevant background information, thereby reducing the risk of learning spurious correlations. For example, this approach ensures that a model associates a bird's species with its defining characteristics
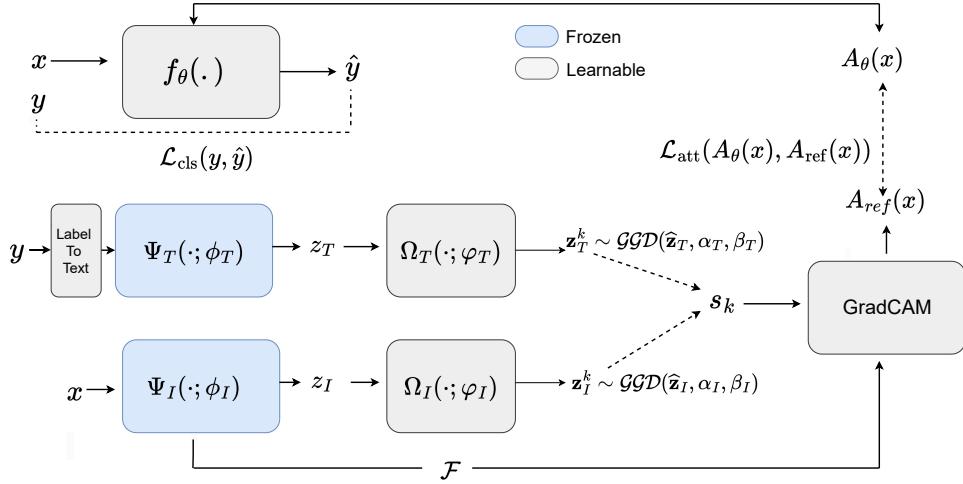
Figure 1: **PARIC Workflow.** The classifier $f_\theta$ predicts the label $\hat{y}$ and generates an attention map $A_\theta(x)$ for input image $x$. The ProbVLM pipeline expects image-text pairs, so the label $y$ is first converted into text prompts. These text prompts, along with the image, $x$, are processed by frozen CLIP encoders $\Psi_T$ (text) and $\Psi_I$ (image) to produce deterministic embeddings $\mathbf{z}_T$ and $\mathbf{z}_I$. Trainable adapters $\Omega_T$ and $\Omega_I$ model these embeddings as Generalized Gaussian Distributions (GGDs), from which $K$ samples are drawn to compute similarity scores. Grad-CAM combines these scores with CLIP's image feature map $\mathcal{F}$ to generate $K$ saliency maps, which are aggregated using mean or median into a reference map $A_{\text{ref}}(x)$. This map guides $A_\theta(x)$ via $\mathcal{L}_{\text{att}}$, complementing $\mathcal{L}_{\text{cls}}$ to improve the robustness and interpretability of the classifier.

rather than unrelated contextual factors like its habitat. Recent frameworks, such as GALS Petryk et al. [2022], have demonstrated the efficacy of modulating visual attention using language cues—a principle that also forms the foundation of this work.

**Information Grounding in VLMs**  refers to the process of associating textual descriptions with specific regions or attributes of visual inputs Yao et al. [2024], Yarom et al. [2023], enabling models to better understand the relationship between text and fine-grained visual features for improved accuracy and contextual awareness in predictions. This is typically achieved using pre-trained VLMs like CLIP or BLIP, which align visual and textual data in a shared embedding space to inform downstream tasks such as classification, segmentation, or retrieval. However, grounding becomes challenging in scenarios with mismatches between linguistic and visual modalities Kamath et al. [2023], particularly in fine-grained tasks where subtle distinctions are crucial, highlighting the need for robust mechanisms to ensure semantically meaningful correspondences between text and image.

**Supervising Attention in Visual Models**  Attention mechanisms plays an important role in identifying the most relevant regions of an image for a specific task that the model should focus on during training Selvaraju et al. [2019]. By explicitly supervising attention, models can be guided to prioritize the most relevant input areas, leading to enhanced performance on tasks requiring fine-grained discrimination. Approaches such as the "Right for the Right Reasons" Ross et al. [2017] enforce alignment between attention and task-specific requirements, ensuring the model focuses on appropriate visual features rather than spurious correlations. This supervision can be further strengthened by leveraging saliency maps derived from pre-trained vision-language models, which highlight semantically meaningful regions in images. In frameworks like GALS, attention supervision is integrated with language-guided cues to create robust models capable of handling both high-level task specifications and detailed visual features, improving interpretability and task accuracy.

## 3   Problem Formulation

Let $\mathcal{X} \in \mathbb{R}^{H \times W \times C}$ be the input space of images, where $H$ and $W$ denote spatial dimensions and $C$ the number of channels, and let $\mathcal{Y}$ be be the set of class labels. Consider a neural network classifier $f_\theta : \mathcal{X} \to \mathcal{Y}$, parameterized by $\theta$, which generates spatial attention map $A_\theta(x) \in \mathbb{R}^{H \times W}$ for each image $x \in \mathcal{X}$. Our goal is to regularize $A_\theta(x)$ by

leveraging a probabilistic attention prior derived from a frozen CLIP encoder, realized through probabilistic adapters akin to ProbVLM.

