

IA 323 - MS DATA AI

Computer vision - EXAM

(3 hours)

24th January 2025

Please be advised that you are permitted to bring any relevant documents to the exam. The examination comprises 10 questions, each correct answer will give 2 points or more. Your objective is to thoroughly read the provided paper and formulate responses to the questions. Here's a suggested approach :

1. Begin by reviewing all the questions.
2. Proceed to read the paper, not focusing too much on the introduction and related works within the allocated time of 1 hour and 30 minutes.
3. Dedicate the following hour to answering the questions.
4. Spend the extra 30 minutes to review and correct your responses.

1 Question 1

Describe what is **BLIP** and how **BLIP** is working.

2 Question 2

Summarize the primary motivation behind revisiting language priors in vision-language models (VLMs).

3 Question 3

Explain the concept of the **Visual Generative Pre-Training Score (VisualGPSTScore)**. How is it computed? Be clear this question is important.

4 Question 4

Why does the paper argue for debiasing generative scores in image-text retrieval tasks?

5 Question 5

Explain the probabilistic derivation of $P_{test}(t|i)$ in equation (4)? Please also try to explain what is α equation (7)?

6 Question 6

Discuss how Monte Carlo sampling is used to estimate $P_{train}(t)$

7 Question 7

Discuss the limitations of blind language models in solving vision-language tasks and explain how enhanced VLMs as proposed in the paper can overcome those limitations?

8 Question 8

Suggest techniques for capturing VQA uncertainty

9 Question 9

Propose a practical application of VisualGPTScore outside of research benchmarks linked to multimodal image generation.

10 Question 10

How would you adapt the concepts in this paper to evaluate text-to-video generative models? Please try to take into account that a video is a sequence of Actions.

Revisiting the Role of Language Priors in Vision-Language Models

Zhiqiu Lin^{*1} Xinyue Chen^{*1} Deepak Pathak¹ Pengchuan Zhang² Deva Ramanan¹

1. Introduction

Vision-language models (VLMs) trained on web-scale datasets will likely serve as the foundation for next-generation visual understanding systems. One reason for their widespread adoption is their ability to be used in an “off-the-shelf” (OTS) or zero-shot manner without fine-tuning for specific target applications. In this study, we explore their OTS use on the task of image-text retrieval (e.g., given an image, predict the correct caption out of K options) across a suite of nine popular benchmarks.

Challenges. While the performance of foundational VLMs is impressive, many open challenges remain. Recent analyses (Kamath et al., 2023; Yuksekgonul et al., 2022) point out that leading VLMs such as CLIP (Radford et al., 2021) may often degrade to “bag-of-words” that confuse captions such as “the horse is eating the grass” and “the grass is eating the horse”. This makes it difficult to use VLMs to capture *compositions* of objects, attributes, and their relations. But somewhat interestingly, large-scale language models (LLMs) trained for autoregressive next-token prediction (Brown et al., 2020) seem to be able to discern such distinctions, which we investigate below. A related but under-appreciated difficulty is that of *benchmarking* the performance of visio-linguistic reasoning. Perhaps the most well-known example in the community is that of the influential VQA benchmarks (Antol et al., 2015), which could be largely solved by exploiting linguistic biases in the dataset – concretely, questions about images could often be answered by “blind” language-only models that did not look at the image (Goyal et al., 2017). Notably, we find that such blind algorithms still excel on many contemporary image-text retrieval benchmarks where VLMs may struggle.

Generative models for discriminative tasks. We tackle the above challenges by revisiting the role of language priors through a probabilistic lens. To allow for a probabilistic treatment, we focus on generative VLMs that take an image as input and stochastically generate text via next-token pre-

diction (Li et al., 2022; 2023). We first demonstrate that such models can be easily repurposed for discriminative tasks (such as retrieval) by setting the match score for an image-text pair to be the probability that the VLM would generate that text from the given image, or $P(\text{text}|\text{image})$. We call this probability score the Visual Generative Pre-Training Score, or VisualGPTScore. Computing the VisualGPTScore is even more efficient than next-token generation since given an image, all tokens from a candidate text string can be evaluated in parallel. Though conceptually straightforward, such an approach is not a common baseline. In fact, the generative VLMs (Li et al., 2022) that we analyze train *separate* discriminative heads for matching/classifying image-text pairs, but we find that their language generation head itself produces better scores for matching (since it appears to better capture compositions). Indeed, the OTS VisualGPTScore performs surprisingly well on many benchmarks, even producing near-perfect accuracy on ARO (Yuksekgonul et al., 2022). But it still struggles on other benchmarks such as Winoground (Thrush et al., 2022). We analyze this below.

The role of language priors. We analyze the discrepancy in performance across benchmarks from a probabilistic perspective. Our key insight is that many benchmark biases can be formalized as mismatching distributions over text between foundational pre-training data and benchmark test data – $P_{\text{train}}(\text{text})$ versus $P_{\text{test}}(\text{text})$. We use a first-principles analysis to account for distribution shift by simply reweighting the VisualGPTScore with the Bayes factor $P_{\text{test}}(\text{text})/P_{\text{train}}(\text{text})$, a process we call *debiasing*. To compute the Bayes reweighting factor, we need access to both the train and test language prior. We compute $P_{\text{train}}(\text{text})$ from an OTS VLM by drawing Monte-Carlo samples of $P_{\text{train}}(\text{text}|\text{image})$ from the trainset or Gaussian noise images. Because $P_{\text{test}}(\text{text})$ may require access to the test set, we explore practical variants that assume P_{test} is (a) identical to $P_{\text{train}}(\text{text})$, (b) uninformative/uniform, or (c) learnable from a small held-out valset. Our analysis helps explain the strong performance of the VisualGPTScore on certain benchmarks and its poor performance on others. Moreover, our analysis offers simple strategies to improve performance through debiasing without requiring any re-training. We conclude by showing a theoretical connection between debiasing and mutual information, which can be seen as a method for removing the effect of marginal priors

^{*}Equal contribution ¹CMU ²Meta. Correspondence to: Zhiqiu Lin <zhiqiu@andrew.cmu.edu>.

