

An Introduction to Vision-Language Modeling

Florian Bordes^{*}, Richard Yuanzhe Pang^{*^}, Anurag Ajay^{**♣}, Alexander C. Li^{**♠}, Adrien Bardes^{*}, Suzanne Petryk[△], Oscar Mañas^{*‡}, Zhiqiu Lin[♠], Anas Mahmoud[†], Bargav Jayaraman^{*}, Mark Ibrahim^{*}, Melissa Hall^{*}, Yunyang Xiong^{*}, Jonathan Lebensold^{*♡}, Candace Ross^{*}, Srihari Jayakumar^{*}, Chuan Guo^{*}, Diane Bouchacourt^{*}, Haider Al-Tahan^{*}, Karthik Padthe^{*}, Vasu Sharma^{*}, Hu Xu^{*}, Xiaoqing Ellen Tan^{*}, Megan Richards^{*}, Samuel Lavoie^{*‡}, Pietro Astolfi^{*}, Reyhane Askari Hemmat^{*}, Jun Chen^{**◇}, Kushal Tirumala^{*}, Rim Assouel^{*‡}, Mazda Moayeri[▽], Arjang Talattof^{*}, Kamalika Chaudhuri^{*}, Zechun Liu^{*}, Xilun Chen^{*}, Quentin Garrido^{*}, Karen Ullrich^{*}, Aishwarya Agrawal^{‡•}, Kate Saenko^{*}, Asli Celikyilmaz^{*} and Vikas Chandra^{*}

^{*}Meta

^{**}Work done while at Meta

[‡]Université de Montréal, Mila

[♡]McGill University, Mila

[†]University of Toronto

[♠]Carnegie Mellon University

[♣]Massachusetts Institute of Technology

[^]New York University

[△]University of California, Berkeley

[▽]University of Maryland

[◇]King Abdullah University of Science and Technology

[•]Canada CIFAR AI Chair

Core contributors, random ordering

Additional contributors, random ordering

Senior contributors, random ordering

Contents

| | | |
|----------|--|-----------|
| 1 | Introduction | 4 |
| 2 | The Families of VLMs | 5 |
| 2.1 | Early work on VLMs based on transformers | 6 |
| 2.2 | Contrastive-based VLMs | 6 |
| 2.2.1 | CLIP | 9 |
| 2.3 | VLMs with masking objectives | 9 |
| 2.3.1 | FLAVA | 9 |
| 2.3.2 | MaskVLM | 10 |
| 2.3.3 | Information theoretic view on VLM objectives | 10 |
| 2.4 | Generative-based VLMs | 11 |
| 2.4.1 | An example of learning a text generator: CoCa | 11 |
| 2.4.2 | An example of multi-modal generative model: Chameleon and CM3leon | 11 |
| 2.4.3 | Using generative text-to-image models for downstream vision-language tasks | 12 |
| 2.5 | VLMs from Pretrained Backbones | 14 |
| 2.5.1 | Frozen | 15 |
| 2.5.2 | The example of MiniGPT | 15 |
| 2.5.3 | Other popular models using pretrained backbones | 16 |
| 3 | A Guide to VLM Training | 16 |
| 3.1 | Training data | 17 |
| 3.1.1 | Improving the training data with synthetic data | 19 |
| 3.1.2 | Using data augmentation | 19 |
| 3.1.3 | Interleaved data curation | 20 |
| 3.1.4 | Assessing multimodal data quality | 21 |
| 3.1.5 | Harnessing human expertise: the power of data annotation | 21 |
| 3.2 | Software | 21 |
| 3.2.1 | Using existing public software repositories | 22 |
| 3.2.2 | How many GPUs do I need? | 22 |
| 3.2.3 | Speeding up training | 22 |
| 3.2.4 | Importance of other hyper-parameters. | 23 |
| 3.3 | Which model to use? | 23 |
| 3.3.1 | When to use contrastive models like CLIP? | 23 |
| 3.3.2 | When to use masking? | 23 |
| 3.3.3 | When to use a generative model? | 24 |
| 3.3.4 | When to use LLM on pretrained backbone? | 24 |
| 3.4 | Improving grounding | 25 |
| 3.4.1 | Using bounding boxes annotations | 25 |
| 3.4.2 | Negative captioning | 25 |
| 3.5 | Improving alignment | 26 |
| 3.5.1 | A LLaVA story | 26 |
| 3.5.2 | Multimodal in-context learning | 27 |
| 3.6 | Improving text-rich image understanding | 28 |

| | | |
|----------|---|-----------|
| 3.7 | Parameter-Efficient Fine-Tuning | 29 |
| 4 | Approaches for Responsible VLM Evaluation | 31 |
| 4.1 | Benchmarking visio-linguistic abilities | 31 |
| 4.1.1 | Image captioning | 31 |
| 4.1.2 | Text-to-image consistency | 33 |
| 4.1.3 | Visual question answering | 33 |
| 4.1.4 | Text-centric Visual Question Answering | 34 |
| 4.1.5 | Zero-shot image classification | 35 |
| 4.1.6 | Visio-linguistic compositional reasoning | 36 |
| 4.1.7 | Dense captioning and crop-caption matching | 37 |
| 4.1.8 | Synthetic data based visio-linguistic evaluations | 38 |
| 4.2 | Benchmarking Bias and disparities in VLMs | 38 |
| 4.2.1 | Benchmarking bias via classifications | 38 |
| 4.2.2 | Benchmarking bias via embeddings | 39 |
| 4.2.3 | Language biases might impact your benchmark! | 40 |
| 4.2.4 | Evaluating how specific concepts in the training data impact down- stream performances | 40 |
| 4.3 | Benchmarking hallucinations | 40 |
| 4.4 | Benchmarking memorization | 41 |
| 4.5 | Red Teaming | 42 |
| 5 | Extending VLMs to Videos | 42 |
| 5.1 | Early work on Videos based on BERT | 43 |
| 5.2 | Enabling text generation using an early-fusion VLM | 44 |
| 5.3 | Using a pretrained LLM | 44 |
| 5.4 | Opportunities in evaluations | 45 |
| 5.5 | Challenges in leveraging video data | 45 |
| 6 | Conclusion | 46 |
| | Acronyms | 47 |

Abstract

Following the recent popularity of Large Language Models (LLMs), several attempts have been made to extend them to the visual domain. From having a visual assistant that could guide us through unfamiliar environments to generative models that produce images using only a high-level text description, the vision-language model (VLM) applications will significantly impact our relationship with technology. However, there are many challenges that need to be addressed to improve the reliability of those models. While language is discrete, vision evolves in a much higher dimensional space in which concepts cannot always be easily discretized. To better understand the mechanics behind mapping vision to language, we present this introduction to VLMs which we hope will help anyone who would like to enter the field. First, we introduce what VLMs are, how they work, and how to train them. Then, we present and discuss approaches to evaluate VLMs. Although this work primarily focuses on mapping images to language, we also discuss extending VLMs to videos.

1 Introduction

In recent years, we have seen impressive developments in language modeling. Many **Large Language Models (LLMs)** such as Llama or ChatGPT are now able to solve such a large variety of tasks that their usage is becoming more and more popular. Such models that were mostly limited to text inputs are now extended to having visual inputs. Connecting vision to language will unlock several applications that will be key to the current AI-based technological revolution. Even though several works have already extended large language models to vision, connecting language to vision is not completely solved. For example, most models struggle to understand spatial relationships or count without complicated engineering overhead that relies on additional data annotation. Many **Vision Language Models (VLMs)** also lack an understanding of attributes and ordering. They often ignore some part of the input prompt, leading to significant prompt engineering efforts to produce the desired result. Some of them can also hallucinate and produce content that is neither required nor relevant. As a consequence, developing reliable models is still a very active area of research.

In this work, we present an introduction to **Vision Language Models (VLMs)**. We explain what VLMs are, how they are trained, and how to effectively evaluate VLMs depending on different research goals. This work *should not be considered as a survey or a complete guide on VLMs*¹. Hence, we do not aim to cite every work from the VLM research field²; nor does this work capture every best practice in this space. Instead, we aim to provide a *clear and easy-to-understand introduction* to VLM research and highlight effective practices for research in this space. This introduction should be especially useful for students or researchers in other areas who want to enter the field.

We start by presenting the different VLM training paradigms. We discuss how contrastive methods have changed the field. Then, we present methods that leverage masking

¹For complete and more technical surveys on VLMs, please refer to Zhang et al. [2024a], Ghosh et al. [2024], Zhou and Shimada [2023], Chen et al. [2023a], Du et al. [2022], Uppal et al. [2022], and Liang et al. [2024].

²Nevertheless, if you find errors or have any comments about missing important references, please address them to fbordes@meta.com. We will try to include them in the next revision of this work.

strategies or generative components. Lastly, we present VLMs which use pre-trained backbones (such as LLMs). Categorizing VLMs into different families is not an easy task, since most of them have overlapping components. However, we hope that our categorization will help new researchers navigate the field and shed light on the inner mechanisms behind VLMs.

Next, we present typical recipes for training VLMs. For example, we cover: Which datasets are appropriate given different research goals? Which data curation strategy? Do we need to train a text encoder, or can we leverage a pre-trained LLM? Is a contrastive loss enough for vision understanding or is a generative component key? We also present common techniques used to improve model performance as well as grounding and better alignment.

While providing the recipes for training models is a crucial step for better understanding VLMs' needs, providing robust and reliable evaluation of those models is equally important. Many benchmarks that are used to evaluate VLMs have been introduced recently. However, some of these benchmarks have essential limitations that researchers should be aware of. By discussing the strengths and weaknesses of VLM benchmarks, we hope to shed light on the challenges ahead to improve our understanding of VLMs. We start by discussing the benchmarks that evaluate the visio-linguistic abilities of VLMs, and then we present how to measure biases.

The next generation of VLMs will be able to understand videos by mapping video to language. However, there are different challenges with videos that are not present with images. The computational cost is of course much higher but there are also other considerations on how to map the temporal dimension through text. By shedding light on the current methods that learn from videos, we hope to highlight the current research challenges to tackle on.

By lowering the barrier to entry into VLM research, we hope to provide the foundations for more responsible development of VLMs while pushing the boundaries of vision understanding.

2 The Families of VLMs

Given the impressive progress powered by deep learning in the fields of computer vision and natural language processing, there have been several initiatives to bridge the two domains. In this paper we **focus on the most recent techniques based on transformers** [Vaswani et al., 2017]. We categorize these recent initiatives into four different training paradigms (Figure 1). The first one around **contrastive** training is a commonly used strategy which leverages pairs of positive and negative examples. The VLM is then trained to predict similar representations for the positive pairs while predicting different representations for the negative pairs. The second initiative, **masking**, leverages reconstruction of masked image patches given some unmasked text. Similarly, by masking words in a caption, it is possible to train a VLM to reconstruct those words given an unmasked image. VLMs based

on **pretrained backbones** often leverage open-source LLMs like Llama [Touvron et al., 2023] to learn a mapping between an image encoder (which could also be pre-trained) and the LLM. Learning a mapping between pre-trained models is often less computationally expensive than training text and image encoders from scratch. While most of those approaches leverage intermediate representations or partial reconstructions, **generative VLMs** are trained in such a way that they can generate images or captions. Given the nature of those models, they are often the most expensive to train. We highlight that **these paradigms are not mutually exclusive; many approaches rely on a mix of contrastive, masking, and generative criteria**. For each of these paradigms, we present only one or two models to give the reader some high-level insights on how those models are designed.

2.1 Early work on VLMs based on transformers

By using a transformer architecture [Vaswani et al., 2017], **Bidirectional Encoder Representations from Transformers (BERT)** [Devlin et al., 2019] significantly outperformed all language modelling approaches at that time. Unsurprisingly, researchers have extended BERT to process visual data. Two of them are visual-BERT [Li et al., 2019] and ViL-BERT [Lu et al., 2019] that combine text with images tokens. The models are trained on two objectives: 1) a classical masked modelling task that aims to predict the missing part in a given input; and 2) a sentence-image prediction task that aims to predict if a caption is actually describing an image content. By leveraging these two objectives, the models achieve strong performance across several vision-language tasks, mostly explained by the ability of the transformer model to learn to associate words with visual clues through the attention mechanisms.

2.2 Contrastive-based VLMs

Contrastive-based training is often better explained through an **Energy-Based Models (EBM)** point of view [LeCun et al., 2006] in which a model E_θ , parameterized by θ , is trained to assign low energy to observed variables and high energy to unobserved ones. Data from a target distribution should have low energy while *any other data points* should have higher energy. To train these models, we consider input data x with an energy function $E_\theta(x)$ of parameters θ . The corresponding Boltzman distribution density function to learn can be written as:

$$p_\theta(x) = \frac{e^{-E_\theta(x)}}{Z_\theta}$$

with normalization factor $Z_\theta = \sum_x e^{-E_\theta(x)}$. To estimate the target distribution P_D from which input data are drawn, we can in principle use the traditional maximum likelihood objective:

$$\arg \min_{\theta} \mathbb{E}_{x \sim P_D(x)} [-\log P_\theta(x)]$$

whose gradient is:

$$\frac{\partial \mathbb{E}_{x \sim P_D(x)} [-\log P_\theta(x)]}{\partial \theta} = \mathbb{E}_{x^+ \sim P_D(x)} \frac{\partial E_\theta(x^+)}{\partial \theta} - \mathbb{E}_{x^- \sim P_\theta(x)} \frac{\partial E_\theta(x^-)}{\partial \theta}$$