**Probabilistic Encoder Integration**   Pre-trained zero-shot models like CLIP are adapted to generate probabilistic embeddings, modeling image-text alignments as random variables rather than fixed points. With probabilistic adapters, any image or caption $\zeta$ is mapped to a random variable $\mathbf{z}$ following a Generalized Gaussian Distribution (GGD), parameterized by the output of the probabilistic adapter $\Omega(\cdot; \varphi)$ applied to the frozen CLIP encoder $\Psi(\cdot; \phi)$, where $\varphi$ and $\phi$ denote the parameters of the adapter and the CLIP encoder, respectively. Formally, for any image or caption $\zeta$, its embedding $\mathbf{z}$, the corresponding probability density function (PDF), and the parameters of the distribution are given by:

$$\mathbf{z} \sim \mathcal{GGD}(\widehat{\mathbf{z}}, \alpha, \beta), \tag{1}$$

$$p(\mathbf{z}; \widehat{\mathbf{z}}, \alpha, \beta) \propto \exp\big(-|\mathbf{z} - \widehat{\mathbf{z}}|^{\beta}/\alpha^{\beta}\big), \tag{2}$$

$$(\widehat{\mathbf{z}}, \alpha, \beta) = [\Omega(\cdot; \varphi) \circ \Psi(\cdot; \phi)](\zeta), \tag{3}$$

where $\mathcal{GGD}(\cdot)$ denotes a generalized Gaussian distribution, $\widehat{\mathbf{z}} \in \mathbb{R}^d$ is the mean embedding with dimension $d$ and $\alpha, \beta \in \mathbb{R}^+$ capture scale and shape of the embedding distribution. This design reflects both aleatoric (via $\alpha, \beta$) and epistemic (via Monte Carlo dropout) uncertainties. The adapters use intra-modal and cross-modal alignment objectives to ensure consistent image-text embeddings.

**Sampling from Probabilistic Embeddings**   To account for the uncertainty inherent in the multimodal task, we sample $K$ points from the learned image and text probabilistic distributions, i.e., $\mathbf{z}_I^{(k)} \sim \mathcal{GGD}(\widehat{\mathbf{z}}_I, \alpha_I, \beta_I)$ and $\mathbf{z}_T^{(k)} \sim \mathcal{GGD}(\widehat{\mathbf{z}}_T, \alpha_T, \beta_T)$. Each sample represents a different instantiation of the embeddings, reflecting distinct possible interpretations of the image-text pair.

**Saliency Map Generation and Aggregation**   To localize semantically relevant image regions, we generate saliency maps from $K$ sampled probabilistic embeddings. For each sample $k$, we compute a similarity score $s^{(k)} = \mathbf{z}_I^{(k)} \cdot \mathbf{z}_T^{(k)}$ between the image and text embeddings. Using this score and the image encoder's final convolutional feature map $\mathcal{F} \in \mathbb{R}^{H \times W \times C}$, we derive attention maps via Grad-CAM Selvaraju et al. [2017]. Channel-wise weights $\{w_c^{(k)}\}_{c=1}^{C}$ are obtained by spatial averaging of the gradients $\partial s^{(k)}/\partial \mathcal{F}$. The saliency map $A^{(k)}(x)$ is then constructed as a weighted sum of feature map channels, with ReLU activation ensuring non-negative saliency values: $A^{(k)}(x) = \text{ReLU}\left(\sum_{c=1}^{C} w_c^{(k)} \mathcal{F}_c\right)$. This results in a set of $K$ saliency maps $\{A^{(k)}(x)\}_{k=1}^{K}$, which are consolidated into a robust reference attention map $A_{\text{ref}}(x)$. To achieve this consolidation, we explore two distinct aggregation strategies: mean aggregation, which produces a smooth central tendency, and median aggregation, which provides an outlier-resistant summary of the sampled maps. Consequently, we define two PARIC model variants:

$$A_{\text{ref}}(x) = \begin{cases} \frac{1}{K} \sum_{k=1}^{K} A^{(k)}(x) & \text{(PARIC mean)} \\ \text{median}(\{A^{(k)}(x)\}_{k=1}^{K}) & \text{(PARIC median)} \end{cases}. \tag{4}$$

**Guiding the Visual Classifier**   To promote alignment between the classifier's attention map $A_\theta(x)$ and the reference attention map $A_{\text{ref}}(x)$, we introduce an attention regularization term:

$$\mathcal{L}_{\text{att}}(A_\theta(x), A_{\text{ref}}(x)) = \sum_{i,j} \big|(1 - A_{\text{ref}}(x)_{i,j}) \cdot A_\theta(x)_{i,j}\big|, \tag{5}$$

where $(1 - A_{\text{ref}}(x))$ represents regions deemed irrelevant by the reference attention map, thus discouraging the classifier from attending to these areas. We combine this with the standard cross-entropy classification loss:

$$\mathcal{L}_{\text{cls}}(f_\theta(x), y) = -\sum_{c=1}^{C} \mathbf{1}\{c = y\} \log p_\theta(c \,|\, x), \tag{6}$$

where $p_\theta(c \mid x)$ is the predicted probability for class $c$. The final objective function for PARIC is then formulated as a weighted sum of the classification and the attention regularization loss:

$$\mathcal{L}_{\text{total}}(x, y) = \mathcal{L}_{\text{cls}}(f_\theta(x), y) + \lambda \mathcal{L}_{\text{att}}(A_\theta(x), A_{\text{ref}}(x)), \tag{7}$$

where $\lambda > 0$ balances classification accuracy against attention alignment. Minimizing $\mathcal{L}_{\text{total}}$ guides the classifier towards stable, high-importance regions identified via the probabilistic embeddings, resulting in more interpretable and robust predictions. Figure 1 provides a schematic overview of the PARIC model.