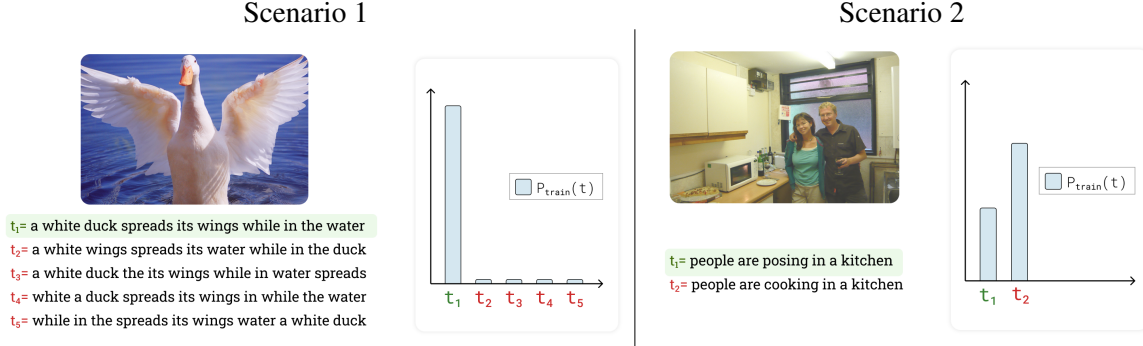


Figure 1. Two train-test shifts encountered in image-to-text retrieval tasks. Scenario 1 (left) constructs negative captions by shuffling words in the true caption (as in ARO-Flickr (Yuksekgonul et al., 2022)), but this produces implausible text such as “white a duck spreads its wings in while the water”. Here, exploiting the language bias of the training set will help since it will downweight the match score for such implausible negative captions. In fact, we show that a blind language-only model can easily identify the correct caption. Scenario 2 (right) constructs negative captions that are curated to be plausible (as in SugarCrep (Hsieh et al., 2023)). Here, the language bias of the training set may hurt, since it will prefer to match common captions that score well under the language prior; i.e., the incorrect caption of “people are cooking in a kitchen” is slightly more likely than the true caption of “people are posing in a kitchen” under the language prior, and so removing the language bias improves performance. We present simple training-free approaches for removing such language biases, and show this significantly improves performance on challenging benchmarks that fall into Scenario 2.

when computing joint probability scores.

Empirical analysis. We conduct a thorough empirical evaluation of the OTS VisualGPTScore (and its debiased variants) for open-sourced image-conditioned language models (Li et al., 2022; 2023; Liu et al., 2023) across nine popular vision-language benchmarks. We first point out that the VisualGPTScore by itself produces SOTA accuracy on certain benchmarks like ARO (Yuksekgonul et al., 2022) where their inherent language biases help remove incorrect captions that are also unnatural (such as “a white duck the its wings while in water” as shown in Fig. 1). In fact, we show that blind baselines also do quite well on these benchmarks, since language-only models can easily identify such implausible captions. However, such language biases do not work well on benchmarks where incorrect captions are carefully constructed to be realistic. Here, VisualGPTScore should be debiased so as not to naively prefer more common captions that score well under its language prior. Debiasing consistently improves performance on benchmarks such as Flickr30K (Young et al., 2014) and Winoground (Thrush et al., 2022). Interestingly, we find that debiasing can also improve accuracy on the *train* set used to learn the generative VLMs, indicating that such models learn biased estimates of the true conditional distribution $P_{train}(\text{text}|\text{image})$. Finally, our approach sets a new state-of-the-art on image-text alignment (Thrush et al., 2022; Wang et al., 2023), showing potential to replace the widely-used CLIPScore (Hessel et al., 2021) in text-to-image evaluation. In fact, our latest work (Lin et al., 2024; Li et al., 2024) extends VisualGPTScore to more pow-

erful vision-language models trained on visual-question-answering (VQA) data, achieving further improvements.

Contributions:

- We introduce VisualGPTScore to repurpose generative VLMs for discriminative (image-text retrieval) tasks.
- Our analysis shows that language priors play a key role in addressing train-test distribution shifts, leading to a zero-shot debiasing technique that significantly improves performance on challenging benchmarks.
- We find that many recent benchmarks for foundational VLMs like ARO can be largely solved by blind solutions (e.g., $P(\text{text})$) that ignore images. This underscores the need to reevaluate language priors in vision-language benchmarks.

2. The role of language priors

In this section, we present a simple probabilistic treatment for analyzing the role of language priors in image-conditioned language models (or generative VLMs). Motivated by their strong but inconsistent performance across a variety of image-text retrieval benchmarks, we analyze their behavior when there exists a mismatch between training and test distributions, deriving simple schemes for addressing the mismatch with reweighting. We emphasize that the training data that we refer to is the foundational pre-training dataset, while the test data is always a given benchmark dataset; in fact, most benchmarks we analyze do not even provide a trainset. We conclude by exposing a connection to related work on mutual information.

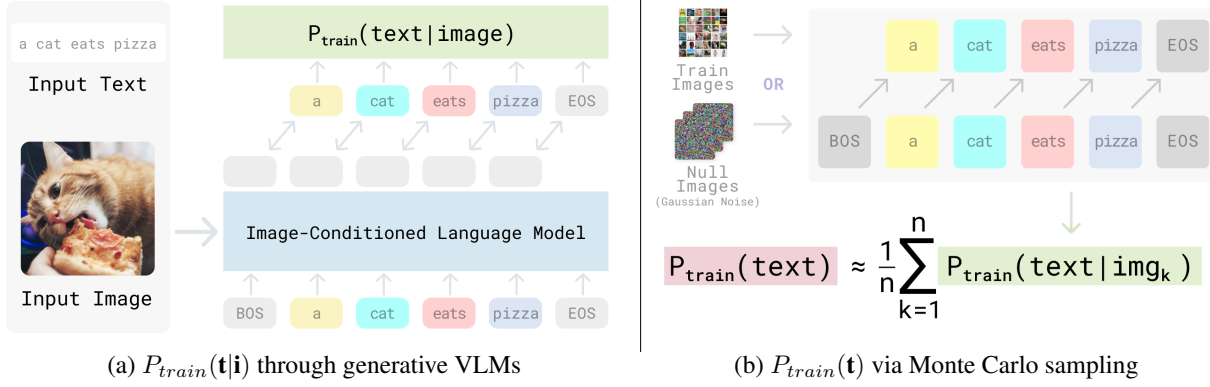


Figure 2. Estimating $P_{train}(\mathbf{t}|\mathbf{i})$ and $P_{train}(\mathbf{t})$ from generative VLMs. Figure (a) shows how image-conditioned language models such as Li et al. (2022) that generate text based on an image can be repurposed for computing $P_{train}(\mathbf{t}|\mathbf{i})$, which is factorized as a product of $\prod_{k=1}^m P(t_k|t_{<k}, \mathbf{i})$ for a sequence of m tokens. These terms can be efficiently computed in *parallel*, unlike *sequential* token-by-token prediction for text generation. Figure (b) shows two approaches for Monte Carlo sampling of $P_{train}(\mathbf{t})$. While the straightforward approach is to sample trainset images, we find that using “null” (Gaussian noise) images can also achieve robust estimates.

Computing $P(\mathbf{t}|\mathbf{i})$. To begin our probabilistic treatment, we first show that image-conditioned language models (that probabilistically generate text based on an image) can be repurposed for computing a score between a given image \mathbf{i} and text caption \mathbf{t} . The likelihood of a text sequence $\mathbf{t} = \{t_1, t_2, \dots, t_m\}$ conditioned on image \mathbf{i} is naturally factorized as an autoregressive product (Bengio et al., 2003):

$$P(\mathbf{t}|\mathbf{i}) = \prod_{k=1}^m P(t_k|t_{<k}, \mathbf{i}) \quad (1)$$

Image-conditioned language models return back m softmax distributions corresponding to the m terms in the above expression. Text generation requires *sequential* token-by-token prediction, since token t_k must be generated before it can be used as an input to generate the softmax distribution over token t_{k+1} . Interestingly, given an image \mathbf{i} and a text sequence \mathbf{t} , the above probability can be computed in *parallel* because the entire sequence of tokens $\{t_k\}$ is already available as input. Figure 2-a shows a visual illustration.