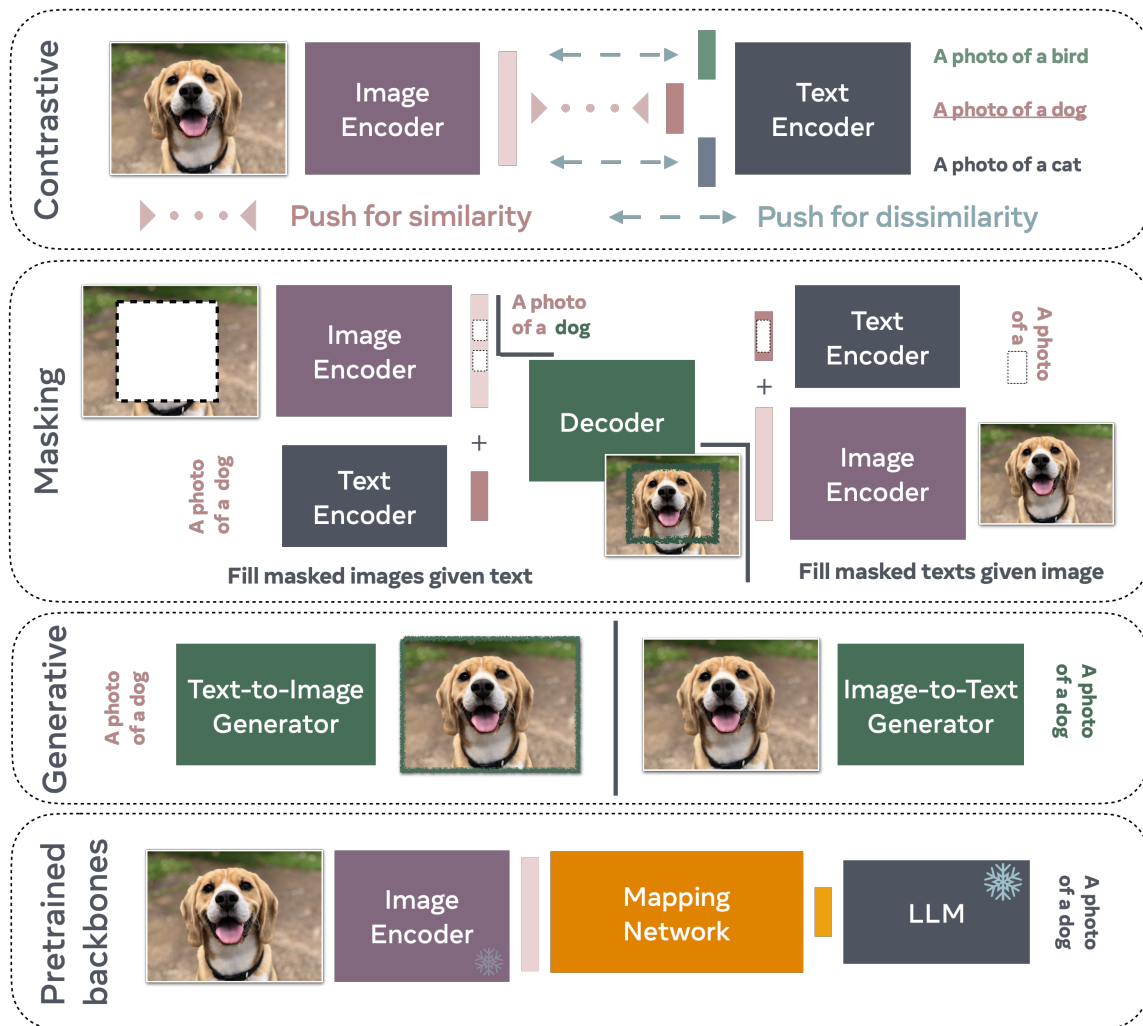


Figure 1: Families of VLMs. **Contrastive** training is a commonly use strategy that uses pairs of positive and negative examples. The VLM is trained to predict similar representations for the positive pairs while predicting different representations for the negative pairs. **Masking** is another strategy that can be leveraged to train VLMs by reconstructing the missing patches given an unmasked text caption. Similarly, by masking words in a caption, it is possible to train a VLM to reconstruct those words given an unmasked image. While most of those approaches leverage intermediate representations or partial reconstructions, **generative** VLMs are trained in such a way they can generate entire images or very long captions. Given the nature of those models, they are often the most expensive to train. **Pretrained backbones** based VLMs often leverage open-source LLMs like Llama to learn a mapping between an image encoder (which could also be pre-trained) and the LLM. It is important to highlight that these paradigms are not mutually exclusive; many approaches rely on a mix of contrastive, masking, and generative criteria.

However, the above requires $x^- \sim P_\theta(x)$, which corresponds to a sample from the model distribution that can be intractable. There are several techniques to approximate such a distribution. One relies on **Markov Chain Monte Carlo (MCMC)** techniques to find examples that minimize the predicted energy through an iterative process. A second one relies on Score Matching [Hyvärinen, 2005] and Denoising Score Matching [Vincent, 2011] criteria which remove the normalization factor by learning only the gradient of the probability density with respect to the input data. Another class of method, on which most of the recent works on Self-Supervised Learning and VLM are based, is **Noise Contrastive Estimation (NCE)** [Gutmann and Hyvärinen, 2010].

Instead of using the model distribution to sample negative examples, the intuition behind NCE is that sampling from a noise distribution $u' \sim p_n(u')$ might approximate samples from the model distribution well enough in certain instances. Even if it can be theoretically difficult to justify why such an approach might work, there is ample empirical evidence of the success of NCE-based methods in recent **Self-Supervised Learning (SSL)** literature [Chen et al., 2020]. The original NCE framework can be described as a binary classification problem in which a model should predict the label $C = 1$ for samples from the real data distribution and $C = 0$ for those coming from the noise distribution. By doing so, the model learns to discriminate between the real data points and the noisy ones. Thus the loss function can be defined as a binary classification with cross-entropy:

$$\mathcal{L}_{\text{NCE}}(\theta) := - \sum_i \log P(C_i = 1|x_i; \theta) - \sum_j \log P(C_j = 0|x_j; \theta) \quad (1)$$

with x_i sampled from the data distribution and $x_j \sim p_n(x), j \neq i$ sampled from the noise distribution.

Wu et al. [2018] introduced NCE without positive pairs with a non-parametric softmax using explicit normalization and a temperature parameter τ . Oord et al. [2018, **CPC**] kept the non-parametric softmax while using positive pairs and coined this approach as **InfoNCE** such that:

$$\mathcal{L}_{\text{InfoNCE}} = - \sum_{(i,j) \in \mathbb{P}} \log \left(\frac{e^{\text{CoSim}(z_i, z_j)/\tau}}{\sum_{k=1}^N e^{\text{CoSim}(z_i, z_k)/\tau}} \right), \quad (2)$$

Instead of predicting a binary value, the InfoNCE loss leverages a distance metric, such as cosine similarity, computed in a model representation space. This requires computing this distance between the positive pairs of examples and between all of the negative pairs of examples. The model learns to predict, through the softmax, the most likely pair of examples that is closest in the representation space while associating lower probability to all other pairs of negative examples. For SSL methods such as SimCLR [Chen et al., 2020], a positive pair of examples is defined as one image and its corresponding handcrafted data-augmented version (such as after applying grayscaleing on the original image) while the negative pairs of examples are built using one image and all other images that are present in a mini-batch. The major drawback of InfoNCE-based methods is the introduction of a dependence on mini-batch content. This often requires large mini-batches to make the contrastive training criterion between the positive and negative samples more effective.

2.2.1 CLIP

A common contrastive method using the InfoNCE loss is **Contrastive Language–Image Pre-training (CLIP)** [Radford et al., 2021]. The positive pairs of examples are defined as one image and its corresponding ground truth caption while the negative examples are defined as the same image but with all the other captions contained in the mini-batch that described the other images. One novelty of CLIP is training a model to incorporate vision and language in a shared representation space. CLIP trains randomly initialized vision and text encoders to map the representation of an image and its caption to similar embedding vectors using a contrastive loss. The original CLIP model trained on 400 million caption-image pairs collected from the web showed remarkable zero-shot classification transfer capabilities. Specifically, a ResNet-101 CLIP matched the performance of a supervised ResNet [He et al., 2015] model (attaining 76.2% zero-shot classification accuracy) and surpassed it on several robustness benchmarks.

SigLIP [Zhai et al., 2023b] is similar to CLIP with the exception that it uses the original NCE loss based on a binary cross-entropy instead of using the CLIP’s multi-class objective based on InfoNCE. This change enables better 0-shot performances on smaller batch sizes than CLIP.

Latent language image pretraining (Llip) [Lavoie et al., 2024] accounts for the fact that an image can be captioned in several different ways. It proposes to condition the encoding of an image on the target caption via a cross-attention module. Accounting for the caption diversity increases the representation’s expressivity and it generally improves the downstream zero-shot transfer classification and retrieval performance.

2.3 VLMs with masking objectives

Masking is a commonly used technique in deep learning research. It can be viewed as a specific form of denoising autoencoder [Vincent et al., 2008] in which the noise has a spatial structure. It is also related to *inpainting* strategies that are notably used by Pathak et al. [2016] to learn strong visual representations. More recently, BERT [Devlin et al., 2019] used **Masked Language Modeling (MLM)** during training to predict missing tokens in a sentence. Masking is particularly well-suited for the transformer architecture [Vaswani et al., 2017] since the tokenization of an input signal makes it easier to randomly drop specific input tokens. There have also been several works on the vision side to learn representations by using **Masked Image Modeling (MIM)** such as MAE [He et al., 2022] or I-JEPA [Assran et al., 2023]. Naturally, there have been works that combined both techniques to train VLMs. A first one is **FLAVA** [Singh et al., 2022] that leverages several training strategies including masking to learn text and image representations. A second one is MaskVLM [Kwon et al., 2023] which is a standalone model. Lastly, we make some connections between information theory and masking strategies.

2.3.1 FLAVA

A first example of the masking-based approach is **Foundational Language And Vision Alignment (FLAVA)** [Singh et al., 2022]. Its architecture comprises three core components,

each based on a transformer framework and tailored to process-specific modalities. The Image Encoder employs the **Vision Transformer (ViT)** [Dosovitskiy et al., 2021] to process images into patches for linear embedding and transformer-based representation, including a classification token ($[\text{CLS}_I]$). The Text Encoder tokenizes textual input using a transformer [Vaswani et al., 2017] and embeds them into vectors for contextual processing and outputting hidden state vectors alongside a classification token ($[\text{CLS}_T]$). Both of those encoders are trained using masking approaches. Building upon these, the Multimodal Encoder fuses hidden states from both the image and text encoders, leveraging learned linear projections and cross-attention mechanisms within the transformer framework to integrate visual and textual information, highlighted by an additional multimodal classification token ($[\text{CLS}_M]$). The model employs a comprehensive training regimen that combines multimodal and unimodal masked modeling losses along with a contrastive objective. It is pretrained on a dataset of 70 million publicly-available image and text pairs. Through this approach, FLAVA demonstrates remarkable versatility and efficacy, achieving state-of-the-art performance across an array of 35 diverse tasks which span vision, language, and multimodal benchmarks, thereby illustrating the model’s ability to understand and integrate information across different domains.

2.3.2 MaskVLM

One limitation of FLAVA is the use of pre-trained vision encoders such as dVAE [Zhang et al., 2019]. To make a VLM that is less dependent on third-party models, Kwon et al. [2023] introduced MaskVLM which applies masking directly in the pixel space and in the text token space. One of the keys to make it work across both text and image is to use the flow of information coming from one modality to the other; the text reconstruction task receives the information coming from the image encoder and vice versa.

2.3.3 Information theoretic view on VLM objectives

Federici et al. [2020] first show that VLMs can be understood to solve a rate-distortion problem, by reducing superfluous information and maximizing predictive information. Dubois et al. [2021] show more specifically, that we can understand any transformation $f(X)$ on data X to implicitly induce an equivalence relationship which partitions the space $f(\mathcal{X})$ into disjoint equivalence classes. We aim to constrain conditional densities to be constant within one region, i.e., $f(x) \sim f(x') \Rightarrow p(z|f(x)) = p(z|f(x'))$, where Z is the learned representation of X . This view unifies masking and other forms of augmentation as well as a choice function between two data modalities; all can be represented as some transformation of the data.

We can formulate the related rate-distortion problem [Shwartz Ziv and LeCun, 2024]:

$$\arg \min_{p(z|x)} I(f(X); Z) + \beta \cdot H(X|Z). \quad (3)$$

To recover the masked VLM objective, we bound Equation (3);

$$\mathcal{L} = - \sum_{x \in \mathcal{D}} \mathbb{E}_{p(f)p(Z|f(x))} [\log q(z) + \beta \cdot \log q(x|z)]. \quad (4)$$

where $\log q(z)$ is an entropy bottleneck, bounding the rate $I(f(X); Z)$, removing superfluous information. Note that the entropy bottleneck in masking VLMs is typically bounded by a constant that depends on the amount information removed by masking. For multimodal VLMs, the amount of information in Z is reduced to the minimum amount of information from either source. The term $\log q(x|z)$ bounds the distortion $H(Z|X)$ and ensures the preservation of information and hence maximizes predictive information. Practically, this term is realized by auto-encoding. In contrast, contrastive losses can be seen as compression without data reconstruction. Here the distortion, see (2), scores the equivalence of two representations. InfoNCE retains the necessary information by classifying which Z is associated with an equivalent example X .

As a result of the information theoretic view, we understand the contrastive loss and auto-encoding loss as implementations of distortions, whereas the rate is mostly determined by the data transformation used.

2.4 Generative-based VLMs

In contrast to previous training paradigms which mostly operate on latent representations to build images or text abstractions that are then mapped between each other, the generative paradigm considers the generation of text and/or images. Some methods like CoCa [Yu et al., 2022b] learn a complete text encoder and decoder which enable image captioning. Some others, like Chameleon Team [2024] and CM3leon [Yu et al., 2023], are multi-modal generative models that are explicitly trained to generate both text and images. Lastly, some models are only trained to generate images based on text such as Stable Diffusion [Rombach et al., 2022], Imagen [Saharia et al., 2022], and Parti [Yu et al., 2022c]. However, even if they are trained to only generate images, they can also be leveraged to solve several vision-language understanding tasks.

2.4.1 An example of learning a text generator: CoCa

Besides the contrastive loss that works well in CLIP, **Contrastive Captioner (CoCa)** [Yu et al., 2022b] also employs a generative loss, which is the loss corresponding to captions generated by a multimodal text decoder that takes in (1) image encoder outputs and (2) representations produced by the unimodal text decoder as inputs. The new loss allows the ability to perform new multimodal understanding tasks (e.g., VQA) without the need for further adaptation using multimodal fusion modules. CoCa is pretrained from scratch by simply treating annotated image labels as text. Pretraining relies on two datasets: ALIGN which contains ~ 1.8 B images with alt-text, as well as JFT-3B which is an internal dataset that consists of >29.5 k classes as labels but treating labels as alt-text.

2.4.2 An example of multi-modal generative model: Chameleon and CM3leon

Yu et al. [2023] introduce CM3Leon, a foundation model for text-to-image and image-to-text generation. CM3Leon borrows the image tokenizer from Gafni et al. [2022] which encodes a 256×256 image into 1024 tokens from a vocabulary of 8192. It borrows the text tokenizer from Zhang et al. [2022] with a vocabulary size of 56320. It introduces a special

token `<break>` to indicate transitions between modalities. This tokenization approach allows the model to process interleaved text and images. The tokenized images and texts are then passed to a decoder-only transformer model [Brown et al., 2020, Zhang et al., 2022] which parameterizes the CM3Leon model.