# 4 Experiments

We compare PARIC with the well-established GALS framework across three diverse datasets (MS-COCO Lin et al. [2014b], Waterbird Sagawa* et al. [2020], Food101 Bossard et al. [2014]) over five independent trials to evaluate whether integration of probabilistic encoders leads to improved classification accuracy and better generalization performance in task-specific classifiers. We assess two variants of PARIC— *PARIC-mean* and *PARIC-median* which differ in their aggregation strategies for combining attention maps derived from sampled embeddings. We outline the specific implementation details below.

**Model Architecture**    We adopt CLIP with a ResNet50 backbone, pre-trained on ImageNet, as the foundation of our vision-language model. CLIP is selected for its excellent performance in aligning vision and language representations, as well as its architectural compatibility with the existing GALS framework Petryk et al. [2022]. This compatibility ensures a consistent baseline, allowing us to isolate and evaluate the effects of incorporating probabilistic layers. For the downstream classification task, we use the same image classifier architecture as in the GALS framework Petryk et al. [2022].

**Probabilistic Adapters**    We augment CLIP's deterministic encoders with post-hoc probabilistic adapters, implemented as small MLPs with dropout layers following the structure outlined in Upadhyay et al. [2023]. These adapters transform embeddings into parameters of a Generalized Gaussian Distribution, enabling uncertainty modeling in image-text alignments.

**Attention Map Generation**    To extract attention via language specification, we employ prompts similar to those used in CLIP (e.g., "an image of *category*" or "a photo of *category*"), where *category* corresponds to task-relevant concepts. For each image-text pair, 50 embeddings are sampled from the learned probabilistic distribution. These embeddings are processed through GradCAM Selvaraju et al. [2017], applied at the final convolutional layer (Layer 4) of the ResNet-50 CLIP encoder, to produce individual saliency maps highlighting the most relevant image regions. The resulting saliency maps are aggregated using either mean or median pooling to create a representative attention map for each image-text pair with uncertainty estimates, which can be calculated by pixel-wise standard deviation of the map samples. a These aggregated maps are then used to regularize the attention mechanisms of task-specific classifiers with the aim of ensuring alignment with semantically meaningful regions.

**Implict vs Explicit Dataset Bias**    We examine two variants of the Waterbirds dataset to study explicit bias. The first, Waterbirds-$100\%$, enforces perfect correlation between bird type and background during training (e.g., waterbirds on water, landbirds on land). The second, Waterbirds-$95\%$, introduces $5\%$ of training samples that break this correlation, testing robustness to slight bias deviations. For dataset details, see Sagawa* et al. [2020].

The Food-101 dataset is used to study implicit biases from its uncurated training set, which includes noisy labels, co-occurring elements, and visual artifacts. For instance, certain ingredients (e.g., sauces) are spuriously correlated with specific dishes, introducing biases. The evaluation set, however, is curated, creating a distribution shift. We use two setups: *Red Meat Subset*: following Petryk et al. [2022], we use a five-way classification task (baby back ribs, filet mignon, pork chop, prime rib, steak) to predict red meat categories. *Meat Subset*: We introduce a three-way classification task (red meat, white meat, fish) by filtering 50 animal-based meat classes, analyzing biases like co-occurring ingredients and image quality variations.

For MS-COCO Lin et al. [2014a], we extract a subset of images labeled with the "Person" category, we further restrict this subset to images included in the MS-COCO-ApparentGender dataset, which introduces implicit bias through gender labels inferred from captions ("Man," "Woman," or "Person"). This filtering reduces the size of the train, validation, and test sets, allowing us to study the impact of implicit bias on classification performance. Implicit bias in the MSCOCO-ApparentGender setup stems from the gender labels which are based on outward appearances described in image captions rather than objective features of individuals. This reliance on subjective descriptions introduces societal stereotypes into the dataset, where certain activities, clothing, or contexts are disproportionately associated with specific genders (e.g., sports with men or domestic settings with women). Additionally, MS-COCO itself is inherently biased, having an unbalanced $1:3$ women to men ratio Tang et al. [2021]. Such under-representation amplifies the risk that the model learns spurious correlations from majority-class features, favoring men over women in its predictions. This imbalance, combined with the subjective nature of the apparent gender labels and the underuse of neutral terms like "Person," aggravates the model's tendency to internalize and reinforce societal stereotypes, influencing its performance across gender categories.