Train-test shifts. Given the image-conditioned model of $P(\mathbf{t}|\mathbf{i})$ above, we now analyze its behavior when applied to test data distributions that differ from the trainset, denoted as P_{test} versus P_{train} . Recall that any joint distribution over images and text can be factored into a product over a language prior and an image likelihood $P(\mathbf{t}, \mathbf{i}) = P(\mathbf{t})P(\mathbf{i}|\mathbf{t})$. Our analysis makes the strong assumption that the image likelihood $P(\mathbf{i}|\mathbf{t})$ is identical across the train and test data, but the language prior $P(\mathbf{t})$ may differ. Intuitively, this assumes that the visual appearance of entities (such as a “white duck”) remains consistent across the training and test data, but the frequency of those entities (as manifested in the set of captions $P(\mathbf{t})$) may vary. We can now

derive $P_{test}(\mathbf{t}|\mathbf{i})$ via Bayes rule:

$$P_{test}(\mathbf{t}|\mathbf{i}) \propto P(\mathbf{i}|\mathbf{t})P_{test}(\mathbf{t}) \quad (2)$$

$$= P(\mathbf{i}|\mathbf{t}) \frac{P_{train}(\mathbf{t})}{P_{train}(\mathbf{t})} P_{test}(\mathbf{t}) \quad (3)$$

$$\propto P_{train}(\mathbf{t}|\mathbf{i}) \frac{P_{test}(\mathbf{t})}{P_{train}(\mathbf{t})} \quad (4)$$

The above shows that the generative pre-training score $P_{train}(\mathbf{t}|\mathbf{i})$ need simply be weighted by the *ratio* of the language priors in the testset versus trainset. Intuitively, if a particular text caption appears *more* often in the testset than the trainset, one should *increase* the score reported by the generative model. However, one often does not have access to the text distribution on the testset. For example, real-world deployments and benchmark protocols may not reveal this. In such cases, one can make two practical assumptions; either the language distribution on test is identical to train, or it is uninformative/uniform (see Figure 1):

Scenario 1:

$$P_{test}(\mathbf{t}) = P_{train}(\mathbf{t}) \Rightarrow \text{Optimal score is } P_{train}(\mathbf{t}|\mathbf{i}) \quad (5)$$

Scenario 2:

$$P_{test}(\mathbf{t}) \text{ is uniform.} \Rightarrow \text{Optimal score is } \frac{P_{train}(\mathbf{t}|\mathbf{i})}{P_{train}(\mathbf{t})} \quad (6)$$

Tunable α . In reality, a testset might be a mix of both scenarios. To model this, we consider a soft combination where the language prior on the testset is assumed to be a flattened version of the language prior on the trainset, for

some temperature parameter $\alpha \in [0, 1]$:

$$P_{test}(\mathbf{t}) \propto P_{train}(\mathbf{t})^{1-\alpha} \Rightarrow \text{Optimal score is } \frac{P_{train}(\mathbf{t}|\mathbf{i})}{P_{train}(\mathbf{t})^\alpha} \quad (7)$$

By setting α to 0 or 1, one can obtain the two scenarios described above. Some deployments (or benchmarks) may benefit from tuning α on a held-out valset, if available.

Implications for retrieval benchmarks. We speculate some benchmarks like ARO-Flickr (Yuksekonul et al., 2022) are close to Scenario 1 because they include negative captions that are *implausible*, such as “a white duck the its wings while in water spreads”. Such captions will have a low score under the language prior $P_{train}(\mathbf{t})$ and so reporting the raw generative score $P_{train}(\mathbf{t}|\mathbf{i})$ (that keeps its language prior or bias) will improve accuracy. In fact, we show that applying a *blind* language model (that ignores all image evidence) can itself often identify the correct caption. On the other hand, for test datasets with more *realistic* negative captions (Scenario 2), it may be useful to remove the language bias of the trainset, since that will prefer to match to common captions (even if they do not necessarily agree with the input image). This appears to be the case for Sugar-Crepe (Hsieh et al., 2023), which uses LLMs like ChatGPT to ensure that the negative captions are realistic.

An information-theoretic derivation of α -debiasing. Our approach to debiasing is reminiscent of mutual information, which can also be seen as a method for removing the effect of marginal priors when computing joint probability scores (Daille, 1994). In fact, α -debiasing (Eq. 7) is equivalent to a form of pointwise mutual information (PMI) known as PMI^k (Role & Nadif, 2011). PMI is a classic information-theoretic measure that quantifies the association between two variables (Yao et al., 2010; Henning & Ewerth, 2017; Shrivastava et al., 2021). In the context of image-text retrieval, PMI measures how much more or less likely the image-text pair co-occurs than if the two were independent:

$$\text{pmi}_P(\mathbf{t}, \mathbf{i}) = \frac{P(\mathbf{t}, \mathbf{i})}{P(\mathbf{t})P(\mathbf{i})} = \frac{P(\mathbf{i}|\mathbf{t})}{P(\mathbf{i})} = \frac{P(\mathbf{t}|\mathbf{i})}{P(\mathbf{t})} \quad (8)$$

However, directly applying PMI (Eq. 8) for retrieval tends to overly inflate scores for rarer texts (Role & Nadif, 2011). Consequently, the PMI^k approach was introduced to control the strength of debiasing. Below, we rewrite the Eq. 7 using

the language of PMI^k :

$$\frac{P_{train}(\mathbf{t}|\mathbf{i})}{P_{train}(\mathbf{t})^\alpha} = \frac{P_{train}(\mathbf{t}, \mathbf{i})}{P_{train}(\mathbf{i})P_{train}(\mathbf{t})^\alpha} \quad (9)$$

$$\propto \frac{P_{train}(\mathbf{t}, \mathbf{i})^{\frac{1}{\alpha}}}{P_{train}(\mathbf{i})P_{train}(\mathbf{t})} \quad (10)$$

, as $P_{train}(\mathbf{i})$ is constant in I-to-T

$$= \text{pmi}_{P_{train}}^k(\mathbf{t}, \mathbf{i}) \quad , \text{ where } k = \frac{1}{\alpha} \geq 1 \quad (11)$$

Eq. 11 shows that our α -debiasing is equivalent to PMI^k for $k = \frac{1}{\alpha}$. PMI^k is widely adopted in information retrieval tasks (Li et al., 2016; Li & Jurafsky, 2016; Wang et al., 2020). This alternative derivation could explain why α -debiasing remains effective across various testing benchmarks (as we show next), even when our previous probabilistic assumptions may not hold.

3. Experimental results on I-to-T retrieval

In this section, we verify our hypothesis on I-to-T retrieval benchmarks using state-of-the-art multimodal generative VLMs. In particular, we adopt image-conditioned language models such as BLIP (Li et al., 2022) as the learned estimator of $P_{train}(\mathbf{t}|\mathbf{i})$. Then, we discuss how we perform Monte Carlo estimation of $P_{train}(\mathbf{t})$, including a novel efficient sampling method based on “content-free” Gaussian noise images. Finally, we show the state-of-the-art results of our generative approach on recent I-to-T retrieval benchmarks.