The CM3Leon model undergoes a two-stage training process. The first stage is retrieval-augmented pretraining. This phase uses a CLIP-based encoder [Radford et al., 2021] as a dense retriever to fetch relevant and diverse multimodal documents and prepends these documents to the input sequence. The model is then trained using next token prediction on the input sequence. The retrieval augmentation effectively increases the tokens available during pretraining thereby increasing data-efficiency. The second stage involves supervised fine-tuning (SFT), where the model undergoes multi-task instruction tuning. This stage allows the model to process and generate content across different modalities, significantly improving its performance on a variety of tasks including text-to-image generation and language-guided image editing. These stages collectively enable CM3Leon to achieve state-of-the-art performance in multi-modal tasks, demonstrating a significant advancement in the capabilities of autoregressive models for handling complex interactions between text and images.

An extension to this work is Chameleon, a new series of mixed-modal foundation models [Team, 2024] that can generate and reason with mixed sequences of interleaved textual and image content. This capability allows for comprehensive multimodal document modeling, extending beyond typical multimodal tasks like image generation, image comprehension, and text-only language models. Chameleon is uniquely designed to be mixed-modal from the beginning, utilizing a uniform architecture trained from scratch in an end-to-end manner on a blend of all modalities—images, text, and code. This integrated approach employs fully token-based representations for both images and text. By converting images into discrete tokens, similar to words in text, the same transformer architecture can be applied to sequences of both image and text tokens without needing separate encoders for each modality. This early-fusion strategy, where all modalities are mapped into a shared representational space from the outset, enables seamless reasoning and generation across different modalities. However, this also introduces significant technical challenges, especially in terms of optimization stability and scaling. These challenges are addressed through a combination of architectural innovations and training techniques, including novel modifications to the transformer architecture such as query-key normalization and revised layer norm placements, which are crucial for stable training in a mixed-modal environment. Additionally, they demonstrate how to adapt supervised fine-tuning approaches used for text-only language models to the mixed-modal context, achieving strong alignment at scale.

2.4.3 Using generative text-to-image models for downstream vision-language tasks

Large advancements have recently been made on language-conditioned image generative models [Bie et al., 2023, Zhang et al., 2023a], from diffusion models like Stable Diffusion [Rombach et al., 2022] and Imagen [Saharia et al., 2022] to autoregressive models like Parti [Yu et al., 2022c]. While the focus has been on their *generative* abilities, they can

actually be directly used for *discriminative* tasks like classification or caption prediction without any retraining.

These generative models are trained to estimate $p_\theta(x | c)$, the conditional likelihood of the image x given a text prompt c . Then, given an image x and a set of n text classes $\{c_i\}_{i=1}^n$, classification can be easily done via Bayes' theorem:

$$p_\theta(c_i | x) = \frac{p(c_i) p_\theta(x | c_i)}{\sum_j p(c_j) p_\theta(x | c_j)} \quad (5)$$

Performing discriminative tasks with conditional generative models is not a new idea – generative classification, or “analysis by synthesis” [Yuille and Kersten, 2006], has been a core idea behind foundational methods like Naive Bayes [Rubinstein et al., 1997, Ng and Jordan, 2001] and linear discriminant analysis [Fisher, 1936]. These generative approaches to classification have traditionally been limited by weak generative modeling capabilities; however, today’s generative models are so good that generative classifiers are becoming competitive again.

Likelihood estimation with autoregressive models. Most state-of-the-art autoregressive models in other modalities (such as language or speech) act on discrete tokens as opposed to raw inputs. This is relatively simple for modalities such as language and speech, which are inherently discrete, but difficult for continuous modalities such as images. In order to effectively leverage techniques from auto-regressive modeling such as LLMs, practitioners generally train an image tokenizer, which maps an image to a sequence of discrete tokens (t_1, \dots, t_K) . After turning an image into a sequence of discrete tokens (e.g., tokenizing the image), estimating the image likelihood is straightforward:

$$\log p_\theta(x | c_i) = \sum_{j=1}^K \log p_\theta(t_j | t_{<j}, c_i) \quad (6)$$

where p_θ is parameterized by the autoregressive VLM. Given that this tokenization is a crucial part of auto-regressive VLMs, one might ask: how do we train image tokenizers? Many current image tokenizers are based on the **Vector Quantised-Variational AutoEncoder (VQ-VAE)** [Van Den Oord et al., 2017] framework, which stitches together an auto-encoder (responsible for creating good compressed *continuous* representations) with a Vector Quantization layer (responsible for mapping continuous representations to discrete representations). The architecture is generally a **Convolutional Neural Network (CNN)** [LeCun and Bengio, 1998] encoder, followed by a Vector Quantization layer, followed by a CNN decoder. The actual discretization step occurs in the vector quantization layer, which maps encoder outputs to the closest embedding in a learned embedding table (“learned” here means that the embedding table is updated throughout training). The loss function for the tokenizer is a combination of reconstruction loss in pixel space (e.g., L2 distance between input and reconstructed pixels) as well as codebook commitment losses to encourage encoder outputs and codebook embeddings to be close to each other. Most modern image tokenizers improve upon this VQ-VAE framework, by either adding different losses or changing the architecture of the encoder/decoder. Notably, VQ-GAN [Esser et al.,

2021] adds perceptual losses and adversarial losses (which involve including a discriminator between ground truth and reconstructed images) to capture more fine-grained details. VIT-VQGAN [Yu et al., 2022a] uses a Vision Transformer instead of CNN for the encoder and decoder architecture.

Likelihood estimation with diffusion models. Obtaining density estimates with diffusion models is more challenging, as they do not directly output $p_\theta(\mathbf{x} \mid \mathbf{c})$. Instead, these networks ϵ_θ are typically trained to estimate the noise ϵ in a noisy image \mathbf{x}_t . Thus, diffusion-based classification techniques [Li et al., 2023a, Clark and Jaini, 2023] estimate a (typically reweighted) variational lower bound for the conditional image likelihood:

$$\log p_\theta(\mathbf{x} \mid \mathbf{c}_i) \propto -\mathbb{E}_{t,\epsilon} [\|\epsilon - \epsilon_\theta(\mathbf{x}_t, \mathbf{c}_i)\|^2] \quad (7)$$

The lower the noise prediction error, the higher the conditional likelihood $p_\theta(\mathbf{x} \mid \mathbf{c})$ is. Measuring the bound in Equation (7) relies on repeated sampling to obtain a Monte Carlo estimate. Li et al. [2023a] and Clark and Jaini [2023] develop techniques for reducing the number of samples required, dynamically allocating samples to the most likely classes and ensuring that the added noise ϵ is matched across all potential classes. However, even with these techniques, classification with conditional diffusion models is still computationally expensive, scaling with the number of classes and requiring hundreds or thousands of network evaluations per test image. Thus, while classification performance with diffusion models is quite good, inference is impractical until further optimizations are developed.

Advantages of generative classifiers. Though inference with these generative classifiers is more expensive, they do have significant advantages. Generative classifiers have more “effective robustness,” which means that they have better out-of-distribution performance for a given in-distribution accuracy [Li et al., 2023a]. On compositional reasoning tasks like Winoground [Thrush et al., 2022], generative classifiers far outperform discriminative methods like CLIP [Li et al., 2023a, Clark and Jaini, 2023]. Generative classifiers, whether autoregressive (Parti) or diffusion-based (Imagen), have been shown to have more shape bias and align better with human judgement [Jaini et al., 2024]. Finally, generative classifiers can be jointly adapted with discriminative models at test-time using only unlabeled test examples [Prabhudesai et al., 2023]. This has been shown to improve performance on classification, segmentation, and depth prediction tasks, especially in online distribution shift scenarios.

2.5 VLMs from Pretrained Backbones

A downside of VLMs is that they are costly to train from scratch. They often require hundreds to thousands of GPUs while having to use hundreds of millions of images and text pairs. Thus, there has been much research work that instead of training models from scratch tried to leverage existing large-language models and/or existing visual extractors. Most of those works are motivated by the fact that many large language models are open-source and thus can be easily used. By leveraging such models, it is possible to then learn a mapping only between the text modality and the image modality. Learning such a mapping enables the LLMs to answer visual questions while requiring a low amount

of compute resources. In this section, we present only two of those models, the first one being Frozen [Tsimpoukelli et al., 2021] which is a first model that leverages pretrained LLMs. Then, we introduce the family of model Mini-GPT [Zhu et al., 2023a].

2.5.1 Frozen

Frozen [Tsimpoukelli et al., 2021] is a first example of a model leveraging a pretrained LLM. This work proposes to connect vision encoders to *frozen* language models through a lightweight mapping network which projects visual features to text token embeddings. The vision encoder (NF-ResNet-50 [Brock et al., 2021]) and the linear mapping are trained from scratch, while the language model (a 7 billion-parameter transformer trained on C4 [Raffel et al., 2020]) is kept frozen (this is crucial to maintain the features that the pre-trained model had already learned). The model is supervised with a simple text generation objective on Conceptual Captions [Sharma et al., 2018b]. At inference time, the language model can be conditioned on interleaved text and image embeddings. The authors show the model is capable of rapid adaptation to new tasks, fast access to general knowledge, and fast binding of visual and linguistic elements. While achieving only modest performance, Frozen has been an important first step toward the current Multimodal LLMs capable of open-ended multimodal zero/few-shot learning.

2.5.2 The example of MiniGPT

Starting from models like Flamingo [Alayrac et al., 2022], a recent trend is to train multimodal language models where the input contains text and images, and the output contains text (and optionally images). MiniGPT-4 [Zhu et al., 2023a] accepts text input and image input, and it produces text output. In MiniGPT-4, a simple linear projection layer is used in order to align image representation (using the same visual encoder in BLIP-2 [Li et al., 2023e], which is based on Q-Former and a ViT backbone) with the input space of the *Vicuna language model* [Chiang et al., 2023]. Given that the visual encoder and Vicuna language model are already pretrained and used as off-the-shelf from prior work, MiniGPT-4 requires only training the linear project layer which is done in two rounds. The first involves 20k training steps (with a batch size of 256), corresponding to around 5M image-text pairs from Conceptual Caption [Sharma et al., 2018b], SBU [Ordonez et al., 2011], and LAION [Schuhmann et al., 2021]. The authors only used four A100 GPUs for around ten hours given that only the linear projection layer parameters needed to be trained. The second round of training leveraged highly-curated data in an instruction-tuning format, needing only 400 training steps (with a batch size of 12).

MiniGPT-5 [Zheng et al., 2023] extends MiniGPT-4 so that the output can contain text interleaved with images. To generate images as well, MiniGPT-5 used generative tokens which are special visual tokens that can be mapped (through transformer layers) to feature vectors, which in turn can be fed into a frozen Stable Diffusion 2.1 model [Rombach et al., 2021]. The authors used supervised training on downstream tasks (e.g., multi-modal dialogue generation and story generation).

LLMs have served as a universal interface for many language-related applications, e.g.,

a general chatbot. Inspired by this, MiniGPT-v2 [Chen et al., 2023b] proposed to perform various vision-language tasks such as image captioning, visual question answering, and object grounding, through a unified interface. To achieve the goal of performing these effectively, MiniGPT-v2 introduced unique identifiers for different tasks when training, enabling the model to distinguish each task instruction effortlessly and also learn efficiently. The experimental results on visual question answering and visual grounding benchmarks show that MiniGPT-v2 demonstrates strong vision-language understanding abilities.

2.5.3 Other popular models using pretrained backbones

Qwen. Similar to MiniGPT-4, Qwen-VL and Qwen-VL-Chat [Bai et al., 2023b] models rely on an LLM, a visual encoder, and a mechanism that aligns the visual representation to the input space of the LLM. In Qwen, the LLM is *initialized from Qwen-7B* [Bai et al., 2023a], the visual encoder is based on ViT-bigG, and a one-layer cross-attention module is used to compress visual representation to a sequence of fixed length (256) which is fed into the LLM.

BLIP2. Li et al. [2023e] introduce BLIP-2, a vision-language model that takes images as input and generates text output. It leverages pretrained, frozen models to greatly shorten training time: a vision encoder (such as CLIP) produces image embeddings that are mapped into the input space of an *LLM such as OPT*. A relatively small (~ 100 -200M parameters) component called a Q-Former is trained for this mapping – it is a Transformer that takes in a fixed number of randomly-initialized “query” vectors; in the forward pass, the queries interact with image embeddings via cross-attention in the Q-Former, followed by a linear layer that projects the queries to the LLM’s input space.

There are many more models based on pretrained LLMs in the literature. Each LLM ends up being extended to a VLM version which means that the scope of a specific survey on such topic would be very large. In this introduction, we aim to present a select few as they all rely on the same principles of learning mappings between representations.

3 A Guide to VLM Training

Several works [Henighan et al., 2020b,a] have shed light on the importance of scaling to push further the performances of deep neural networks. Motivated by these scaling laws, most recent works have focused on increasing compute and scale to learn better models. This led to a model like CLIP [Radford et al., 2021] which was trained on 400M images using a remarkably high compute budget. Even its corresponding open-source implementation, OpenCLIP [Ilharco et al., 2021] was trained using between 256 and 600 GPUs across multiple days or weeks depending on the model size. However, recent work [Sorscher et al., 2022] has shown that it is possible to beat the scaling law using a data curation pipeline. In this section, we first discuss the importance of data when training models and present some of the recipes that are used to create datasets for training VLMs. Then, we discuss the common software, tools and tricks that practitioners might use to train VLMs more efficiently. Since there are different methods to train VLMs, we also discuss

what type of models to choose in specific situations. Lastly, we present some tricks on how to improve grounding (the ability to correctly map text with visual clues). We also introduce techniques to improve alignment using human preferences. VLMs are often used to read and translate text, so we also present some of the techniques that can be used to push further the OCR capabilities of VLMs. Lastly, we discuss the common fine-tuning methods.