5

### 4.1 Results

**Explicit Bias on Waterbirds** On the Waterbirds-100% dataset, *PARIC-Median* achieves an overall accuracy of (96.80%), followed by *PARIC-Mean* (96.69%) and GALS (96.65%) (Table 1a). GALS achieves the highest accuracy for the Waterbird class (95.67% ± 0.71), but also exhibits higher variance compared to both *PARIC-Mean* (±0.33) and *PARIC-Median* (±0.66). For the Landbird class, *PARIC-Mean* achieves the best accuracy (97.08% ± 0.50), with lower variance (57% reduction) compared to GALS (96.81% ± 1.16). These results highlight *PARIC*'s ability to reduce variance and improve stability, particularly in the presence of strong background-label correlations.

On the more challenging Waterbirds-95% dataset, where 5% of training samples break the background-label correlation, GALS achieves the highest overall (96.93%) accuracy. *PARIC-Mean* achieves the best Waterbird accuracy (94.66% ± 0.87), while *PARIC-Median* demonstrates the lowest variance for the Waterbird class (±0.32) and overall (±0.15). Table 1b indicates that *PARIC-Median* is more robust to outliers and provides more stable predictions, even when the dataset contains counterexamples to the dominant bias. Notably, both GALS and PARIC do not perform worse for the Waterbird class on Waterbird-100 which is more challenging, as compared to Waterbird-95. This could be attributed to a combination of regularization via noise, a greater diversity of training data samples in Waterbirds-100 or the noise affecting non-critical regions of the test images, which the attention mechanism is able to mitigate.

Overall, *PARIC* demonstrates consistent performance across both datasets, with *PARIC-Median* showing particular strength in reducing variance and improving robustness. GALS performed well on the Landbird class in the Waterbirds-95% dataset, and Waterbird class in the Waterbirds-100% dataset indicating gains due to more model expressiveness (and less regularization).

(a) Waterbirds 100%

| Method | Waterbird | Landbird | Overall |
|---|---|---|---|
| GALS | **95.67** ± 0.7101 | 96.81 ± 1.1617 | 96.65 ± 0.9319 |
| PARIC mean | 95.04 ± **0.3292** | **97.08** ± **0.5003** | 96.69 ± 0.5195 |
| PARIC median | 95.24 ± 0.6620 | 96.92 ± 0.6690 | **96.80** ± **0.5059** |

(b) Waterbirds 95%

| Method | Waterbird | Landbird | Overall |
|---|---|---|---|
| GALS | 93.69 ± 0.6691 | **97.46** ± 0.2487 | **96.93** ± 0.1870 |
| PARIC mean | **94.66** ± 0.8683 | 96.90 ± 0.5459 | 96.59 ± 0.3582 |
| PARIC median | 94.13 ± **0.3220** | 97.35 ± **0.2180** | 96.91 ± **0.1513** |

Table 1: PARIC vs. GALS: Classification accuracy on (a) Waterbirds 100% and (b) Waterbirds 95% test sets (mean ± std. dev.).

| Method | Man | Woman | Overall | Outcome Divergence |
|---|---|---|---|---|
| GALS | 73.72 ± 8.9108 | 64.15 ± 8.0289 | 68.93 ± 2.0765 | 0.0729 ± 0.0406 |
| PARIC mean | 75 ± 5.6892 | **65.42** ± 5.5482 | **70.21** ± 1.8194 | **0.0597** ± 0.0249 |
| PARIC median | **75.21** ± 4.7765 | 63.93 ± **2.7205** | 69.57 ± 2.8968 | 0.0642 ± **0.0119** |

Table 2: **MSCOCO.** Classification accuracy of PARIC and the baseline GALS on the MS-COCO test set. Highlighted in **bold** are the best results in each column (mean ± std. dev.).

**Implicit bias on COCO Gender.** On the COCO Gender dataset, PARIC outperforms the baseline GALS in both overall accuracy and fairness (Table Table 2). *PARIC-Mean* achieves the highest overall accuracy (70.21%), surpassing GALS (68.93%), and improves classification accuracy for the underrepresented "Woman" class to 65.42%. PARIC also demonstrates greater stability and consistency, with *PARIC-Mean* achieving the lowest standard deviation for overall accuracy (±1.82) and *PARIC-Median* minimizing outcome divergence (Jensen-Shannon divergence Lin [1991] calculated over the score distributions of the two classes; lower values indicate better performance) variability to 0.0119, compared to GALS (0.0406). For the "Woman" class, *PARIC-Median* further reduces variability to (±2.72), ensuring

6

more predictable and trustworthy outcomes. Additionally, PARIC reduces gender disparities, with *PARIC-Mean* achieving a lower outcome divergence (0.0597) compared to GALS (0.0729), indicating more balanced performance across gender categories.

**Robustness to noisy data**   To evaluate PARIC's robustness, we assess its performance on the Food-101 dataset containing implicit biases and image noise. We test the model on two setups: (i) the Red Meat Subset with five classes and (ii) the Meat Dataset with three classes. These tasks assess the model's ability to handle challenging noisy scenarios while accounting for biases.