Preliminaries. We leverage OTS image-conditioned language models to estimate $P_{train}(\mathbf{t})$. Most of our diagnostic experiments focus on the open-sourced BLIP (Li et al., 2022; 2023) model, trained on public image-text corpora using discriminative (ITC and ITM) and generative (captioning) objectives. Discriminative objectives typically model $P(\text{match}|\mathbf{t}, \mathbf{i})$. For example, ITCScore calculates cosine similarity scores between image and text features using a dual-encoder; ITMScore jointly embeds image-text pairs via a fusion-encoder and returns softmax scores from a binary classifier. We term the generative score as **Visual Generative Pre-Training Score (VisualGPTScore)**. While BLIP is pre-trained using all three objectives, this generative score has not been applied to discriminative tasks before our work. Lastly, our approach can be extended to other generative VLMs. We also present some additional results using LLaVA-1.5 (Liu et al., 2023), a recent state-of-the-art VLM (Liu et al., 2023) that produces SOTA accuracy on several challenging benchmarks.

Implementing VisualGPTScore. Our method calculates an average of the log-likelihoods of t_k at each token position k and applies an exponent to cancel the log:

$$\text{VisualGPTScore}(\mathbf{t}, \mathbf{i}) := e^{\frac{1}{m} \sum_{k=1}^m \log(P(t_k|t_{<k}, \mathbf{i}))} \quad (12)$$

To condition on an input image, BLIP uses a multimodal casual self-attention mask (Li et al., 2022) in its image-grounded text decoder, i.e., each text token attends to all its preceding vision and text tokens. We emphasize that VisualGPTScore has the same computational cost as ITM-Score, which uses the same underlying transformer but with a bi-directional self-attention mask to encode an image-text pair.

Estimating $P_{train}(\mathbf{t})$ using Monte Carlo sampling (oracle approach). Given $P_{train}(\mathbf{t}|\mathbf{i})$, we can estimate $P_{train}(\mathbf{t})$ via classic Monte Carlo sampling (Shapiro, 2003), by drawing n images from the train distribution, such as LAION114M (Schuhmann et al., 2021) for BLIP:

$$P_{train}(\mathbf{t}) \approx \frac{1}{n} \sum_{k=1}^n P_{train}(\mathbf{t}|\mathbf{i}_k) \quad (13)$$

Reducing sampling cost with Gaussian noise images (our approach). The above Equation 13 requires many trainset samples to achieve robust estimates. To address this, we draw inspiration from (Zhao et al., 2021), which uses a *content-free* text prompt “N/A” to calibrate the probability of a text from LLMs, i.e., $P(\mathbf{t}|\text{“N/A”})$. To apply this to our generative VLMs, we choose to sample “null” inputs as Gaussian noise images. It turns out Eq. 13 can be estimated using as few as 1-3 Gaussian noise images (with a mean and standard deviation calculated from trainset distribution). We provide a visual illustration of this method in Figure 2-b. We find this method to be less computationally demanding and just as effective as sampling thousands of images from trainset.

Benchmarks and evaluation protocols. We comprehensively report on four recent I-to-T retrieval benchmarks that assess compositionality, including ARO (Yuksekgonul et al., 2022), Crepe (Ma et al., 2022), SugarCrepe (Hsieh et al., 2023), and VL-CheckList (Zhao et al., 2022). In these datasets, each image has a single positive caption and multiple negative captions. ARO (Yuksekgonul et al., 2022) has four datasets: VG-Relation, VG-Attribution, COCO-Order, and Flickr30k-Order. SugarCrepe (Hsieh et al., 2023) has three datasets: Replace, Swap, and Add. For Crepe (Ma et al., 2022), we use the entire productivity set and report on three datasets: Atom, Negate, and Swap. VL-CheckList (Zhao et al., 2022) has three datasets: Object, Attribute, and Relation.

SOTA performance on all four benchmarks. In Table 1, we show that our OTS generative approaches, based on the BLIP model pre-trained on LAION-114M with ViT-L image encoder, achieves state-of-the-art results on all benchmarks. We outperform the best discriminative VLMs, including LAION5B-CLIP, and consistently surpass other heavily-engineered solutions, including NegCLIP, SyViC, MosaiCLIP, DAC, SVLC, SGVL, Structure-CLIP, all of

which fine-tune CLIP on much more data. For reference, we also include results of text-only Vera and Grammar from Hsieh et al. (2023). To show that even the most recent SugarCrepe is not exempt from language biases, we run two more text-only methods:

1. $P_{LLM}(\mathbf{t})$: passing captions into a pure LLM, such as BART-base (Yuan et al., 2021), FLAN-T5-XL (Chung et al., 2022), and OPT-2.7B (Zhang et al., 2022), to compute a text-only GPTScore (Fu et al., 2023).
2. $P_{train}(\mathbf{t})$: passing both captions and Gaussian noise images to BLIP as shown in Figure 2.

Discussion on α -debiasing. Table 2 shows that debiasing affects benchmarks differently depending on their construction; benchmarks with unrealistic negative captions (such as ARO-Flicker) benefit from a language prior that can identify such negative examples. Here, debiasing with large α hurts performance. On the other hand, benchmarks with realistic negative captions (such as SugarCrepe) tend to benefit from debiasing because it reduces the influence of the language prior. Our findings are reminiscent of the lessons from the VQA benchmark (Goyal et al., 2017), known to be solvable by “blind” algorithms that do not look at the image, e.g., questions such as “Is there a clock” have an answer of “Yes” 98% of the time. However, we also find that some recent benchmarks such as Winoground (Thrush et al., 2022) and EqBen (Wang et al., 2023) introduce strict evaluation protocols that aggressively penalize such blind algorithms. We discuss these challenging Scenario 2 benchmarks (with far lower SOTA accuracy) in the next section.

4. Additional Challenging Benchmarks

In this section, we apply our OTS generative approaches to five more Scenario 2 benchmarks: (a) Winoground (Thrush et al., 2022) and EqBen (Wang et al., 2023) for image-text alignment; (b) COCO (Lin et al., 2014) and Flickr30K (Young et al., 2014) for large-scale retrieval; (c) ImageNet (Deng et al., 2009) for zero-shot image classification. While naively applying OTS VisualGPTScore leads to inferior performance on these benchmarks, our training-free α -debiasing consistently improves its performance even with a fixed $\alpha=1$, without accessing the held-out valset (Table 3-a). We also derive the optimal text-to-image (T-to-I) retrieval objective and show that OTS generative scores can achieve robust T-to-I performance (Table 3-b). Lastly, we apply VisualGPTScore and its α -debiased version to a state-of-the-art VLM, LLaVA-1.5 (Liu et al., 2023), and outperform widely-used methods such as CLIPScore (Hessel et al., 2021) on the challenging Winoground and EqBen benchmarks. This suggests that VisualGPTScore is a superior choice for measuring image-text alignment.

Table 1. **OTS generative VLMs are SOTA on image-to-text retrieval benchmarks.** We begin by evaluating blind language models (in red). Surprisingly, this already produces SOTA accuracy on certain benchmarks such as ARO-Flickr, compared to the best discriminative approaches (in gray). We also find that blind inference of generative VLMs, $P_{train}(t)$ via sampling Gaussian noise images (in blue), often performs better and achieve above-chance performance even on the most recent SugarCrepe. Next, we show that simply repurposing a generative VLM’s language generation head for computing image-text scores (VisualGPTScore in yellow), which corresponds to $\alpha = 0$, consistently produces SOTA accuracy across all benchmarks. Finally, debiasing this score by tuning α on valset (in green) further improves performance, establishing the new SOTA.