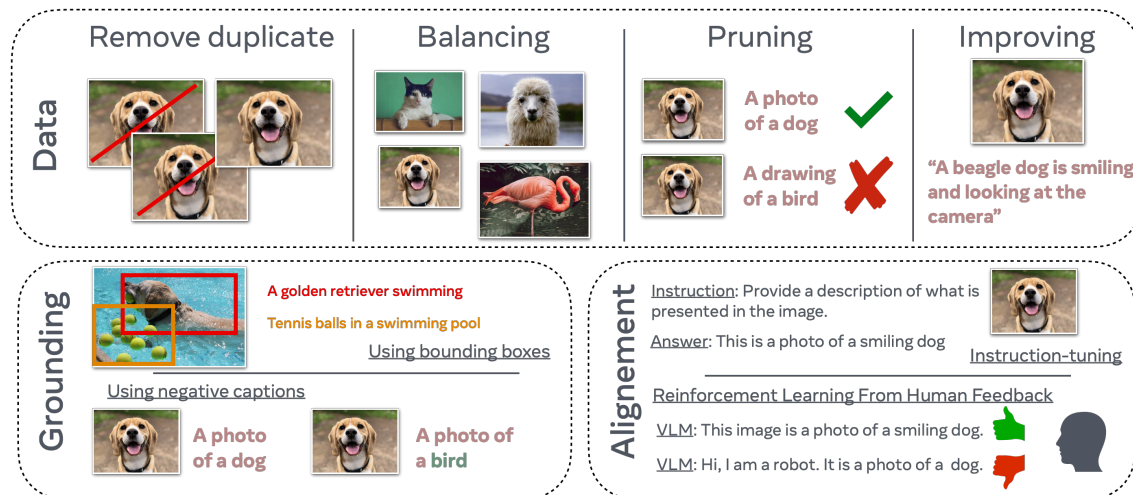


Figure 2: Important considerations to keep in mind when training VLMs. **Data** is one of the most important aspects of training VLMs. Having a diverse and balanced dataset is important for learning good world models that can span enough concepts. It is also important to remove duplicates which occur a lot within large-scale datasets, this will save a lot of compute time and mitigate the risks of memorization. In addition, pruning the data is also an important component since we want to be sure that the captions are indeed related to the image content. Lastly, improving the caption quality is crucial to enhance VLMs performance. **Grounding** VLMs is another important step to ensure that the VLMs correctly associate words with specific concepts. Two common grounding methods leverage either bounding boxes or negative captions. Lastly, **alignment** is a much-needed step to ensure that the model is producing answers that are expected from a human point of view.

3.1 Training data

To evaluate the quality of pretraining datasets, DataComp [Gadre et al., 2023] proposes a benchmark where the model architecture and pretraining hyperparameters of CLIP are fixed. The focus is on designing image-text datasets that achieve strong zero-shot and retrieval performance on 38 downstream tasks. DataComp provides multiple pools of noisy web datasets, ranging from small (1.28 million) to extra-large (12.8 billion) image-text pairs. For each pool, multiple filtering strategies are proposed and evaluated. DataComp demonstrates that data pruning is a crucial step in training highly efficient and performant VLMs.

Data-pruning methods for VLMs can be categorized into three categories: (1) heuristics that eliminate low-quality pairs; (2) bootstrapping methods that utilize pretrained VLMs to rank image-text pairs based on their multimodal alignment, discarding poorly aligned pairs; and finally, (3) methods that aim to create diverse and balanced datasets.

Heuristics: Filters based on heuristics can be further categorized into unimodal and multimodal filters. Unimodal heuristics include removing captions with low text complexity as measured by the number of objects, attributes, and actions [Radenovic et al., 2023a], eliminating non-English alt-text using fastText [Joulin et al., 2017], and removing images based on their resolution and aspect ratio [Gadre et al., 2023]. Multimodal heuristics involve methods that employ image classifiers to filter out image-text pairs for which none of the objects detected in the image map to any of the text tokens [Sharma et al., 2018a]. Additionally, since web-scale datasets often display part of the caption as text in the image, multimodal heuristics, such as text spotting, aim to eliminate image-text pairs with high overlap using off-the-shelf text-spotters [Kuang et al., 2021]. This results in models that learn to extract high-level visual semantics rather than focusing on optical character recognition, thereby preventing low performance on object-centric and scene-centric downstream zero-shot tasks [Radenovic et al., 2023a].

Ranking based on Pretrained VLMs: One of the most effective pruning methods, CLIP-Score [Hessel et al., 2021, Schuhmann et al., 2021], computes the cosine similarity between image and text embeddings using a pretrained CLIP model. This score is then used to rank the alignment of image-text pairs. LAION filtering [Schuhmann et al., 2021] employs an OpenAI CLIP model [Radford et al., 2021] pretrained on 400 million image-text pairs to evaluate the image-text alignment of large web-scale datasets and filter out samples with the lowest CLIPScore. Inspired by text spotting [Radenovic et al., 2023a], T-MARS [Maini et al., 2023] detects and masks text regions in images before computing the CLIPScore, resulting in a more accurate alignment score. Sieve by Mahmoud et al. [2024] demonstrates that false positives and negatives resulting from CLIPScore ranking can be minimized by relying on generative image captioning models pretrained on small but curated datasets.

Diversity and Balancing: Pretraining Vision-Language Models using a diverse and well-balanced dataset can enhance their generalization capabilities [Radford et al., 2021]. To create such a dataset, DataComp [Gadre et al., 2023] suggests sampling image-text pairs that are semantically similar to diverse and curated datasets like ImageNet [Deng et al., 2009]. Text-based sampling retains image-text pairs whose captions overlap with one of the ImageNet classes. Meanwhile, image-based sampling methods encode noisy web-scale images using the OpenAI CLIP ViT-L/14 vision encoder and cluster the images into 100,000 groups using FAISS [Johnson et al., 2019]. Subsequently, embeddings of ImageNet training samples are used to select the closest cluster to each sample. While this approach can result in a diverse dataset, sampling images semantically similar to ImageNet images could bias the CLIP model, potentially limiting its generalization to new downstream tasks. MetaCLIP [Xu et al., 2024] utilizes 500,000 queries from Wikipedia/WordNet as metadata to create a pretraining data distribution that captures a wide range of concepts. Their “balanced” sampling algorithm (similar to the one described in Radford et al. [2021]) aims to strike a balance between well-represented and under-represented concepts, by limiting the

number of samples for each query to 20,000. Nonetheless, collecting a perfectly balanced dataset is impractical due to the natural long-tailed distribution of web data. Consequently, all these CLIP variants still exhibit imbalanced performances across downstream visual concepts [Parashar et al., 2024]. Having a **wide range of training data concepts seems to be one of the most important components behind the “zero-shot abilities”** of VLMs. Actually, Udandarao et al. [2024] demonstrate that the zero-shot performances of VLMs depend mostly on how much those zero-shot downstream concepts are present in the training data.

3.1.1 Improving the training data with synthetic data

A line of research focuses on improving the quality of VLM’s training data by improving the captions through filtering and synthetic data generation. Specifically, **Bootstrapping Language-Image Pre-training (BLIP)** [Li et al., 2022b] performs bootstrapping by generating synthetic samples and filtering out noisy captions. Subsequently, in Santurkar et al. [2022], authors leverage BLIP to approximate the descriptiveness of a caption and show that models trained on consistent and complete synthetic captions generated by BLIP outperform a model trained on human-written captions. Nguyen et al. [2023] use large image-captioning models like BLIP2 [Li et al., 2023e] to replace poorly aligned alt-text labels with descriptive synthetic captions. They demonstrate that pretraining CLIP with a mixture of real and synthetic captions is effective. However, they also show that at scale, the improvement provided by synthetic captions is capped by the limited diversity of generated captions compared to the high diversity of noisy text labels. More recently, Chen et al. [2024] demonstrate that by using **Large Language-and-Vision Assistant (LLaVA)** [Liu et al., 2023d,c, 2024a] as captioning model, it is possible to train very efficiently a text-to-image generative model.

Inspired by the great progress of large-scale diffusion models [Rombach et al., 2022, Dai et al., 2023] and considering the promise of using synthetic image samples in other applications such as classification [Hemmat et al., 2023, Azizi et al., 2023, Bansal and Grover, 2023], another line of research is to use generated images from text-to-image generative models. Tian et al. [2023b] demonstrate improved performance of using synthetic data compared to CLIP [Radford et al., 2021] and SimCLR [Chen et al., 2020] using only synthetic samples. Specifically, they use multiple synthetic samples of the same text prompt as multi-positives pairs for the the contrastive representation learning objective. Furthermore, SynCLR [Tian et al., 2023a] and SynthCLIP [Hammoud et al., 2024] also train a VLM without any real datapoints and only leverage synthetic samples. They use an LLM to generate captions, and then give them to a text-to-image model to generate images based on those captions.

3.1.2 Using data augmentation

Can we exploit data augmentation similarly to self-supervised visual models? SLIP [Mu et al., 2022] addresses this question by introducing an auxiliary self-supervised loss term on the vision encoder. As in SimCLR [Chen et al., 2020], the input image is used to generate two augmentations that create a positive pair to be contrasted with all other

images in the batch. The overhead of this addition is rather small, while providing a regularization term that improves the learned representations. However, the use of the SSL loss only for the visual encoder does not fully exploit the important signal coming from text. To this end, CLIP-rocket [Fini et al., 2023] suggests converting SSL losses to be cross-modal. In particular, it shows that the CLIP contrastive loss can be used in presence of multiple augmentations of the image-text pair, and it is better than other non-contrastive alternatives inspired from SSL, e.g., Grill et al. [2020], Caron et al. [2020], and Zbontar et al. [2021]. In CLIP-rocket, the input image-text pair is augmented in an asymmetrical way, with one weak and one strong set of augmentations. The two resulting augmented pairs are embedded with the standard CLIP encoder and then projected to the multimodal embedding space using two different projectors. The projector of the weakly augmented pair is kept the same as in the original CLIP, i.e., a linear layer, while the projector of the strongly augmented pair is a 2-layer MLP to cope with the noisier embeddings. As highlighted in Bordes et al. [2022] it is crucial to separate the two projectors as the *strong* one is learning more invariant, too invariant, representations for downstream tasks. At inference time, weak and strong learnt representations are interpolated to get a single vector.

3.1.3 Interleaved data curation

Autoregressive language models like Flamingo [Alayrac et al., 2022] and MM1 [McKinzie et al., 2024] have shown including interleaved text and image data during training improves few-shot performance of the model. The interleaved datasets used for pre-training are usually crawled from the internet and are curated to improve quality and safety. There are two types of curation strategies that can be used to collect interleaved datasets:

Natural interleaved data: *Open Bimodal Examples from Large filtered Commoncrawl Snapshots (OBELICS)* [Laurençon et al., 2023] dataset is a good example of this category of datasets; OBELICS is constructed by preserving the intrinsic structure and context in which text and images co-occur within web documents offering a more authentic representation of multimodal web content. Multiple curation steps are used to curate this dataset where English data is collected from common crawl and deduplicated followed by pre-processing HTML document where useful DOM nodes are identified and retained, then for each DOM node we apply image filtering to remove logos followed by a paragraph, and we apply document-level text filtering using various heuristics to handle text that is not well-formed or coherent.

Synthetic interleaved data: MMC4 [Zhu et al., 2023b] is a good example of this type of dataset where text only dataset is retrofitted with images collected from the internet, in this process images are paired with text based on contextual relevance enabled by calculating the CLIP based similarity scores. This method provides a means to retrofit existing vast text corpora with visual information, thereby extending their utility for multimodal learning. While this approach may lack the contextual nuance of naturally interleaved datasets, it allows for the scalable creation of multimodal data from well-established text-only resources.

3.1.4 Assessing multimodal data quality

A very active area for research when it comes to VLMs is to identify the quality of the underlying data used to train it. Since quality is a subjective metric, it's hard to determine a priori what qualifies as good data to train these models. Previous works like Flamingo [Alayrac et al., 2022], MM1 [McKinzie et al., 2024], and OBELICS [Laurençon et al., 2023] have demonstrated that high-quality interleaved multimodal data is a critical requirement for obtaining optimal performance for these VLM models which makes it essential to quantify the quality of the data in a fast and scalable manner. The quality itself could be assessed on multiple fronts incorporating the quality of the text itself, the image itself, and the alignment information between the image and text. Methods like QuRating [Wettig et al., 2024], Data efficient LMs [Sachdeva et al., 2024], and text-quality-based pruning [Sharma et al., 2024] have explored ways to quantify textual data quality and use that to identify high-quality data subsets to train LM models in a data efficient manner. Similarly methods like VILA [Ke et al., 2023] and LAION-aesthetics [Schuhmann, 2023] attempt to quantify the aesthetic quality of an image to select high-quality subsets of image data to improve image generation models. For alignment, the CLIP family of approaches [Radford et al., 2021, Xu et al., 2024, Gao et al., 2024] have been the models of choice to evaluate how coherent the textual data is with respect to the provided image. Despite having some relevant work on evaluating text, image, and alignment quality, we lack a holistic way of evaluating the quality of multimodal and interleaved data, which remains an active area of research to further improve training of VLM models.

3.1.5 Harnessing human expertise: the power of data annotation

In recent years, the importance of leveraging human data annotation has become increasingly evident in advancing the field of vision-language modeling. This approach involves strategically selecting images and having humans provide labels or descriptions that capture the intricate relationship between visual elements and language. By learning from more subtle and detailed information, models can better comprehend complex scenes and generate more accurate descriptions. Although there are several popular multimodal datasets available, such as OKVQA [Marino et al., 2019], A-OKVQA [Schwenk et al., 2022], Image Paragraph Captioning [Krause et al., 2017], VisDial [Das et al., 2017], Visual Spatial Reasoning [Liu et al., 2023a], and MagicBrush [Zhang et al., 2024b], many of these rely on older image benchmarks like COCO [Lin et al., 2014] or Visual Genome [Krishna et al., 2017], which highlights the need for more diverse and contemporary imagery sources. More recently, Urbanek et al. [2023] introduce the DCI dataset which contains fine-grained human annotations for some images from the SA-1B dataset [Kirillov et al., 2023]. A limitation of human-annotated data is that it is often costly to get, especially when requesting fine-grained annotations. In consequence, the number of images with highly detailed annotations is often low which makes often those datasets more suited for evaluation or fine-tuning than for large-scale pre-training.

3.2 Software

In this section, we discuss some of the existing software that people can leverage to evaluate and train VLMs as well as the resources needed to train them.

3.2.1 Using existing public software repositories

There exist several software such as OpenCLIP (https://github.com/mlfoundations/open_clip) or transformers (<https://github.com/huggingface/transformers>) that implement most VLMs. Those tools are extremely useful when making benchmarks or comparing different models. If one's goal is to try and compare different pre-trained VLM on a given downstream task, then those software provide a good platform to do that.

3.2.2 How many GPUs do I need?

The question around the compute resources needed is very important since it will mostly determine the budget one will need to train such model. CLIP [Radford et al., 2021] and OpenCLIP [Ilharco et al., 2021] have leveraged more than 500 GPUs to train their models. When looking at the public cloud prices for such resources, they are equivalent to hundreds of thousands of dollars which is inaccessible to most companies or academic labs. But, when using the right ingredients such as having a high-quality dataset and leveraging masking strategies when using bigger models, training a contrastive model like CLIP on hundreds of millions of images from scratch should not require more than 64 GPUs (which should be equivalent to spending around 10K USD in compute). If the VLM that is used for training leverages existing pre-trained image or text encoder, or LLM, the cost of learning a mapping should be much lower.