The experimental results for the Red Meat Subset (Table 3a) demonstrate that GALS achieves higher accuracy for 3 classes (Filet Mignon, Pork Chop, Prime Rib), while PARIC-mean achieves lower variance compared to PARIC-median and GALS. In the Steak category, GALS achieves an accuracy of $50.4\%$, which PARIC-mean improves this to $55.6\%$ while significantly reducing variability ($\pm 2.0861$). Across all categories, PARIC-mean exhibits lower variance in accuracy scores, enhanced model stability.

On the Meat Dataset (Table 3b), PARIC demonstrates its ability to handle class imbalance effectively. With class distributions of $38\%$ Fish, $36\%$ Red Meat, and $26\%$ White Meat, PARIC-median reduces prediction variability for the minority class (White Meat) from 0.5824 (GALS) to 0.2486. While PARIC shows slightly higher standard deviations for majority classes (Red Meat and Fish), it maintains higher overall accuracy ($85.90\%$ vs. $85.78\%$ for GALS). This suggests that PARIC explores a broader range of feature representations.

(a) Red Meat Dataset

| Method | Baby Back Ribs | Filet Mignon | Pork Chop | Prime Rib | Steak | Overall |
|---|---|---|---|---|---|---|
| GALS | $86.56 \pm 1.1482$ | $\mathbf{72.4} \pm 1.9919$ | $\mathbf{71.2} \pm 1.9758$ | $\mathbf{85.92} \pm 1.1142$ | $50.4 \pm 3.4409$ | $73.30 \pm \mathbf{0.4422}$ |
| PARIC mean | $\mathbf{87.84} \pm \mathbf{0.8616}$ | $69.52 \pm \mathbf{0.7756}$ | $70.2 \pm \mathbf{0.9666}$ | $84.32 \pm \mathbf{0.6881}$ | $55.6 \pm \mathbf{2.0861}$ | $\mathbf{73.50} \pm 0.4837$ |
| PARIC median | $86.32 \pm 1.5676$ | $70.96 \pm 2.2712$ | $69.68 \pm 1.0552$ | $84.40 \pm 0.9797$ | $\mathbf{55.68} \pm 2.4709$ | $73.41 \pm 0.8654$ |

(b) Meat Dataset

| Method | White meat | Red meat | Fish | Overall |
|---|---|---|---|---|
| GALS | $85.37 \pm 0.5824$ | $81.48 \pm \mathbf{0.2889}$ | $89.12 \pm \mathbf{0.3978}$ | $85.78 \pm 0.2504$ |
| PARIC mean | $\mathbf{85.72} \pm 0.2608$ | $81.15 \pm 0.5898$ | $\mathbf{89.33} \pm 0.4999$ | $\mathbf{85.90} \pm 0.2352$ |
| PARIC median | $84.84 \pm \mathbf{0.2486}$ | $\mathbf{81.49} \pm 0.6732$ | $89.15 \pm 0.8396$ | $85.61 \pm \mathbf{0.1469}$ |

Table 3: Comparison of classification accuracy of PARIC and the baseline GALS on (a) Read Meat Dataset; (b) Meat Dataset. Highlighted in **bold** are the best results in each column (mean $\pm$ std. dev.).

**Attention maps visualization**   The attention maps in Figs. 2, 3a and 4 visually demonstrate how probabilistic embeddings enhance model interpretability. PARIC's attention maps focus on semantically meaningful regions, reducing variability compared to deterministic approaches. For example, in the Waterbirds dataset, PARIC captures the entire bird even when background cues conflict with class labels. Similarly, in the COCO dataset, PARIC produces broader, more context-aware attention maps that focus on the person rather than disjointed body parts, mitigating overfitting and improving consistency. Despite these advantages, the proposed framework has limitations. In some instances like the example attention map for Food-101 (Figure 2 and Fig. 3b, bottom two rows), the regularization can be too strong, limiting model expressiveness. Both guided approaches also rely on CLIP's deterministic embeddings as a foundation, meaning the performance depends on the quality of CLIP's representations. If CLIP fails to capture relevant features, GALS and PARIC may also struggle.