Score	Method	ARO			
		Rel	Attr	COCO	Flickr
Random	-	50.0	50.0	20.0	20.0
Text-Only	Vera	61.7	82.6	59.8	63.5
	Grammar	59.6	58.4	74.3	76.3
$P_{LLM}(t)$	BART	81.1	73.6	95.0	95.2
	Flan-T5	84.4	76.5	98.0	98.2
	OPT	84.7	79.8	97.9	98.6
$P_{train}(t)$	BLIP	87.6	80.7	98.6	99.1
$P(match t, i)$	CLIP	59.0	62.0	46.0	60.0
	LAION2B-CLIP	51.6	61.9	25.2	30.2
	LAION5B-CLIP	46.1	57.8	26.1	31.0
	NegCLIP	81.0	71.0	86.0	91.0
	Structure-CLIP	83.5	85.1	-	-
	SyViC	80.8	72.4	92.4	87.2
	SGVL	-	-	87.2	91.0
	MosaiCLIP	82.6	78.0	87.9	86.3
	DAC-LLM	81.3	73.9	94.5	95.7
	DAC-SAM	77.2	70.5	91.2	93.9
	BLIP-ITC	63.1	81.6	34.3	41.7
	BLIP-ITM	58.7	90.3	45.1	51.3
$P_{train}(t i)$	Ours ($\alpha = 0$)	89.1	95.3	99.4	99.5
$P_{train}(t)^\alpha$	Ours ($\alpha = 1$)	68.1	87.9	32.4	44.5
	Ours ($\alpha = \alpha^*$)	89.1	95.4	99.4	99.5

(a) Accuracy on ARO

Score	Method	SugarCrepe		
		Replace	Swap	Add
Random	-	50.0	50.0	50.0
Text-Only	Vera	49.5	49.3	49.5
	Grammar	50.0	50.0	50.0
$P_{LLM}(t)$	BART	48.4	51.9	61.2
	Flan-T5	51.4	57.6	40.9
	OPT	58.5	66.6	45.8
$P_{train}(t)$	BLIP	75.9	77.1	70.9
$P(match t, i)$	CLIP	80.8	63.3	75.1
	LAION2B-CLIP	86.5	68.6	88.4
	LAION5B-CLIP	85.0	68.0	89.6
	NegCLIP	88.3	76.2	90.2
	BLIP-ITC	85.8	73.8	85.7
	BLIP-ITM	88.7	81.3	87.6
$P_{train}(t i)$	Ours ($\alpha = 0$)	93.3	91.0	91.0
$P_{train}(t)^\alpha$	Ours ($\alpha = 1$)	83.2	85.5	85.9
	Ours ($\alpha = \alpha^*$)	95.1	92.4	97.4

(c) Accuracy on SugarCrepe

Score	Method	VL-CheckList		
		Object	Attribute	Relation
Random	-	50.0	50.0	50.0
Text-Only	Vera	82.5	74.0	85.7
	Grammar	58.0	52.4	68.5
$P_{LLM}(t)$	BART	52.0	51.0	45.1
	Flan-T5	60.3	55.0	49.3
	OPT	59.3	48.8	60.0
$P_{train}(t)$	BLIP	68.2	58.7	75.9
$P(match t, i)$	CLIP	81.6	67.6	63.1
	LAION2B-CLIP	84.7	67.8	66.5
	LAION5B-CLIP	87.9	70.3	63.9
	NegCLIP	81.4	72.2	63.5
	SyViC	-	70.4	69.4
	SGVL	85.2	78.2	80.4
	SLVC	85.0	72.0	69.0
	DAC-LLM	87.3	77.3	86.4
	DAC-SAM	88.5	75.8	89.8
	BLIP-ITC	90.6	80.3	73.5
	BLIP-ITM	89.9	80.7	67.7
$P_{train}(t i)$	Ours ($\alpha = 0$)	92.6	78.7	90.8
$P_{train}(t)^\alpha$	Ours ($\alpha = 1$)	90.4	77.6	77.8
	Ours ($\alpha = \alpha^*$)	94.4	82.1	92.8

(b) Accuracy on VL-CheckList

Score	Method	Crepe		
		Atom	Swap	Negate
Random	-	16.7	16.7	16.7
Text-Only	Vera	43.7	70.8	66.2
	Grammar	18.2	50.9	9.8
$P_{LLM}(t)$	BART	38.8	53.3	44.4
	Flan-T5	43.0	69.5	13.6
	OPT	53.3	72.7	5.0
$P_{train}(t)$	BLIP	55.4	69.7	60.8
$P(match t, i)$	CLIP	22.3	26.6	28.8
	LAION2B-CLIP	23.6	24.8	18.0
	LAION5B-CLIP	24.2	23.9	20.1
	BLIP-ITC	24.8	17.7	26.5
	BLIP-ITM	29.5	20.7	25.5
$P_{train}(t i)$	Ours ($\alpha = 0$)	73.2	78.1	79.6
$P_{train}(t)^\alpha$	Ours ($\alpha = 1$)	20.6	28.3	35.6
	Ours ($\alpha = \alpha^*$)	73.3	78.1	79.6

(d) Accuracy on Crepe

Balanced evaluation protocols for retrieval. Winoground and EqBen evaluate image-text alignment through retrieval tasks, and we find their evaluation protocols discourage blind solutions. We refer the reader to the benchmarks for more details, but in summary, both benchmarks operate on

pairs of image-text pairs $\{(i_0, t_0), (i_1, t_1)\}$ and construct two I-to-T retrieval (text score) tasks with a single image and two candidate captions. The text score is awarded 1 point only if *both* retrieval tasks are correct. Consider the common case where one caption is more likely under a

Table 2. α -debiasing on I-to-T benchmarks and $P_{train}(\mathbf{t})$ frequency charts of both positive and negative captions. Increasing α from 0 to 1 hurts performance on benchmarks with non-sensical negative captions like ARO-Flickr. ARO’s negative captions are easier to identify because of their low score under the language prior $P_{train}(\mathbf{t})$, implying such benchmarks may even be solved with blind algorithms that avoid looking at images. On the other hand, for benchmarks like SugarCrepe with more balanced $P_{train}(\mathbf{t})$ between positive and negative captions, tuning α leads to performance gain.