3.2.3 Speeding up training

There were recent software developments such as the introduction of torch.compile by the PyTorch team (https://pytorch.org/tutorials/intermediate/torch_compile_tutorial.html) that significantly speed up model training. By using more efficient attention mechanisms, the xformers library [Lefaudeux et al., 2022] is also often used to give an additional speed up. However, there is an area that is often overlooked when training vision models which is data loading. By having to load large mini-batch of images, data loading often becomes a bottleneck that significantly slows down training. In addition, because of space constraint, large-scale datasets are often saved in chunks of tar files that have to be uncompressed on the fly (and thus slowing down training). The main recommendation we have is to store as many uncompressed files as possible to speed up training. In addition, one can leverage the **Fast Forward Computer Vision (FFCV)** library [Leclerc et al., 2023] to create data files that are much faster to load. Using FFCV instead of webdataset can significantly speed up VLM training. The only drawback of storing uncompressed files with either webdataset or FFCV is that the storage might be more costly than storing compressed files. However since the training speed will be much faster, the additional storage cost should be compensated quickly by the lower amount of compute needed.

Masking. Masking is another way to quickly improve the training efficiency of large models. When using models with hundreds of millions or billions of parameters, the cost of a single forward and backward might be high. Li et al. [2023f] show that by randomly masking image tokens one can significantly speed up training time while improving model performances.

3.2.4 Importance of other hyper-parameters.

McKinzie et al. [2024] study the most important design choices for training VLMs showing image resolution, visual encoder capacity, and visual pretraining data are the choices that most impact model performance. They also show while there are many ways to connect modalities, this choice is much less important. The authors also discuss the importance of various types of training data from text-only data to interleaved and image-caption paired data, demonstrating the right mix achieves the best performance across both zero-shot classification and visual-question answering tasks.

3.3 Which model to use?

As highlighted in the first part of this introduction, there are several methods to train VLMs. Some of them leverage simple contrastive training criteria, others use masking strategies to predict missing texts or image patches, while some models are using generative paradigms such as autoregression or diffusion. It is also possible to leverage a pre-trained vision or text backbones like Llama or GPT. In that instance, building a VLM model requires learning only a mapping between the LLM and vision encoder representations. So, from all those methods, which one should someone choose? Do we need to train vision and text encoder from scratch like CLIP or is it better to start with pretrained LLM such as Flamingo or MiniGPT?

3.3.1 When to use contrastive models like CLIP?

Contrastive models like CLIP associate text with visual concepts while keeping a simple training paradigm by pushing text and image representation to be matched in the representation space. By doing so, CLIP learns representations that have both *meaning* in the image and text space, which makes it possible to prompt the CLIP text encoder with words such that we can retrieve the images that map to the corresponding text representations. For example, many data curation pipelines such as MetaCLIP [Xu et al., 2024] are using metadata string matching to build datasets to ensure that each word or concept has enough images associated with them. CLIP models are also a good base for building more complex models, especially when trying to improve grounding. For researchers who are looking at trying additional training criteria or different model architectures to better capture relations or a better understanding of concepts, CLIP is a particularly good starting point. However, one should keep in mind that CLIP is not a generative model, thus it is not possible to generate a caption given a specific image. It is only possible to retrieve the best *caption* within a list of already existing captions. In consequence, current CLIP models cannot be used to provide high-level descriptions of a given image. Another drawback is that CLIP usually needs a very large dataset as well as large batch sizes to offer decent performances, which implies that CLIP usually needs significant resources to be trained from scratch.

3.3.2 When to use masking?

Masking is an alternative strategy to train VLMs. By learning to reconstruct data from both masked images and text, it is possible to jointly model their distributions. In contrast

to contrastive models which operate in a representation space, models based on masking might need to leverage a decoder to map back the representation to the input space (and thus to apply a reconstruction loss). Training an additional decoder might add an additional bottleneck which might make these methods less efficient than a purely contrastive one. However, the advantage is that there is no batch dependency anymore since each example can be considered separately (because we do not need negative examples). Removing negative examples can enable the use of smaller mini-batches without the need to fine-tune additional hyper-parameters such as the softmax temperature. Many VLM methods leverage a mix of masking strategies along with some contrastive loss.

3.3.3 When to use a generative model?

Generative models based on diffusion or autoregressive criteria have demonstrated impressive abilities in generating photorealistic images based on text prompt. Most large-scale training efforts on VLM are also starting to integrate image generation components. Some researchers argue that having the ability to generate images given words is an important step towards creating a good world model while other researchers argue that such a reconstruction step is not needed [Balestrierio and LeCun, 2024]. However from an application perspective, it might be easier to understand and assess what the model has learned when it is able to decode abstract representations in the input data space³. While models like CLIP would need extensive k -NN evaluations using millions of image data points to show what the images closest to a given word embedding look like, generative models can just output the most probable image directly without such an expensive pipeline. In addition, generative models can learn an implicit joint distribution between text and images which might be more suited for learning good representations than leveraging pretrained unimodal encoders. However, they are more computationally expensive to train than their contrastive learning counterpart.

3.3.4 When to use LLM on pretrained backbone?

Using already pretrained text or vision encoder can be a good alternative when having access to limited resources. In that case, only the mapping between the text representation and vision representation should be learned. However, the main issue with this approach is that the VLM will be impacted by the potential hallucination of the LLM. It could also be impacted by any bias coming from the pretrained models. In consequence, there might be an additional overhead in trying to correct the defect of the vision model or of the LLM. Some might argue that it is important to leverage independent image and text encoder to project the information into a lower dimension manifold on which we can learn a mapping while others might argue that it is important to learn the distribution of image and text jointly. To summarize leveraging a pre-trained model is interesting when having limited access to compute resources and when researchers are interested in learning mapping in representation spaces.

³It is also possible to learn a decoder on top of a trained join-embedding architecture [Bordes et al., 2022].

3.4 Improving grounding

Grounding is an important challenge in the VLM and generative model literature. It mostly aims to solve the problem of models not understanding well the text prompt which could either lead to ignoring some part of the prompt or to hallucinating something that is not even part of the prompt. Some of those challenges are related to understanding relations such as an object being on the left or right, negations, counting, or understanding attributes (such as colors or textures). Improving grounding is an active area of research and for now there isn't a single simple method that can solve that. Nevertheless, in this section, we present some of the tricks that are typically used to improve grounding performances.

3.4.1 Using bounding boxes annotations

Models like X-VLM [Zeng et al., 2022] leverage bounding box annotations and incorporate box regression and **Intersection over Union (IoU)** loss to accurately locate and align visual concepts with their corresponding textual descriptions. By knowing where the objects are on the images and what are the captions associated with each object, it is easier for the model to associate text to the right visual clues, and thus improve grounding. X-VLM is trained on a comprehensive collection of datasets, including COCO [Lin et al., 2014], Visual Genome [Krishna et al., 2017], SBU, and Conceptual Captions [Changpinyo et al., 2021], amassing up to 16 million images. This extensive training catalog of data with bounding boxes annotations enables X-VLM to outperform existing methods across a variety of vision-language tasks such as image-text retrieval, visual reasoning, visual grounding, and image captioning.

Instead of using already annotated data, some methods like Kosmos-2 [Peng et al., 2024] rely on public models to create their own image-text datasets. They make a web-scale grounded image-text pairs from web-crawl data by first extracting the nouns from the text captions using spaCy [Honnibal and Montani, 2017] and then use the grounded model GLIP [Li et al., 2022c] to predict bounding boxes associated with the nouns extracted from the captions. Then they use spaCy to extract the expression associated with a given words such that to produce captions that can be associated with each of the bounding boxes that have been detected. Doing so enable the use of very large-scale web-annotated datasets. However such an approach is limited by how strong the grounding model for bounding box detection is. It is likely that if this base model fails on some rare nouns or instances, the downstream model would make similar mistakes.

3.4.2 Negative captioning

Negative samples within the realm of contrastive objectives have been extensively used to mitigate collapse, enhance generalization, and discriminative feature learning [Chen et al., 2020, Liu et al., 2023c, Grill et al., 2020, He et al., 2020, Caron et al., 2021]. By contrasting positive pairs (similar or related samples) with negative pairs (dissimilar or unrelated samples), models are forced to develop a nuanced understanding of the data, going beyond mere superficial features to grasp the underlying patterns that distinguish

different classes or categories.

In the same vein, recent works on VLMs have shown that similar techniques (negative samples) can be adopted to mitigate various problems in vision-language models [Yuksekgonul et al., 2023, Li et al., 2021, Goel et al., 2022, Radenovic et al., 2023b]. For instance, the ARO benchmark [Yuksekgonul et al., 2023] evaluates VLMs on their ability to correctly associate images with captions, using negative samples to test the model’s understanding of incorrect or nonsensical pairings. This approach has demonstrated that VLMs can significantly benefit from the nuanced differentiation capabilities fostered by exposure to negative samples, leading to more accurate and contextually aware models.

3.5 Improving alignment

Motivated by the success of instruction tuning in the language domain [Chung et al., 2024], vision-language models have also begun to incorporate instruction-fine-tuning and **Reinforcement Learning from Human Feedback (RLHF)** in vision-language models to improve multimodal chat capabilities and align outputs with desired responses.

Instruction-tuning involves fine-tuning a vision-language model on supervised data containing instructions, inputs, and the desired response. Typically instruction tuning datasets are much smaller compared to pretraining data—with instruction tuning data sizes ranging from a few to one hundred thousand samples (see Li et al. [2023d] for further discussion of instruction tuning). LLaVa, InstructBLIP [Liu et al., 2023d], and OpenFlamingo [Awadalla et al., 2023] are three prominent vision-language models that incorporate instruction tuning.

RLHF also aims to align model outputs with human preferences. For RLHF a reward model is trained to match human preferences for what humans consider a good or bad model response. While instruction tuning requires supervised training samples, which can be costly to gather, RLHF takes advantage of an auxiliary reward model to mimic human preferences. The primary model, whether a language-only or a vision-language model, is then fine-tuned with the reward model to align outputs with human preferences. LLaVa-RLHF is one prominent example of vision-language models incorporating RLHF to improve model output alignment with factual information [Sun et al., 2023].

3.5.1 A LLaVA story

Motivated by the success of instruction tuning in the language domain, **LLaVA** [Liu et al., 2023d] was among the first models to incorporate instruction-fine-tuning in vision-language models to improve multimodal chat capabilities. The authors generate 150k synthetically generated visual instruction samples for fine-tuning. The original LLaVa model incorporates a pretrained Vicuna language model encoder and a pretrained CLIP ViT-L/14 vision encoder. The encoder outputs are fused into the same dimensional space with a linear projector. Along with improved qualitative chat interactions, LLaVA also shows improvements on synthetic instruction following and Science QA benchmarks [Lu et al., 2022].

LLaVA 1.5. Liu et al. [2023c] improves on LLaVA’s instruction fine-tuning by using a cross-modal fully connected multi-layer perceptron (MLP) layer and incorporating academic VQA instruction data. LLaVA 1.5 is trained on 600k image-text pairs making it much more efficient to train compared to other instruction-tuned models such as InstructBLIP or Qwen-VL. Training takes approximately one day on 8-A100 GPUs. LLaVA 1.5 performs well on a suite of academic VQA and instruction-following benchmarks.

LLaVA-RLHF. Due to the scarcity of high-quality visual instruction tuning data for vision language model training, VLLMs such as LLaVA [Liu et al., 2023d] may misalign the vision and language modalities and generate hallucinated outputs. To address this issue, LLaVA-RLHF [Sun et al., 2023] was proposed to improve multimodal alignment with a novel **RLHF** algorithm, Factually Augmented RLHF. The idea is based on adapting RLHF from text domain to vision-language task and augmenting the reward model with extra factual information of image captions and ground-truth multi-choice to reduce reward hacking. LLaVA-RLHF also uses GPT4-generated training data and human-written image-text pairs for further improving its general capabilities. On LLaVA-Bench, LLaVA-RLHF achieves 94% performance level of GPT-4 [Achiam et al., 2023]. On MMHAL-BENCH with a special focus on penalizing hallucinations, LLaVA-RLHF outperforms baselines by 60%.

LLaVA-NeXT (v1.6). LLaVA-NeXT [Liu et al., 2024a] improves over LLaVA-v1.5 on several fronts. First, the image resolution is increased by concatenating visual features from the full image and smaller image patches, which are separately fed through the vision encoder. Second, the visual instruction tuning data mixture is improved with better visual reasoning, OCR, world knowledge, and logical reasoning examples. Third, the largest model variant uses a 34B-parameter LLM backbone (Nous-Hermes-2-Yi-34B). LLaVA-NeXT achieves state-of-the-art performance compared to open-source multimodal LLMs such as CogVLM [Hong et al., 2023, Wang et al., 2023b] or Yi-VL [AI et al., 2024], and closes the gap with commercial models such as Gemini Pro [Reid et al., 2024].

3.5.2 Multimodal in-context learning

Otter [Li et al., 2023c] shows that *multimodal in-context learning* is possible: A few examples (e.g., instruction-image-answer tuples) are provided as the context, and the model could successfully follow instructions in the test examples without extra fine-tuning. This ability is analogous to text-only LLM in-context learning. The multimodal in-context learning ability can be attributed to fine-tuning on the newly proposed multimodal instruction tuning dataset MIMIC-IT [Li et al., 2023b] that contains around 2.8M multimodal instruction-response pairs with in-context examples. Each sample in MIMIC-IT contains in-context instruction-image-answer tuples as well as a test example (where given the instruction and an image, the goal is to generate the answer in the test example). The in-context tuples are relevant to the test example in one of the three ways: (1) the in-context instructions are similar but the images are different; (2) the images are the same but the instructions are different; (3) the images are in a sequential fashion but the instructions are different, where the sequential images are taken from video repositories like Yang et al. [2023]. Fine-tuning OpenFlamingo [Awadalla et al., 2023] on MIMIC-IT results in the model Otter,

and Otter exhibits stronger instruction following ability as well as multimodal in-context learning ability.