## 5   Discussion

PARIC demonstrates improvements owing to the probabilistic approach in specific cases involving imbalanced and noisy datasets, by reducing variability and providing more consistent and equitable predictions. In the Waterbirds dataset, PARIC effectively reduces the influence of background correlations, enabling the model to focus on the bird's features. This leads to more robust generalization, especially in challenging scenarios where spurious biases may mislead predictions. In the COCO Gender dataset, PARIC mitigates gender biases, offering more balanced results
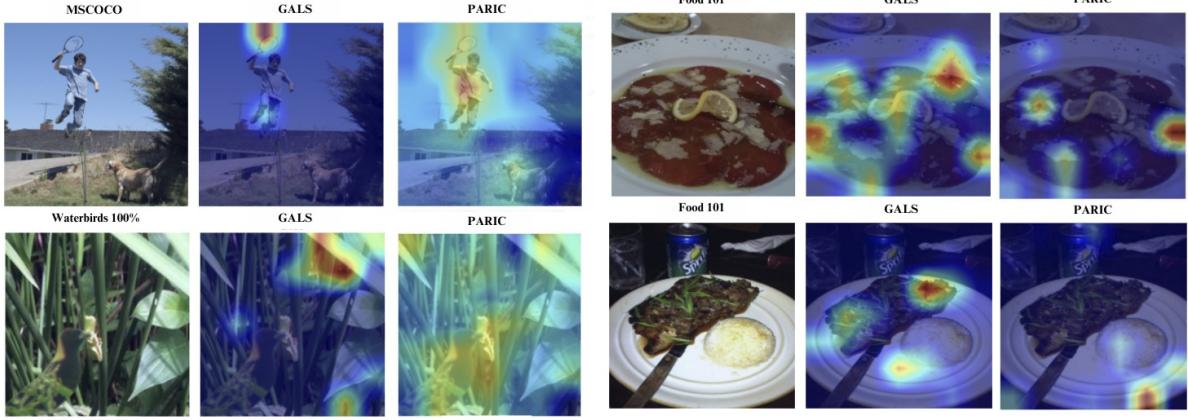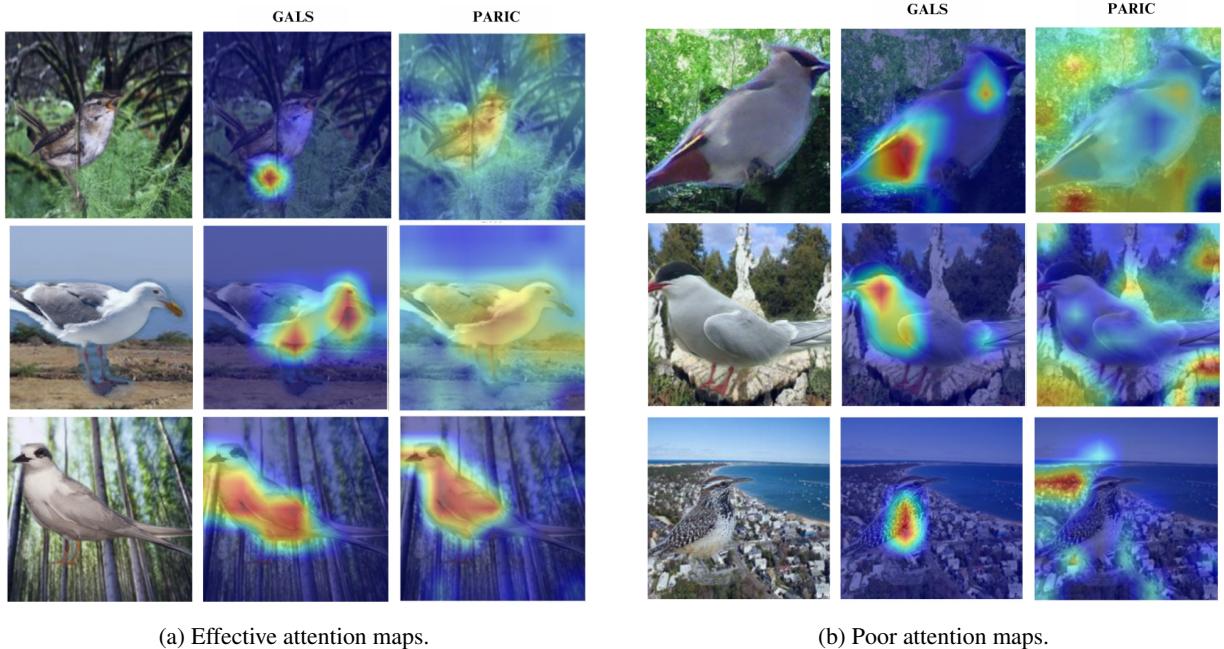
Figure 2: *Attention Map Visualization.* for three instances: COCO and Waterbirds 100% (GALS vs. PARIC Mean) and Food-101 (GALS vs. PARIC Median). Each row shows the original image, the attention map from frozen CLIP, and the refined map after integrating probabilistic layers. The first two instances show the strengths of the probabilistic approach, where the attention maps are more accurate, while the case of Food-101 shows an experiment where PARIC performs worse, with the regularization being too strong and limiting.



(a) Effective attention maps.                                          (b) Poor attention maps.

Figure 3: **Comparison of Effective vs. Poor Attention Maps on the Waterbirds 100% Dataset.** Each subfigure presents three instances from the dataset: the first image is the original input, and the third image is the final attention map obtained after integrating probabilistic layers. In (a), the first two rows use PARIC Mean and the third uses PARIC Median, whereas in (b), the first row uses PARIC Mean and the subsequent rows use PARIC Median.

across classes and reducing variability. Leveraging its probabilistic embeddings, PARIC increases stability and fairness, making it a reliable framework for biased or imbalanced datasets.