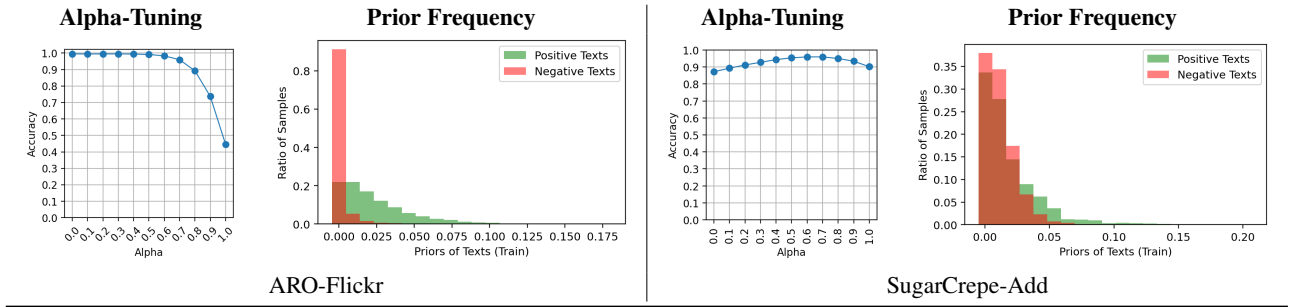


Table 3. Additional results on Winoground/EqBen/COCO/Flickr30K/ImageNet1K. Table (a) shows the importance of α -debiasing on these compositionality and large-scale retrieval benchmarks. While OTS generative scores do not work well, debiasing with a larger α close to 1 can consistently and often significantly improve I-to-T performance. To highlight the improvement, we mark results without debiasing ($\alpha = 0$) (in yellow), debiasing with a fixed $\alpha = 1$ (in pink), and cross-validation using held-out valsets ($\alpha = \alpha_{val}^*$) (in green). Table (b) shows that OTS generative scores can obtain favorable results on all T-to-I retrieval tasks, competitive with the ITMScore.

Metric	Benchmark	ITMScore	$\frac{P_{train}(\mathbf{t} \mathbf{i})}{P_{train}(\mathbf{t})^\alpha}$			
			$\alpha=0$	$\alpha=1$	$\alpha=\alpha_{val}^*$	α_{val}^*
Text Score	Winoground	35.5 _(2.4)	27.5 _(2.3)	33.7 _(2.4)	36.6 _(2.6)	0.855 _(0.023)
	EqBen	26.1 _(0.3)	9.6 _(0.2)	19.8 _(0.3)	19.8 _(0.3)	0.992 _(0.007)
R@1 / R@5	COCO	71.9 / 90.6	19.7 / 40.6	46.2 / 73.1	48.0 / 74.2	0.819
	Flickr30k	88.8 / 98.2	34.6 / 59.0	58.7 / 88.0	63.6 / 89.2	0.719
Accuracy	ImageNet1K	37.4	18.6	36.2	40.0	0.670

(a) α -debiasing on valsets for I-to-T retrieval

Metric	Benchmark	ITMScore	$P_{train}(\mathbf{t} \mathbf{i})$
Image Score	Winoground	15.8	21.5
	EqBen	20.3	26.1
R@1 / R@5	COCO	54.8 / 79.0	55.6 / 79.2
	Flickr30k	77.8 / 93.9	76.8 / 93.4

(b) T-to-I retrieval

language prior; here the common caption will be correctly retrieved for one of the tasks but will be incorrectly retrieved for the other, implying *no* points will be awarded. Similarly stringent metrics are used for T-to-I retrieval (image score). The final group score is awarded 1 point only if all 4 retrieval tasks are correct.

α -debiasing consistently improves I-to-T retrieval. Table 3-a shows that simply debiasing VisualGPTScore with a fixed $\alpha = 1$ significantly improves performance on challenging I-to-T benchmarks. One can also do slightly better by using a held-out valset to tune for the optimal $\alpha \in [0, 1]$. For Winoground and EqBen, we sample half of the data as a valset and perform a grid search for α_{val}^* (using a step size of 0.001), reporting the performance on the other half. We repeat this process 10 times and report the mean and standard deviation. For COCO and Flickr30K, we perform α -debiasing using Recall@1 (R@1) on the official valset. We report the zero-shot classification accuracy on ImageNet1K, which can be viewed as an I-to-T retrieval task that retrieves the best textual label (out of 1000) for each image. We simply use one-shot samples from Lin et al. (2023) to cross

validate on ImageNet, which incurs negligible costs.

Lastly, we observe that generative approaches still lag behind the ITMScore of BLIP for the two large-scale retrieval benchmarks.

Extending to T-to-I retrieval. Though not the focus of our work, we show that image-conditioned language models can be applied to T-to-I retrieval. Given a text caption \mathbf{t} , we can rewrite the Bayes optimal T-to-I retrieval objective as:

$$P_{test}(\mathbf{i}|\mathbf{t}) \propto P_{train}(\mathbf{t}|\mathbf{i}) * P_{train}(\mathbf{i}) \quad (14)$$

Equation 14 is hard to implement because we do not have access to $P_{train}(\mathbf{i})$. However, when $P_{train}(\mathbf{i})$ is approximately uniform, one can directly apply $P_{train}(\mathbf{t}|\mathbf{i})$ for optimal performance. We report T-to-I performance in Table 3-b, where our generative approach obtains competitive results compared against ITMScore, likely because T-to-I retrieval is less affected by language biases.

Table 4. Superior performance of VisualGPTScore on challenging image-text alignment benchmarks. We compare VisualGPTScore (and its $\alpha=1$ version) against popular image-text scoring methods such as CLIPScore and those that combine VLMs with additional LLMs like ChatGPT. On Winoground and EqBen, our VisualGPTScore ($\alpha=0$) outperforms all methods using only a state-of-the-art VLM (LLaVA-1.5). Moreover, debiasing with $\alpha=1$ (using a single Gaussian noise image) consistently improves I-to-T retrieval, thereby increasing the text and group score. To ensure a fair comparison, we use the publicly available model checkpoints and corresponding code of prior works.

Method	LLMs used	Winoground			EqBen		
		Text	Image	Group	Text	Image	Group
Random Chance	–	25.0	25.0	16.7	25.0	25.0	16.7
<i>Official implementation</i>							
CLIPScore	–	31.3	11.0	8.8	35.0	33.6	21.4
VPEval	ChatGPT	12.8	11.0	6.3	34.3	25.7	21.4
LLMScore	ChatGPT	21.3	17.8	12.5	32.9	27.9	22.9
<i>Our results based on LLaVA-1.5</i>							
TIFA	Llama-2	22.8	18.5	15.5	30.0	30.0	21.4
VQ2	FlanT5	14.0	27.3	10.0	22.9	40.7	20.0
Davidsonian	ChatGPT	21.0	16.8	15.5	26.4	20.0	20.0
VisualGPTScore ($\alpha=0$)	–	36.3	37.0	24.8	25.7	42.1	21.4
VisualGPTScore ($\alpha=1$)	–	44.3	37.0	27.5	42.9	42.1	29.3