3.6 Improving text-rich image understanding

Understanding text is a crucial aspect of visual perception in our daily lives. The success of **Multimodal Large Language Models (MLLMs)** paved the way for the ability to handle extraordinary applications of VLMs in zero-shot tasks transferred to many real-world scenarios. Liu et al. [2023e] show that MLLMs exhibit excellent zero-shot **Optical Character Recognition (OCR)** performance in the wild, without explicitly training on the OCR domain-specific data. However, these models often struggle with interpreting texts within images when presented with complex relationships between the datatypes, possibly due to the prevalence of natural images in their training data (for instance, Conceptual Captions [Changpinyo et al., 2021] and COCO [Lin et al., 2014]). Some common, non-exhaustive challenges with text understanding and models tackling them:

Instruction tuning with fine-grained text-rich data : LLaVAR [Zhang et al., 2023c]

To address issues with comprehending textual details within an image, LLaVaR enhances the current visual instruction tuning pipeline with text-rich images such as movie posters and book covers. The authors used publicly available OCR tools to collect results on 422K text-rich images from the LAION dataset [Schuhmann et al., 2022]. They then prompted text-only GPT-4 [Achiam et al., 2023] with recognized text and image captions to generate 16K conversations, each containing question-answer pairs for text-rich images. By combining this collected data with previous multimodal instruction-following data, the LLaVAR model was able to substantially improve the capability of the LLaVA model [Liu et al., 2023d]. with up to a 20% accuracy improvement on text-based VQA datasets and a slight improvement on natural images.

Dealing with fine-grained text in high resolution images : Monkey [Li et al., 2023h]

Currently, most MM-LLMs have their input images limited to a resolution of 224 x 224, consistent with the input size of the visual encoder used in their architecture. These models struggle to extract detailed information in complex text-centric tasks such as Scene Text-Centric Visual Question Answering (VQA), Document-Oriented VQA, and Key Information Extraction (KIE) with high-resolution input and detailed scene understanding. To address these challenges, a new approach, Monkey [Li et al., 2023h], has been introduced.

Monkey’s architecture is designed to enhance the capabilities of LLMs by processing input images in uniform patches using a sliding window method, each matching the size used in the original training of the well-trained vision encoder. Each patch is processed independently by a static visual encoder, enhanced with LoRA adjustments and a trainable visual resampler. This allows Monkey to handle higher resolutions up to 1344×896 pixels, enabling the detailed capture of complex visual information. It also employs a multi-level description generation method, enriching the context for scene-object associations. This two-part strategy ensures more effective learning from generated data. By integrating the unique capabilities of these systems, Monkey offers a comprehensive and layered

approach to caption generation, capturing a wide spectrum of visual details.

Decoupled Scene Text Recognition Module and MM-LLM : Lumos [Shenoy et al., 2024] Lumos proposes a multimodal assistant with text understanding capabilities that leverages a combination of on-device and cloud computation. Lumos uses a decoupled **Scene text recognition (STR)** module which then feeds into the multimodal LLM. Lumos' STR module contains four sub-components: **Region of Interest (ROI)** detection, Text detection, Text recognition, and Reading-order reconstruction. ROI detection effectively detects salient areas in the visual, and then crops the salient area as STR input. Text detection takes the cropped image from ROI detection as input, detects words, and outputs the identified bounding box coordinates for each word. Text recognition takes the cropped image from ROI detection and the word bounding box coordinates from Text detection as input, and returns the recognized words. Reading-order reconstruction organizes recognized words into paragraphs and in reading order within each paragraph based on the layout.

The cloud hosts a multimodal LLM module, which takes in the recognized text and coordinates from the STR module. This decoupled STR module can be run on-device, reducing power and latency from transferring high-resolution images to the cloud. As mentioned above, one of the key challenges has been capturing fine-grained text from the scene due to limitations of LLM's encoders. Lumos's STR module works on 3kx4k sized image which would yield enhanced performance in complex text understanding tasks similar to Monkey.

3.7 Parameter-Efficient Fine-Tuning

Training VLMs has shown great effectiveness in cross-domain vision and language tasks. However, as the size of pre-trained models continues to grow, fine-tuning the entire parameter set of these models becomes impractical due to computational constraints. To address this challenge, **Parameter-Efficient Fine-Tuning (PEFT)** methods have been developed to address the high computational cost associated with fine-tuning large-scale models. These methods focus on training a subset of parameters, rather than the entire model, to adapt to downstream tasks. Existing PEFT methods can be categorized into four main groups, namely **Low Rank Adapters (LoRa)** based methods, Prompt-based methods, Adapter-based methods, and Mapping-based methods.

LoRA-based methods. LoRA [Hu et al., 2022] is recognized as a popular method for parameter fine-tuning. LoRA can be applied to both pure language models and vision-language models. Several variants of LoRA have been developed to enhance its functionality and efficiency. One such variant is QLoRA [Dettmers et al., 2023], which integrates LoRA with a quantized backbone and enables the back-propagation of gradients through a frozen, 4-bit quantized pre-trained language model into LoRA. Another variant is VeRA [Kopiczko et al., 2024], which is designed to reduce the number of trainable parameters in comparison to LoRA, while maintaining equivalent performance levels. This is achieved by utilizing a single pair of low-rank matrices shared across all layers and learning small scaling vectors instead. Lastly, DoRA [Liu et al., 2024b] decomposes the

pre-trained weight into two components, magnitude and direction, for fine-tuning. DoRA has demonstrated the capability to generalize Low-rank adaptation methods from language models to Vision-Language benchmarks through empirical experiments.

Prompt-based methods. The process of vision-language pre-training involves the alignment of images and texts within a shared feature space, enabling zero-shot transfer to subsequent tasks via prompting. Consequently, another method for efficient fine-tuning is linked with prompting. Zhou et al. [2022] introduce Context Optimization (CoOp), a technique designed to adapt large pre-trained vision-language models, such as CLIP, for downstream image recognition tasks, eliminating the need for manual prompt engineering. CoOp optimizes the context words of the prompt using learnable vectors during the training process. The method provides two implementations: unified context and class-specific context. Experimental results from 11 datasets indicate that CoOp outperforms hand-crafted prompts and linear probe models in few-shot learning. Additionally, it exhibits superior domain generalization capabilities compared to zero-shot models that utilize manual prompts. Then Jia et al. [2022] present Visual Prompt Tuning (VPT), for adapting large-scale Transformer models in vision. Contrary to the conventional approach of full fine-tuning, which updates all backbone parameters, VPT introduces a minimal amount (less than 1% of model parameters) of trainable parameters in the input space. This is achieved while keeping the model backbone frozen, and in many instances. VPT demonstrates comparable or even superior accuracy to full fine-tuning.

Adapter-based methods. Adapters refer to new modules added between layers of a pre-trained network [Houlsby et al., 2019]. Specifically, in the vision-language model domain, CLIP-Adapter [Gao et al., 2024] fine-tunes with feature adapters on either the visual or language branch. It adopts an additional bottleneck layer to learn new features and performs residual-style feature blending with the original pre-trained features. In addition, VL-adapter [Sung et al., 2022] evaluates various adapter-based methodologies, within a unified multi-task framework across a diverse range of image-text and video-text benchmark tasks. The study further delves into the concept of weight-sharing between tasks as a strategy to augment the efficiency and performance of these adapters. Empirical results indicate that the application of the weight-sharing technique in conjunction with adapters can effectively rival the performance of full fine-tuning, while necessitating updates to only a minimal fraction of the total parameters (4.18% for image-text tasks and 3.39% for video-text tasks). Subsequently, LLaMA-Adapter V2 [Gao et al., 2023] proposes a parameter-efficient visual instruction model that enhances large language models' multi-modal reasoning capabilities without requiring extensive parameters or multi-modal training data. It proposes unlocking more learnable parameters (e.g., norm, bias, and scale) and an early fusion method to incorporate visual tokens into LLM layers. Compared to other full-fine-tuning approaches like MiniGPT-4 and LLaVA, LLaMA-Adapter V2 involves much fewer additional parameters.

Mapping-based methods. Injecting trainable modules into pretrained models through adapters or LoRA requires some knowledge of the network's architecture to decide where to insert or adapt parameters. In the context of VLMs, Mañas et al. [2023] and Merullo et al. [2022] propose a simpler approach which only requires training a mapping between

pretrained unimodal modules (i.e., vision encoders and LLMs), while keeping them completely frozen and free of adapter layers. In addition, this method requires fewer trainable parameters and leads to increased data-efficiency [Vallaeyts et al., 2024]. LiMBer [Merullo et al., 2022] uses a linear layer that projects visual features to have the same LLM hidden state dimension. This projection is independently applied to each feature vector, which means the length of the sequence passed to the LLM is the same as the number of visual feature vectors, increasing the computational cost of training and inference. MAPL [Mañas et al., 2023] designs a mapping network which addresses this issue by aggregating the visual feature vectors into a smaller set. The input feature vectors are projected and concatenated to a sequence of learnable query tokens, and only the outputs of the query tokens are fed to the LLM.

4 Approaches for Responsible VLM Evaluation

As the main ability of VLMs is to map text with images, it is crucial to measure visio-linguistic abilities so as to ensure that the words are actually mapping to visual clues. Early tasks used to evaluate VLMs were image captioning and **Visual Question Answering (VQA)** [Antol et al., 2015]. In this section, we also discuss the task of text-centric VQA that assesses the ability of the model to understand and read text from images. Another common evaluation introduced by Radford et al. [2021] is based on zero-shot predictions such as the ImageNet [Deng et al., 2009] classification task. Such classification tasks are important to assess if a VLM has a good enough knowledge of the world. More recent benchmarks such as Winoground [Thrush et al., 2022] measure visio-linguistic compositional reasoning. Since VLM models are known to display biases or hallucinations, it is important to assess those two components.

4.1 Benchmarking visio-linguistic abilities

A first way to evaluate VLMs is to leverage visio-linguistic benchmarks. These are designed in such a way to assess whether the VLMs are able to associate specific words or phrases with the corresponding visual clues. These benchmarks are at the forefront of VLM evaluation since they assess how well a visio-linguistic mapping is learned. From visual question answering to zero-shot classification, there are many methods that are often used to evaluate VLMs. Some of them focus on the detection of simple visual clues such as “Is a dog visible in the image?” to much more complex scenes in which we would try to assess whether the VLM is able to give the correct answer to questions such as “How many dogs are in the images, and what are they looking at?” By starting from simple captions that highlight clear visual clues to more complex captions that require some level of spatial understanding and reasoning, these benchmarks allow us to assess the strengths and weaknesses of most VLMs.

4.1.1 Image captioning

Introduced by Chen et al. [2015], the COCO captioning dataset and challenge evaluate the caption quality generated by a given VLM. By leveraging an external evaluation server,



Figure 3: Different methods to evaluate VLMs. **Visual Question Answering (VQA)** has been one of the most common methods, though the model and ground truth answers are compared via exact string matching, which may underestimate the model performance. **Reasoning** consists of giving VLMs a list of captions and making it select the most probable one within this list. Two popular benchmarks in this category are Winoground [Diwan et al., 2022] and ARO [Yuksekgonul et al., 2023]. More recently, **Dense** human annotations can be used to evaluate how well the model is able to map the captions to the correct parts of an image [Urbanek et al., 2023]. Lastly, one can use synthetic data like PUG [Bordes et al., 2023] to generate images in different configurations to evaluate VLM robustness to specific variations.

researchers could send the caption generated by their models and have them evaluated by the server that used scores like BLEU [Papineni et al., 2002] or ROUGE [Lin, 2004] to compare the generated caption to a set of reference captions. However, such scores are still heuristics that try to approximate the similarity of those captions. Many works such as Mathur et al. [2020] have advocated for the retirement of scores like BLEU.

To avoid the issue of having to compare a caption with a set of reference captions, Hessel et al. [2021] introduce the CLIPScore that leverages CLIP to predict how close a caption is to an image. The higher the score, the more likely the caption is to actually describe the image content. However, there is a significant limitation of CLIPScore which is the underlying performances of the CLIP model used.

4.1.2 Text-to-image consistency

In addition to evaluating the ability to generate a caption for a given image, one might also want to evaluate the ability to generate an image given a caption. There are end-to-end approaches that use a single model to produce a consistency score. Though it was initially proposed for image captioning, CLIPScore is also used in image generation to measure the alignment between a generated image and a text prompt. Lin et al. [2024b] and Li et al. [2024a] apply another approach that formats the text prompt as a question (e.g., “Does this figure show {text caption}”) and gets the probability of a VQA model answering *yes*. There are also a series of metrics that leverage a **Language Model (LM)** to generate questions given a text caption. TIFA [Hu et al., 2023] and Davidsonian Scene Graph (DSG) [Cho et al., 2024] both use an **LM** to generate natural language binary and multiple choice questions, and a **Visual Question Answering (VQA)** model to evaluate the questions. DSG additionally addresses hallucinations in LLMs and VLMs – the generated questions are organized into a scene graph based on their dependencies and a question is counted as correct if and only if the questions it depends on are also correct. For example, assume a VQA model is given the questions “Is there a car?”, “What color is the car?” and “How many wheels does the car have?”. If the model incorrectly answers “*no*” to the first question, the rest of the questions are deemed incorrect regardless of their answers because the model did not recognize the car. VPEval [Cho et al., 2023] is another metric that also generates questions but instead of being in natural language, the questions are visual programs. These visual programs are executable by different visual modules, such as a counting module, a VQA module or an **Optical Character Recognition (OCR)** module. Lin et al. [2024b] and Li et al. [2024a] introduce VQAScore, another VQA-based method for text-to-image evaluation. Instead of generating questions using an LM, they instead take the text prompt and pass that directly to a VQA model. For instance, given a prompt *a red dog next to blue flower*, VQAScore computes the probability of a VQA model generating *yes* given the question Does this figure show a red dog next to a blue flower?.

4.1.3 Visual question answering

Visual Question Answering (VQA) is the task of answering natural language questions about images. Due to its simplicity and generality, VQA is one of the main tasks used to evaluate VLMs. In fact, most VLM tasks can be reformulated as VQA (e.g., “what is

in the image?” for captioning, “where is this?” for phrase grounding, etc.). The task was originally proposed [Antol et al., 2015] in two flavors: multiple-choice and open-ended answers. Popular benchmarks based on the VQA task include VQAv2 [Goyal et al., 2017], TextVQA [Singh et al., 2019], GQA [Hudson and Manning, 2019], Visual Genome QA [Krishna et al., 2017], VizWiz-QA [Gurari et al., 2018], OK-VQA [Marino et al., 2019], ScienceQA [Lu et al., 2022], MMMU [Yue et al., 2023] (see Figure 3). VQA is traditionally evaluated with VQA Accuracy, which is based on exact string match between a candidate answer generated by a model and a set of reference answers annotated by humans. This metric has worked well so far in the multiple-choice and IID training settings. However, the community is transitioning towards generative models (capable of generating free-form, open-ended answers) and OOD evaluation (e.g., zero-shot transfer). In these new settings, the traditional VQA Accuracy metric is too stringent and tends to underestimate the performance of current VLM systems [Agrawal et al., 2023]. To overcome this limitation, some works have resorted to artificially constraining [Li et al., 2023e] or rephrasing [Awal et al., 2023] the output of VLM to match the format of reference answers. However, this precludes a fair comparison among VLM as their perceived performance is largely dependent on answer formatting tricks. To enable a truthful and fair evaluation of VLM, Mañas et al. [2024] propose to leverage LLMs as judges for VQA.