In PARIC, mean aggregation provides strong overall performance by smoothing predictions across runs, leading to improved generalization. However, it does not account for potential variability in predictions caused by ambiguous or biased data. Median aggregation, on the other hand, is better suited for mitigating extreme variations, resulting in more stable and consistent predictions. This stability is particularly beneficial in datasets like COCO Gender, where fairness metrics are critical for underrepresented classes.
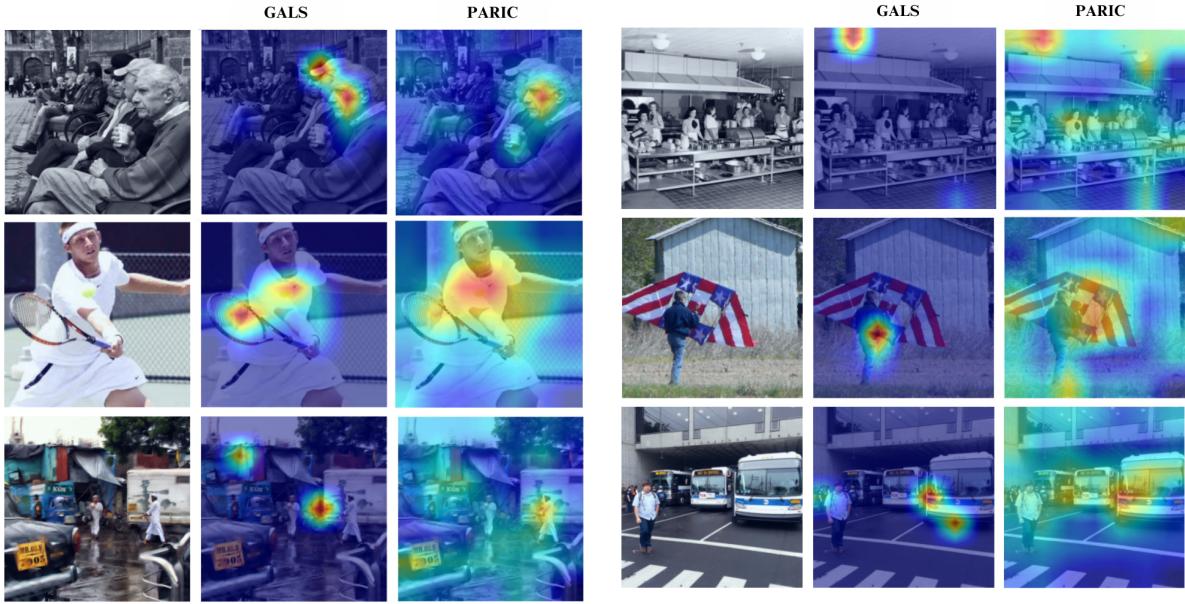
8

Figure 4: *Attention maps for MSCOCO*. In each row, the first image represents the original input, followed by the attention map from the frozen CLIP model. The third image shows the attention map obtained after integrating probabilistic layers and sampling 50 embedding using the mean aggregation method.

Integrating uncertainty-aware aggregation into PARIC could further harness the strengths of its probabilistic embeddings. Dynamically weighting predictions based on confidence or uncertainty measures can enable PARIC to better adapt to noisy or ambiguous data. For instance, in a dataset prone to inherent subjectivity, such as COCO Gender, incorporating uncertainty measures could help the model identify and prioritize ambiguous cases, further improving fairness and robustness. Additionally, as future work, increasing the number of samples generated from probabilistic encoders could refine the attention maps. More accurate and detailed attention maps would offer better guidance for the classifier, potentially enhancing both interpretability and predictive performance. This approach would allow PARIC to fully leverage its probabilistic nature for improved consistency, fairness, and overall reliability.

## 6 Conclusions

This work introduces a natural language-guided attention framework leveraging probabilistic embeddings adapted from pre-trained vision-language models such as CLIP. The probabilistic treatment of embeddings allows treatment of intrinsic multivaluedness and ill-posedness of cross-modal mappings. The proposed probabilistic framework of guiding visual attention through language specifications is validated on three challenging benchmark test problems consisting of noise, implicit bias and class imbalance, where it delivers improved classification performance and model stability in a majority of cases.

### Acknowledgments

### References

Anton Baumann, Rui Li, Marcus Klasson, Santeri Mentu, Shyamgopal Karthik, Zeynep Akata, Arno Solin, and Martin Trapp. Post-hoc probabilistic vision-language models. *arXiv*, 2024.

Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural network. In *International conference on machine learning*, pages 1613–1622. PMLR, 2015.

Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101 – mining discriminative components with random forests. In *European Conference on Computer Vision*, 2014.

Sanghyuk Chun, Seong Joon Oh, Rafael Sampaio de Rezende, Yannis Kalantidis, and Diane Larlus. Probabilistic embeddings for cross-modal retrieval. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8411–8420, 2021. URL https://api.semanticscholar.org/CorpusID:231592523.

Yifan Du, Zikang Liu, Junyi Li, and Wayne Xin Zhao. A survey of vision-language pre-trained models. In *International Joint Conference on Artificial Intelligence*, 2022. URL https://api.semanticscholar.org/CorpusID:247026006.

Chenquan Gan, Xiang Fu, Qingdong Feng, Qingyi Zhu, Yang Cao, and Ye Zhu. A multimodal fusion network with attention mechanisms for visual–textual sentiment analysis. *Expert Systems with Applications*, 242:122731, 2024. ISSN 0957-4174. doi: https://doi.org/10.1016/j.eswa.2023.122731. URL https://www.sciencedirect.com/science/article/pii/S0957417423032335.

Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc V. Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *International Conference on Machine Learning*, 2021. URL https://api.semanticscholar.org/CorpusID:231879586.

Amita Kamath, Jack Hessel, and Kai-Wei Chang. What's "up" with vision-language models? investigating their struggle with spatial reasoning. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, Singapore, December 2023. Association for Computational Linguistics. URL https://aclanthology.org/2023.emnlp-main.568/.

Hao Li, Jingkuan Song, Lianli Gao, Pengpeng Zeng, Haonan Zhang, and Gongfu Li. A differentiable semantic metric approximation in probabilistic embedding for cross-modal retrieval. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS '22, Red Hook, NY, USA, 2022a. Curran Associates Inc. ISBN 9781713871088.

Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *ICML*, 2022b.

Jianhua Lin. Divergence measures based on the shannon entropy. *IEEE Transactions on Information theory*, 37(1): 145–151, 1991.

Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer vision–ECCV 2014: 13th European conference, zurich, Switzerland, September 6-12, 2014, proceedings, part v 13*, pages 740–755. Springer, 2014a.

Tsung-Yi Lin, Michael Maire, Serge J. Belongie, Lubomir D. Bourdev, Ross B. Girshick, James Hays, Pietro Perona, Deva Ramanan, Piotr Doll'a r, and C. Lawrence Zitnick. Microsoft COCO: common objects in context. *CoRR*, abs/1405.0312, 2014b. URL http://arxiv.org/abs/1405.0312.

Andrei Neculai, Yanbei Chen, and Zeynep Akata. Probabilistic compositional embeddings for multimodal image retrieval. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 4546–4556, 2022. URL https://api.semanticscholar.org/CorpusID:248118624.

Seong Joon Oh, Andrew C. Gallagher, Kevin P. Murphy, Florian Schroff, Jiyan Pan, and Joseph Roth. Modeling uncertainty with hedged instance embeddings. In *International Conference on Learning Representations*, 2019. URL https://openreview.net/forum?id=r1xQQhAqKX.

Suzanne Petryk, Lisa Dunlap, Keyan Nasseri, Joseph Gonzalez, Trevor Darrell, and Anna Rohrbach. On guiding visual attention with language specification. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18071–18081, 2022. doi: 10.1109/CVPR52688.2022.01756.

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 8748–8763, 2021.

Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "why should i trust you?" explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1135–1144, 2016.

Andrew Slavin Ross, Michael C. Hughes, and Finale Doshi-Velez. Right for the right reasons: Training differentiable models by constraining their explanations. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI-17*, pages 2662–2670, 2017.

Shiori Sagawa*, Pang Wei Koh*, Tatsunori B. Hashimoto, and Percy Liang. Distributionally robust neural networks. In *International Conference on Learning Representations*, 2020. URL `https://openreview.net/forum?id=ryxGuJrFvS`.

C. Schuhmann, R. Vencu, R. Beaumont, R. Kaczmarczyk, C. Mullis, A. Katta, T. Coombes, J. Jitsev, and A. Komatsuzaki. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. *arXiv*, 2021.

R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. *International Journal of Computer Vision*, 128(2):336–359, 2019.

Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 618–626, 2017. doi: 10.1109/ICCV.2017.74.

Andrew Stirn, Harm Wessels, Megan Schertzer, Laura Pereira, Neville Sanjana, and David Knowles. Faithful heteroscedastic regression with neural networks. In *International Conference on Artificial Intelligence and Statistics*, pages 5593–5613. PMLR, 2023.

Ruixiang Tang, Mengnan Du, Yuening Li, Zirui Liu, Na Zou, and Xia Hu. Mitigating gender bias in captioning systems. In *Proceedings of the Web Conference 2021*, pages 633–645, 2021.

A Torralba and AA Efros. Unbiased look at dataset bias. In *Proceedings of the 2011 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1521–1528, 2011.

Uddeshya Upadhyay, Shyamgopal Karthik, Massimiliano Mancini, and Zeynep Akata. Probvlm: Probabilistic adapter for frozen vison-language models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1899–1910, 2023.

Tianshi Wang, Fengling Li, Lei Zhu, Jingjing Li, Zheng Zhang, and Heng Tao Shen. Cross-modal retrieval: A systematic review of methods and future directions. *Proceedings of the IEEE*, 112(11):1716–1754, 2024.

Gengcong Yang, Jingyi Zhang, Yong Zhang, Baoyuan Wu, and Yujiu Yang. Probabilistic modeling of semantic ambiguity for scene graph generation. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12522–12531, 2021. URL `https://api.semanticscholar.org/CorpusID:232168713`.

Yuan Yao, Ao Zhang, Zhengyan Zhang, Zhiyuan Liu, Tat-Seng Chua, and Maosong Sun. Cpt: Colorful prompt tuning for pre-trained vision-language models. *AI Open*, 5:30–38, 2024.

Michal Yarom, Yonatan Bitton, Soravit Changpinyo, Roee Aharoni, Jonathan Herzig, Oran Lang, Eran Ofek, and Idan Szpektor. What you see is what you read? improving text-image alignment evaluation. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL `https://openreview.net/forum?id=j5AoleAIru`.