References

- Agrawal, H., Desai, K., Wang, Y., Chen, X., Jain, R., Johnson, M., Batra, D., Parikh, D., Lee, S., and Anderson, P. Nocaps: Novel object captioning at scale. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 8948–8957, 2019.
- Alayrac, J.-B., Donahue, J., Luc, P., Miech, A., Barr, I., Hasson, Y., Lenc, K., Mensch, A., Millican, K., Reynolds, M., et al. Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems*, 35:23716–23736, 2022.
- Antol, S., Agrawal, A., Lu, J., Mitchell, M., Batra, D., Zitnick, C. L., and Parikh, D. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pp. 2425–2433, 2015.
- Bengio, Y., Ducharme, R., Vincent, P., and Jauvin, C. A neural probabilistic language model. *Journal of Machine Learning Research*, 3:1137–1155, 2003.
- Bertolini, L., Weeds, J., and Weir, D. Testing large language models on compositionality and inference with phrase-level adjective-noun entailment. In *Proceedings of the 29th International Conference on Computational Linguistics*, pp. 4084–4100, 2022.
- Betker, J., Goh, G., Jing, L., Brooks, T., Wang, J., Li, L., Ouyang, L., Zhuang, J., Lee, J., Guo, Y., et al. Improving image generation with better captions. <https://cdn.openai.com/papers/dall-e-3.pdf>, 2023.
- Brendel, W. and Bethge, M. Approximating cnns with bag-of-local-features models works surprisingly well on imagenet. *arXiv preprint arXiv:1904.00760*, 2019.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33: 1877–1901, 2020.
- Cascante-Bonilla, P., Shehada, K., Smith, J. S., Doveh, S., Kim, D., Panda, R., Varol, G., Oliva, A., Ordonez, V., Feris, R., et al. Going beyond nouns with vision & language models using synthetic data. *arXiv preprint arXiv:2303.17590*, 2023.
- Cho, J., Hu, Y., Garg, R., Anderson, P., Krishna, R., Baldridge, J., Bansal, M., Pont-Tuset, J., and Wang, S. Davidsonian scene graph: Improving reliability in fine-grained evaluation for text-image generation. *arXiv preprint arXiv:2310.18235*, 2023a.
- Cho, J., Zala, A., and Bansal, M. Visual programming for text-to-image generation and evaluation. *arXiv preprint arXiv:2305.15328*, 2023b.
- Chung, H. W., Hou, L., Longpre, S., Zoph, B., Tay, Y., Fedus, W., Li, E., Wang, X., Dehghani, M., Brahma, S., Webson, A., Gu, S. S., Dai, Z., Suzgun, M., Chen, X., Chowdhery, A., Valter, D., Narang, S., Mishra, G., Yu, A. W., Zhao, V., Huang, Y., Dai, A. M., Yu, H., Petrov, S., hsin Chi, E. H., Dean, J., Devlin, J., Roberts, A., Zhou, D., Le, Q. V., and Wei, J. Scaling instruction-finetuned language models. *ArXiv*, abs/2210.11416, 2022.
- Daille, B. *Approche mixte pour l’extraction automatique de terminologie: statistiques lexicales et filtres linguistiques*. PhD thesis, Ph. D. thesis, Université Paris 7, 1994.
- Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- Diwan, A., Berry, L., Choi, E., Harwath, D., and Mahowald, K. Why is winoground hard? investigating failures in visuolinguistic compositionality. *arXiv preprint arXiv:2211.00768*, 2022.
- Doveh, S., Arbelle, A., Harary, S., Panda, R., Herzig, R., Schwartz, E., Kim, D., Giryes, R., Feris, R., Ullman, S., et al. Teaching structured vision&language concepts to vision&language models. *arXiv preprint arXiv:2211.11733*, 2022.
- Doveh, S., Arbelle, A., Harary, S., Alfassy, A., Herzig, R., Kim, D., Giryes, R., Feris, R., Panda, R., Ullman, S., et al. Dense and aligned captions (dac) promote compositional reasoning in vl models. *arXiv preprint arXiv:2305.19595*, 2023.
- Fang, Y., Wang, W., Xie, B., Sun, Q., Wu, L., Wang, X., Huang, T., Wang, X., and Cao, Y. Eva: Exploring the limits of masked visual representation learning at scale. *arXiv preprint arXiv:2211.07636*, 2022.
- Fu, J., Ng, S.-K., Jiang, Z., and Liu, P. Gptscore: Evaluate as you desire. *arXiv preprint arXiv:2302.04166*, 2023.
- Goyal, Y., Khot, T., Summers-Stay, D., Batra, D., and Parikh, D. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 6904–6913, 2017.
- Guo, C., Zhao, B., and Bai, Y. Deepcore: A comprehensive library for coreset selection in deep learning. In *Database and Expert Systems Applications: 33rd International Conference, DEXA 2022, Vienna, Austria, August*

- 22–24, 2022, *Proceedings, Part I*, pp. 181–195. Springer, 2022.
- Henning, C. A. and Ewerth, R. Estimating the information gap between textual and visual representations. In *Proceedings of the 2017 ACM on International Conference on Multimedia Retrieval*, pp. 14–22, 2017.
- Herzig, R., Mendelson, A., Karlinsky, L., Arbelle, A., Feris, R., Darrell, T., and Globerson, A. Incorporating structured representations into pretrained vision & language models using scene graphs. *arXiv preprint arXiv:2305.06343*, 2023.
- Hessel, J. and Schofield, A. How effective is bert without word ordering? implications for language understanding and data privacy. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pp. 204–211, 2021.
- Hessel, J., Holtzman, A., Forbes, M., Bras, R. L., and Choi, Y. Clipscore: A reference-free evaluation metric for image captioning. *arXiv preprint arXiv:2104.08718*, 2021.
- Hsieh, C.-Y., Zhang, J., Ma, Z., Kembhavi, A., and Krishna, R. Sugarcrepe: Fixing hackable benchmarks for vision-language compositionality. *arXiv preprint arXiv:2306.14610*, 2023.
- Hu, Y., Liu, B., Kasai, J., Wang, Y., Ostendorf, M., Krishna, R., and Smith, N. A. Tifa: Accurate and interpretable text-to-image faithfulness evaluation with question answering. *arXiv preprint arXiv:2303.11897*, 2023.
- Huang, Y., Tang, J., Chen, Z., Zhang, R., Zhang, X., Chen, W., Zhao, Z., Lv, T., Hu, Z., and Zhang, W. Structureclip: Enhance multi-modal language representations with structure knowledge. *arXiv preprint arXiv:2305.06152*, 2023.
- Kamath, A., Hessel, J., and Chang, K.-W. Text encoders are performance bottlenecks in contrastive vision-language models. *arXiv preprint arXiv:2305.14897*, 2023.
- Kang, B., Xie, S., Rohrbach, M., Yan, Z., Gordo, A., Feng, J., and Kalantidis, Y. Decoupling representation and classifier for long-tailed recognition. *arXiv preprint arXiv:1910.09217*, 2019.
- Keskar, N. S., McCann, B., Varshney, L. R., Xiong, C., and Socher, R. Ctrl: A conditional transformer language model for controllable generation. *arXiv preprint arXiv:1909.05858*, 2019.
- Li, B., Lin, Z., Pathak, D., Li, J., Fei, Y., Wu, K., Xia, X., Zhang, P., Neubig, G., and Ramanan, D. Evaluating and improving compositional text-to-visual generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2024.
- Li, J. and Jurafsky, D. Mutual information and diverse decoding improve neural machine translation, 2016.
- Li, J., Galley, M., Brockett, C., Gao, J., and Dolan, B. A diversity-promoting objective function for neural conversation models. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 110–119, San Diego, California, June 2016. Association for Computational Linguistics. doi: 10.18653/v1/N16-1014. URL <https://aclanthology.org/N16-1014>.
- Li, J., Selvaraju, R., Gotmare, A., Joty, S., Xiong, C., and Hoi, S. C. H. Align before fuse: Vision and language representation learning with momentum distillation. *Advances in neural information processing systems*, 34: 9694–9705, 2021.
- Li, J., Li, D., Xiong, C., and Hoi, S. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning*, pp. 12888–12900. PMLR, 2022.
- Li, J., Li, D., Savarese, S., and Hoi, S. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*, 2023.
- Lin, T.-Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., and Zitnick, C. L. Microsoft coco: Common objects in context. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pp. 740–755. Springer, 2014.
- Lin, Z., Yu, S., Kuang, Z., Pathak, D., and Ramana, D. Multimodality helps unimodality: Cross-modal few-shot learning with multimodal models. *arXiv preprint arXiv:2301.06267*, 2023.
- Lin, Z., Pathak, D., Li, B., Li, J., Xia, X., Neubig, G., Zhang, P., and Ramanan, D. Evaluating text-to-visual generation with image-to-text generation. *arXiv preprint arXiv:2404.01291*, 2024.
- Liu, H., Li, C., Wu, Q., and Lee, Y. J. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*, 2023.
- Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., and Stoyanov, V. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.