Selective prediction. Besides answer correctness, another dimension of evaluation is *selective prediction* for VQA – how well a VLM can abstain from answering questions it would otherwise get incorrect, and achieve high accuracy on questions it chooses to answer. This is important for applications where accuracy is critical, and incorrect answers could mislead users who place trust in the model. Whitehead et al. [2022] formalize this framework for VQA, defining evaluation in terms of coverage (the fraction of questions answered) at a specified risk (level of error tolerated), as well as a cost-based metric (Effective Reliability) that penalizes incorrect answers more than abstentions. The decision to abstain can be determined by thresholding an uncertainty measure, such as the answer probability directly, a learned correctness function [Whitehead et al., 2022, Dancette et al., 2023], or agreement among expert models (e.g., Si et al. [2023] in the unimodal language space).

Visual Dialog. Das et al. [2017] introduced VisDial, a dataset and benchmark that extends VQA by using a series of questions about an image. Its goal is to measure the ability of an agent to hold a discussion about a given image. In contrast to traditional VQA in which questions can be considered as independent, visual dialog benchmarks evaluate more general intelligence abilities such as being able to understand context from the discussion history.

4.1.4 Text-centric Visual Question Answering

Text-Based VQA is a task that involves providing responses to natural language inquiries about textual content in the image. Beyond understanding the correlation between textual and visual content in generic VQA, these queries require the model to 1) read the text in the scene accurately and devise how it’s structured and ordered and 2) reason about the text in the image in correlation to each other as well as other visual elements in the image.

Text-centric evaluations can be done using a broad spectrum of tasks like Text Recognition, Scene Text-Centric Visual Question Answering (VQA), Document-Oriented VQA, Key Information Extraction (KIE), and Handwritten Mathematical Expression Recognition (HMER). Each of these tasks presents unique challenges and requirements, providing a comprehensive overview of the capabilities and limitations of the LLMs.

Text Recognition is a fundamental task in **Optical Character Recognition (OCR)**, requiring the model to accurately identify and transcribe text from a variety of sources. Scene Text-Centric VQA extends this challenge by requiring the model to not only recognize text within a scene, but also to answer questions about it. Document-Oriented VQA further complicates this by introducing structured documents, such as forms and invoices, into the mix. KIE is a task that focuses on extracting key pieces of information from a document, such as names, dates, or specific values. Finally, HMER is a specialized task that involves recognizing and transcribing handwritten mathematical expressions, a particularly challenging task due to the complexity and variability of handwritten notation. Some popular benchmarks include IIT5K [Mishra et al., 2012], COCOText [Veit et al., 2016], SVT [Shi et al., 2014], IC13 [Karatzas et al., 2013] for text recognition, STVQA [Biten et al., 2019], Text VQA [Singh et al., 2019], OCR VQA [Mishra et al., 2019] and EST VQA [Wang et al., 2020] for scene text-centric VQA, DocVQA [Mathew et al., 2021], Info VQA [Mathew et al., 2022] and ChartQA [Masry et al., 2022] for document-oriented VQA, SROIE [Huang et al., 2019], FUNSD [Jaume et al., 2019] and POIE [Kuang et al., 2023] for KIE and HME100k [Yuan et al., 2022] for HMER. The composition of the datasets varies widely and should be chosen primarily based on purpose of the evaluation – some focus on specific types of text (such as handwritten or artistic text), while others include a mix of text types. Some datasets were specifically designed to challenge the models' ability to handle multilingual text, handwritten text, non-semantic text, and mathematical expression recognition. Some datasets purely focus on a plethora of different infographics and tabular representations.

4.1.5 Zero-shot image classification

Zero-shot classification consists of evaluating a model on a classification task for which the model was not explicitly trained. This should be contrasted with few-shot learning which requires few training data samples of the downstream task of interest for model fine-tuning. Radford et al. [2021] demonstrate that zero-shot classification performance of CLIP can be significantly improved with different types of prompt structures, especially when customized for specific tasks. They were able to show competitive performances on the well-known ImageNet classification benchmark [Deng et al., 2009]. This was the first work to show that VLM approaches might be able to compete with standard classification training. In addition to ImageNet, it is standard to evaluate VLMs on additional classification datasets such as CIFAR10/100 [Krizhevsky, 2009], Caltech 101 [Li et al., 2022a], Food101 [Bossard et al., 2014], CUB [Wah et al., 2011], StanfordCars [Krause et al., 2013], Eurosat [Helber et al., 2019], Flowers102 [Nilsback and Zisserman, 2008], OxfordPets [Parkhi et al., 2012], FGVC-Aircraft [Maji et al., 2013] and Pascal VOC [Everingham et al., 2010].

Since prompt engineering, e.g., using concept names within human-engineered prompt templates, such as “a photo of a {class}” and “a demonstration of a {class}”, can substantially enhance zero-shot performance, recent studies introduce novel approaches [Menon and Vondrick, 2023, Pratt et al., 2023, Parashar et al., 2023] that employ LLMs like ChatGPT to automatically generate prompts, often with rich visual descriptions, e.g., “a tiger, which has sharp claws”. While these methods adopt label names as originally written by CLIP [Radford et al., 2021], Parashar et al. [2024] substitutes these names with their most frequently used synonyms (e.g., replacing `cash machine` with `ATM`) to improve accuracy, irrespective of the prompt templates employed. As highlighted by Udandarao et al. [2024], zero-shot abilities of a VLM depend mostly on whether those concepts are present or not in the training data. Thus, it is not clear whether we should still consider such evaluations as zero-shot since the model might be already trained in some indirect way to solve the downstream task.

Generalization on Out-Of Distribution (OOD) tasks. Using zero-shot evaluation for CLIP on tasks like ImageNet and achieving good performances is only possible because the CLIP training data is large enough that it may contain much of the concepts and class labels that are present in the ImageNet dataset. In consequence, when there are some downstream tasks for which the training CLIP distribution might be too different, it can lead to poor generalization. Samadh et al. [2023] suggests modifying the token distribution of test examples such that they align with ImageNet data distribution (since the original CLIP training data is unknown). They show that such alignment can help improve performances on various OOD benchmarks as well as on different downstream tasks.

4.1.6 Visio-linguistic compositional reasoning

Several recent benchmarks introduce artificially created captions that are designed with ambiguity to attack the model. One easy way to create such captions can be by reordering the words in the ground-truth caption. Then the model is evaluated on its ability to discriminate the correct caption from the perturbed one (which makes this evaluation equivalent to a binary classification problem). In this section, we are presenting some of the benchmarks that are often used that leverage such binary classification setups.

Winoground [Thrush et al., 2022] is a task for evaluating the visiolinguistic abilities of VLMs. Each example in the dataset contains two images and two captions. Each image matches exactly one caption, with the captions differing only in their word order. For example, in Figure 3, there are two captions “*some plants surrounding a lightbulb*” and “*a lightbulb surrounding some plants*”. Models are tasked with scoring the correct image-caption pairs higher than the incorrect pairs. Diwan et al. [2022] additionally explore Winoground and provide insight on why this task is so challenging for VLMs.

More recently, **Attribution, Relation, and Order (ARO)** was introduced by Yuksekgonul et al. [2023] to assess relation, attribute, and order understanding by VLMs. The dataset was built using GQA, COCO and Flickr30k. Then, negative captions were generated by swapping either the relations, attribute or order from the original caption. By doing so, a

caption describing “A horse eating grass” becomes “grass eating a horse” (Figure 3). Then the model is evaluated on its ability to predict a lower probability to the negative caption. In contrast to Winoground that finds real images that correspond to the negative caption, ARO does not come with true “negative” images. Such an approach has the advantage that it is possible to generate a lot of negative captions; however, some of them might not make any sense in the real world.

Hsieh et al. [2023] have observed that recently developed image-to-text retrieval benchmarks [Yuksekgonul et al., 2023, Zhao et al., 2022, Ma et al., 2023], which are designed to assess the detailed compositional abilities of VLMs, can be manipulated. These benchmarks indeed depend on procedurally-generated hard negatives that often lack logical coherence or fluency due to grammatical inaccuracies. To mitigate these issues, Hsieh et al. [2023] instead suggest leveraging ChatGPT to generate more plausible and linguistically correct hard negatives. They have divided the SUGARCREPE dataset [Hsieh et al., 2023], similar to the approach taken in ARO to evaluate different forms of hard-negatives, each measuring a specific compositional aspect (e.g., attribute, relationship, object understanding).

Warning! A major issue with many of the benchmarks relying on the binary classification problem of discriminating the correct caption from the negative one is that they often do not consider the case in which the model outputs an equal probability for both captions. This can occur if the model collapses the information to the same representation vector for both captions. If the model outputs the same probabilities, then the argmax operation used by frameworks like PyTorch will always return the first element of the vector. It happens that many benchmarks put the correct caption as the first element. Thus, a model whose parameters are all equal to zero could achieve 100% accuracy in these benchmarks. We recommend adding a small epsilon random number or keeping track if the captions are assigned the same probabilities.

4.1.7 Dense captioning and crop-caption matching

The current generation of VLMs is often limited to short text descriptions as input due to text tokenizers. The popular Clip Tokenizer (used to train CLIP-based models) generates only a maximum of 77 tokens, equivalent to fifty English words or a small paragraph. Even if it is possible to summarize an image with few words, images are often much richer than that. When using short captions, we lose information about the background and the fine-grained specifics of the object we want to describe. The **Densely Captioned Images (DCI)** dataset [Urbanek et al., 2023] was introduced to provide complete image descriptions. By dividing an image into distinct parts using Segment Anything [Kirillov et al., 2023], the authors asked human annotators to provide detailed descriptions for each segmented part of this image. Using such an approach, they annotated 7805 images with captions above 1000 words. Using the DCI dataset, the authors evaluated VLMs on a new crop-caption matching task. For each image, the VLM should match the correct caption within all the sub-image captions to the correct sub-image. Doing so allowed the author to evaluate how much a given VLM can have a fine-grained understanding of scene details.

4.1.8 Synthetic data based visio-linguistic evaluations

One of the challenges we encounter when using real data is that it might be hard to find an image that could be associated with a negative caption. In addition, it is difficult to distinguish with these benchmarks if the model is failing because it is not able to recognize a specific object in a specific scene or because despite recognizing both objects, it is not able to recognize the relation between them. In addition, most of the time, the captions that describe images are often extremely simple and might come with ambiguity or biases. Many VLM retrieval-based benchmarks rely on real images extracted from well-known datasets such as COCO. However, using real image datasets that were not designed for VLM evaluation can be problematic since such dataset does not provide images that can be associated with negative captions. For example, a “coffee cup” will always be photographed on top of a table. Consequently, a VLM could leverage this *location bias* to consistently predict the correct positive caption in which the “coffee cup is on top of the table” without using the image information. To avoid such a scenario caused by the bias real images and languages have, it is essential to provide the corresponding image in addition to the negative caption. In the “coffee cup” scenario, it will correspond to having it placed under a table and assessing the capability of the VLM to find the correct spatial location. However, manually placing a real object at various locations would be highly costly since it would require human intervention. In contrast, synthetic image datasets offer unmatched advantages for designing and evaluating VLMs: they make it possible to control each scene precisely and yield granular ground truth labels (and captions). Using **Photorealistic Unreal Graphics (PUG)**, Bordes et al. [2023] granularly constructed complex scenes by adding a single element at a time. By doing so, the authors assess whether a given VLM can associate the correct caption given a background. Then, they add an animal to the scene and verified if the VLM could detect this specific animal in each background. If the animal were correct, they moved it to the left or right to confirm whether the VLM could still find the proper caption indicating whether it was on the left or right. The authors found that current VLMs are not performing better than random chance when evaluating spatial relations.

4.2 Benchmarking Bias and disparities in VLMs

In recent years, biases have been studied heavily across machine learning systems [Buo-lamwini and Gebru, 2018, Corbett-Davies et al., 2017, de Vries et al., 2019]. We now discuss methods for benchmarking biases in VLMs, including analyses of bias via model classifications and their embedding spaces.

4.2.1 Benchmarking bias via classifications

One of the most common ways to benchmark biases in classification models is via *classifications*. For example, biases related to people-related attributes, such as gender, skin tone, and ethnicity, are frequently measured in the context of classifying occupation and profession [Gustafson et al., 2023, Agarwal et al., 2021]. In addition, classifications of people with concepts that allude to harmful associations is frequently evaluated [Agarwal et al., 2021, Goyal et al., 2022, Berg et al., 2022]. Less common, but still relevant, are evaluations of the rate of classification between seemingly benign objects and concepts

such as clothing items or sport equipment when they co-occur with people from different groups [Srinivasan and Bisk, 2021, Agarwal et al., 2021, Hall et al., 2023b].

With real data. These evaluations are commonly done with real data. As an example, Agarwal et al. [2021] perform an evaluation of potential biases in CLIP. They measure representational harms by analyzing the rates of classification for images of faces containing group labels related to race, gender, and age [Karkkainen and Joo, 2021, Hazirbas et al., 2024] to classes like “thief”, “criminal”, and “suspicious person”. Additionally, they measure the distribution of labels related to clothing, appearance, and occupation between gender groups at different thresholds. Among these experiments, they find notable patterns of harmful associations and disparities among race, gender, and age groups.

It is important to be aware of variations in prevalence between groups in real evaluation data sources, as this may affect disparity evaluations. For example, label quality of evaluation data can vary, with potential bias for certain groups or inconsistent concept assignment between groups [Hall et al., 2023a]. Furthermore, there may be distribution shifts among groups, such as the use of different image sources between people with different attributes [Scheuerman et al., 2023].

With synthetic data. Smith et al. [2023] demonstrate that one can evaluate the biases in VLMs using synthetic, gender-balanced contrast sets, generated using diffusion models that only edit the gender-related information and keep the background fixed. Similarly, Wiles et al. [2022] study the failure modes of a model using off-the-shelf image generation and captioning models. Specifically, a generative model is used to generate synthetic samples using ground-truth labels. Then, a captioning model is used to caption misclassified samples to generate additional samples. This results in a corpus of human-interpretable failure modes of the model. Furthermore, Li and Vasconcelos [2024] propose a framework for quantifying the biases in VLMs by applying causal interventions and generating counterfactual image-text pairs. This allows for measuring the discrepancy of the model’s prediction on original and counter-factual distributions.