- Lu, Y., Yang, X., Li, X., Wang, X. E., and Wang, W. Y. Llm-score: Unveiling the power of large language models in text-to-image synthesis evaluation. *arXiv preprint arXiv:2305.11116*, 2023.
- Ma, Z., Hong, J., Gul, M. O., Gandhi, M., Gao, I., and Krishna, R. Crepe: Can vision-language foundation models reason compositionally? *arXiv preprint arXiv:2212.07796*, 2022.
- Mehrabani, N., Morstatter, F., Saxena, N., Lerman, K., and Galstyan, A. A survey on bias and fairness in machine learning. *ACM Computing Surveys (CSUR)*, 54(6):1–35, 2021.
- Miech, A., Alayrac, J.-B., Laptev, I., Sivic, J., and Zisserman, A. Thinking fast and slow: Efficient text-to-visual retrieval with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9826–9836, 2021.
- OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Papadimitriou, I., Futrell, R., and Mahowald, K. When classifying grammatical role, bert doesn’t care about word order... except when it matters. *arXiv preprint arXiv:2203.06204*, 2022.
- Parashar, S., Lin, Z., Liu, T., Dong, X., Li, Y., Ramanan, D., Caverlee, J., and Kong, S. The neglected tails of vision-language models. *arXiv preprint arXiv:2401.12425*, 2024.
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I., et al. Language models are unsupervised multitask learners. 2019.
- Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.
- Role, F. and Nadif, M. Handling the impact of low frequency events on co-occurrence based measures of word similarity. In *Proceedings of the international conference on Knowledge Discovery and Information Retrieval (KDIR-2011)*. Scitepress, pp. 218–223, 2011.
- Schuhmann, C., Vencu, R., Beaumont, R., Kaczmarczyk, R., Mullis, C., Katta, A., Coombes, T., Jitsev, J., and Komatsuzaki, A. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114*, 2021.
- Schuhmann, C., Beaumont, R., Vencu, R., Gordon, C., Wightman, R., Cherti, M., Coombes, T., Katta, A., Mullis, C., Wortsman, M., et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *arXiv preprint arXiv:2210.08402*, 2022.
- Shapiro, A. Monte carlo sampling methods. *Handbooks in operations research and management science*, 10:353–425, 2003.
- Shrivastava, A., Selvaraju, R. R., Naik, N., and Ordonez, V. Clip-lite: information efficient visual representation learning from textual annotations. *arXiv preprint arXiv:2112.07133*, 2021.
- Singh, H., Zhang, P., Wang, Q., Wang, M., Xiong, W., Du, J., and Chen, Y. Coarse-to-fine contrastive learning in image-text-graph space for improved vision-language compositionality. *arXiv preprint arXiv:2305.13812*, 2023.
- Sinha, K., Jia, R., Hupkes, D., Pineau, J., Williams, A., and Kiela, D. Masked language modeling and the distributional hypothesis: Order word matters pre-training for little. *arXiv preprint arXiv:2104.06644*, 2021.
- Tejankar, A., Sanjabi, M., Wu, B., Xie, S., Khabsa, M., Pirsiavash, H., and Firooz, H. A fistful of words: Learning transferable visual models from bag-of-words supervision. *arXiv preprint arXiv:2112.13884*, 2021.
- Thrush, T., Jiang, R., Bartolo, M., Singh, A., Williams, A., Kiela, D., and Ross, C. Winoground: Probing vision and language models for visio-linguistic compositionality. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5238–5248, 2022.
- Tschannen, M., Kumar, M., Steiner, A., Zhai, X., Houlsby, N., and Beyer, L. Image captioners are scalable vision learners too. *arXiv preprint arXiv:2306.07915*, 2023.
- Wang, T., Lin, K., Li, L., Lin, C.-C., Yang, Z., Zhang, H., Liu, Z., and Wang, L. Equivariant similarity for vision-language foundation models. *arXiv preprint arXiv:2303.14465*, 2023.
- Wang, Z., Feng, B., Narasimhan, K., and Russakovsky, O. Towards unique and informative captioning of images. In *European Conference on Computer Vision (ECCV)*, 2020.
- Wortsman, M., Ilharco, G., Kim, J. W., Li, M., Kornblith, S., Roelofs, R., Lopes, R. G., Hajishirzi, H., Farhadi, A., Namkoong, H., et al. Robust fine-tuning of zero-shot models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7959–7971, 2022.
- Wu, X., Deng, Z., and Russakovsky, O. Multimodal dataset distillation for image-text retrieval. *arXiv preprint arXiv:2308.07545*, 2023.

- Yao, T., Mei, T., and Ngo, C.-W. Co-reranking by mutual reinforcement for image search. In *Proceedings of the ACM international conference on image and video retrieval*, pp. 34–41, 2010.
- Yarom, M., Bitton, Y., Changpinyo, S., Aharoni, R., Herzig, J., Lang, O., Ofek, E., and Szepkter, I. What you see is what you read? improving text-image alignment evaluation. *arXiv preprint arXiv:2305.10400*, 2023.
- Young, P., Lai, A., Hodosh, M., and Hockenmaier, J. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *Transactions of the Association for Computational Linguistics*, 2:67–78, 2014.
- Yuan, W., Neubig, G., and Liu, P. Bartscore: Evaluating generated text as text generation. In Ranzato, M., Beygelzimer, A., Dauphin, Y., Liang, P., and Vaughan, J. W. (eds.), *Advances in Neural Information Processing Systems*, volume 34, pp. 27263–27277. Curran Associates, Inc., 2021. URL <https://proceedings.neurips.cc/paper/2021/file/e4d2b6e6fdeca3e60e0f1a62fee3d9dd-Paper.pdf>.
- Yuksekgonul, M., Bianchi, F., Kalluri, P., Jurafsky, D., and Zou, J. When and why vision-language models behave like bag-of-words models, and what to do about it? *arXiv preprint arXiv:2210.01936*, 2022.
- Zhang, S., Roller, S., Goyal, N., Artetxe, M., Chen, M., Chen, S., Dewan, C., Diab, M., Li, X., Lin, X. V., et al. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*, 2022.
- Zhao, T., Wallace, E., Feng, S., Klein, D., and Singh, S. Calibrate before use: Improving few-shot performance of language models. In *International Conference on Machine Learning*, 2021.
- Zhao, T., Zhang, T., Zhu, M., Shen, H., Lee, K., Lu, X., and Yin, J. VI-checklist: Evaluating pre-trained vision-language models with objects, attributes and relations. *arXiv preprint arXiv:2207.00221*, 2022.