4.2.2 Benchmarking bias via embeddings

Another approach to benchmarking bias focuses on the *embedding space* of VLMs. Instead of evaluating specific end-tasks like classification, these methods analyze the relationships between the representations of text and images.⁴ Embedding space analyses can unveil learned relationships that are difficult to measure in evaluation tasks. To understand these types of relationships, Ross et al. [2020] introduce two tests embedding association tests, Grounded-WEAT and Grounded-SEAT, that measure biases similar to those found in implicit associations in humans. For instance, they showed that pleasant concepts such as flowers are more associated with European American names and lighter skin than with African Americans and darker skin. Similar nuanced findings are that VLMs

⁴The work largely builds on word embedding analyses in NLP. Word embedding relationships such $\overrightarrow{\text{king}} - \overrightarrow{\text{queen}} \approx \overrightarrow{\text{man}} - \overrightarrow{\text{woman}}$ are semantically useful; however, there are also harmful relationships such as $\overrightarrow{\text{man}} - \overrightarrow{\text{woman}} \approx \overrightarrow{\text{computer programmer}} - \overrightarrow{\text{homemaker}}$ [Bolukbasi et al., 2016].

associate being *American* with being white [Wolfe and Caliskan, 2022] and exhibit sexual objectification [Wolfe et al., 2023]. The explosion of CLIP has brought new approaches that leverage its explicit mapping between text and image embeddings. Demographic biases have been discovered when mapping images to the encoding of demographic attributes (e.g., gender, skin tone, age) and for stereotyped words (e.g., *terrorist*, *CEO*) [Garcia et al., 2023, Hamidieh et al., 2023].

4.2.3 Language biases might impact your benchmark!

As the field of VLMs progresses, it is crucial to address the often overlooked yet critical challenge of curating multimodal benchmarks. A notable example is the influential Visual Question Answering (VQA) benchmark [Antol et al., 2015], which is known to be solvable by “blind” algorithms that exploit unimodal (linguistic) biases in the dataset, e.g., questions starting with “Is there a clock” has the answer “yes” 98% of the time [Goyal et al., 2017]. In other words, multimodal benchmarks that are not carefully curated can be susceptible to unimodal shortcut solutions. Indeed, Lin et al. [2024a] discovers that a blind language prior ($P(\text{text})$) estimated using image-captioning models like BLIP [Li et al., 2022b] perform well on contemporary image-text retrieval benchmarks, including ARO [Yuksekgonul et al., 2023], Crepe [Ma et al., 2023], VL-CheckList [Zhao et al., 2022], and SugarCrepe [Hsieh et al., 2023]. In contrast, balanced benchmarks like Winoground [Thrush et al., 2022] and EqBen [Wang et al., 2023a] actually penalize unimodal shortcuts.

4.2.4 Evaluating how specific concepts in the training data impact downstream performances

Recently, Udandarao et al. [2024] show that concepts that are frequent in the training data will enable good downstream performances on those concepts. However, if those concepts are not present or rare, then the model will perform poorly on those. The authors suggest finding a list of concepts that describe a given downstream task (like class names for classification tasks), and then leverage recognition models (such as RAM [Zhang et al., 2023d]) to detect how much of those concepts are present in the training data. Such evaluation approximates the likelihood for the VLMs to be able to solve those downstream tasks after training.

4.3 Benchmarking hallucinations

Hallucinations are a major concern for LLMs [Huang et al., 2023]. They often produce with very high confidence information that might seem true but that is just false. For example, they can argue that the first time a person was walking on the moon was in 1951 while the true answer was 1969. They can also imagine historical events that just never happen. VLMs could potentially hallucinate text or captions that might not be related to the image a user is asking the model to describe. Thus, assessing if VLMs are not hallucinating is a very important research area. Rohrbach et al. [2018] developed the first benchmark (CHAIR) for object hallucination in captions measuring hallucinations within a fixed object set on COCO [Lin et al., 2014]. While it remains popular, especially for evaluating short, single-sentence captions, it can be misleading for evaluating long

generations from recent VLMs (e.g., counting hypothetical statements as hallucinations, or missing hallucinations that are outside the fixed object set), and is limited to COCO data, which is often included in training sets and offers a narrow view of evaluation on its own. Instead, POPE [Li et al., 2023g] evaluates object hallucination with binary polling questions, both positive (using ground-truth objects) and negative (sampling from negative objects). More recent efforts take model-based approaches to expand evaluation, such as using GPT-4 [Achiam et al., 2023] by Liu et al. [2023b] for evaluating instruction-following (GAVIE), by Zhai et al. [2023a] for localizing object hallucinations in captions (CCEval), and by Sun et al. [2023] for evaluating a VLM’s responses to questions targeting hallucination (MMHal-Bench). Additionally, there is always human evaluation, as Gunjal et al. [2024] demonstrate with fine-grained caption annotations.

4.4 Benchmarking memorization

The potential memorization of training data has been extensively investigated for unimodal models such as LLMs [Carlini et al., 2021] and diffusion models [Somepalli et al., 2023, Carlini et al., 2023]. For VLMs, how to measure memorization is more complex for two main reasons: 1) Unlike generative models, joint embedding VLMs such as CLIP do not come with a decoder, which makes it difficult to decode information memorized in the model’s parameters and learned embeddings. 2) For VLMs such as CoCa and LLaVA that have limited generative capabilities, it remains an open question how to expose cross-modal memorization, e.g., how to probe what the model memorizes about its training image through text.

Jayaraman et al. [2024] study the capability of VLMs to memorize the objects in the training images when queried with their respective captions. They call this phenomenon *déjà vu* memorization and show that CLIP models can effectively “remember” objects present in the training images, even if they are not described in the caption. To do so, the authors propose a k -nearest neighbor test where they utilize a public set of images sampled from the underlying training distribution but has no overlap with the training set. For a target training caption, they find the k public set images closest to the caption in the embedding space. These images are then used to decode the different objects present in target training image. However, this step in itself does not distinguish whether the objects are inferred due to the model memorization or due to the model learning general correlations from image–caption pairs. To distinguish this, the authors train another CLIP model (called the reference model) that has not seen the target image–caption pair during training. A similar k -NN test is then performed on this reference model to evaluate the objects inferred by the reference model. Finally, *déjà vu* memorization is quantified in terms of the gap between the object detection precision/recall scores of the target and reference models, whereby a larger gap indicates a higher degree of memorization.

While different regularization techniques can have varying impact on mitigating the memorization, Jayaraman et al. [2024] find text randomization to be the most effective regularization technique that significantly reduces memorization without severely penalizing the model utility. In this technique, a random fraction of text tokens from training captions are masked in each training epoch. This introduces text augmentation, thereby

reducing the model’s ability to overfit the association between a training caption and its corresponding image.

4.5 Red Teaming

Red teaming in the context of foundation models refers to trying to exploit the public interface of the model to have it generate some undesirable output [Perez et al., 2022]. Red teaming efforts typically include some sort of adversarial dataset aimed at eliciting a harm. The dataset will have a pair of prompts with reference answers deemed correct (e.g., refusal to answer) and the model will be scored based on its distance from the correct answer [Vidgen et al., 2023, Bianchi et al., 2024].

To make things concrete, consider how a VLM may be prompted with a sensitive image and then asked to describe it in graphic detail. While the text prompt could be benign (“describe the activity in this image”), the output could be considered harmful. Work by Li et al. [2024b] attempt to characterize the unique red teaming challenges in terms of faithfulness, privacy, safety, and fairness.

In order to anticipate the kind of challenges in evaluating VLMs, it is helpful to consider some of the red teaming work which has already been developed for text-to-text and text-to-image models. In the language domain, red teaming datasets are crafted to surface certain harms. These harms serve as a proxy for a number of potential risks, which can then be organized into a risk taxonomy [Weidinger et al., 2022, Sun et al., 2024, Derczynski et al., 2023]. To organize these efforts, leaderboards have been developed to benchmark language models across a range of adversarial tasks [Liang et al., 2022, Röttger et al., 2024]. The text-to-image work by Lee et al. [2024] offers a similar ranking effort. To be able to map harms to risks, red teaming efforts fix a definition of the risk they wish to mitigate and then probe the model to try and surface said risk. The formalization of these risks (e.g., privacy, toxicity, bias) remains an active area of research.

After performing a red team evaluation, it can become possible to mitigate certain risks using post-processing methods or model fine-tuning methods, such as Reinforcement Learning for Human Feedback [Ouyang et al., 2022].

5 Extending VLMs to Videos

Our focus so far has been on VLMs that are trained and evaluated on static visual data, i.e., images. However, video data brings new challenges and potentially new capabilities to models, such as understanding the motion and dynamics of objects or localizing objects and actions in space and time. Rapidly, text-to-video retrieval, video question answering and generation emerged as fundamental computer vision tasks [Xu et al., 2015, Tapaswi et al., 2016, Brooks et al., 2024]. The temporal space of video challenges storage, GPU memory and training by a factor of frame rate (e.g., a 24 fps video requires 24× storage/processing, if each frame is considered as an image). This requires trade-offs in VLMs for videos, such as videos in compressed form (e.g., H.264 encoding) with an on-the-

fly video decoder in data loader; initializing video encoders from image encoders; video encoder has spatial/temporal pooling/masking mechanism [Fan et al., 2021, Feichtenhofer et al., 2022]; non-end2end VLMs (extracting video features offline and training models that take video features instead of frames of pixels for long videos). Similar to image-text models, early video-text models trained from scratch the visual and text components with a self-supervised criterion [Alayrac et al., 2016]. But contrary to image models, contrastive video-text models were not the go-to approach, and early fusion and temporal alignment of video and text were preferred [Sun et al., 2019], as more temporal granularity in the representation is more interesting compared to computing a global representation of the video. More recently, a trend similar to image-language models is observed for video-language models: pretrained LLMs are used and aligned with a video encoder, augmenting the LLMs with the capability of video understanding. Modern techniques such as visual instruction tuning are also commonly used and adapted to video.

5.1 Early work on Videos based on BERT

Although the initial approaches to video language were highly specific to the task they were designed to solve, such as video retrieval or video question answering, VideoBERT [Sun et al., 2019] was the first successful general approach to video-language modeling. Contrary to CLIP-based approaches using contrastive learning that were successful for image language modeling, VideoBERT is an early fusion approach, similar to Flamingo [Alayrac et al., 2022], where visual and textual tokens representing video caption pairs are fused together with a single transformer network. The video data comes from YouTube, from instructional cooking videos, and the aligned text is obtained using automatic speech recognition (ASR). The videos are processed frame by frame, each frame corresponding to a single visual token. The pretraining objective is then based on the popular BERT language model, where some tokens are masked and reconstructed. VideoBERT demonstrates strong alignment and is the first model able to perform well on video tasks that require generating text, such as zero-shot action classification and open-ended video captioning.

Going beyond global video and text alignment where a descriptive sentence is matched to a video, **Multimodal Event Representation Learning Over Time (MERLOT)** [Zellers et al., 2021] achieves video language alignment where the text is temporally aligned with the video. Contrary to VideoBERT, which is trained on curated instructional cooking videos, MERLOT is trained on a large-scale dataset of YouTube videos that is less curated and also more diverse, and where the corresponding text is obtained by ASR. The model uses a transformer network trained in a purely self-supervised way, with a contrastive objective between local text tokens and frame visual tokens, a masked language modeling objective, and a temporal reordering objective. The model demonstrated at the time impressive capabilities on question answering tasks, particularly visual common sense reasoning. First, it is able to transfer the knowledge it has learned from videos to answer questions about what is going to happen next from an image, which demonstrates how video models are useful for understanding the visual world. Second, it is able to answer particularly difficult questions from videos on a wide set of datasets and benchmarks. The main limitation of MERLOT is that it lacks the ability to generate text, which prevents it from demonstrating advanced visual reasoning capabilities.

5.2 Enabling text generation using an early-fusion VLM

VideoOFA [Chen et al., 2023c] is an early-fusion VLM for video-to-text generation. Many earlier video VLMs either lack the ability to generate texts, or combine a video encoder with a separately trained text decoder leading to suboptimal accuracy. In contrast, VideoOFA proposes a two-stage pre-training framework to adapt a single generative image-text VLM to video-text tasks. In particular, VideoOFA initializes from an image-text VLM that is capable of text generation and jointly pre-trained on massive image-text data to learn fundamental visual-language representations⁵. It then proposes an *intermediate video-text pre-training* step to adapt the backbone VLM to video-text tasks and learn video-specific concepts such as temporal reasoning. The intermediate pre-training stage consists of three training objectives, all reformulated as video-to-text generation tasks: Video Captioning, Video-Text Matching, and Frame Order Modeling. VideoOFA is evaluated on several Video Captioning and Video Question Answering benchmarks and showed improved performance compared to previous models.

5.3 Using a pretrained LLM

Image-language models progressively converged toward leveraging the power of existing LLMs as their ability to understand text. Instead of training a language model to be aligned with a pre-trained visual backbone, the idea is to align the visual backbone with an existing LLM, often using captioning objectives. The same trend was followed for video models, and Video-LLaMA [Zhang et al., 2023b] emerged as a popular approach, demonstrating strong video-language alignment, both of visual and audio signals. The architecture of Video-LLaMA is based on BLIP-2, a Video Q-former and an Audio Q-former are trained separately on Webvid-2M, a curated dataset of videos, in order to align language with video and audio. The LLM is an LLaMA model and the training objective is a captioning loss. As a second step, the model is fine-tuned on visual instructional data from MiniGPT-4, LLaVA, and VideoChat, making it suitable for human interactions. Video-LLaMA is a conversational agent and is therefore not evaluated with standard benchmarks. The model is accessible through a chat API, where a user can dialogue with the model using text prompts, videos, and images, and ask questions related to it. Many follow-up works such as Video-LLaVA [Lin et al., 2023] further explore LLM alignment with videos.

A more recent one, MiniGPT4-Video [Ataallah et al., 2024] extends MiniGPT-v2 for video comprehension with text input. MiniGPT4-Video adapts the scheme from MiniGPT-v2, concatenating every four adjacent visual tokens into one single token, to reduce the number of input tokens without losing much information. Alongside with visual tokens, text tokens from subtitle of each frame are also extracted for a better representation of each video frame. This mixture of visual tokens and text tokens can facilitate the understanding of the video content for LLMs. The architecture of MiniGPT4-Video consists of a vision encoder, a single linear projection layer, and a large language model. To evaluate the effectiveness of MiniGPT4-Video, three types of benchmarks are used for showing its decent performance for video understanding, including Video-ChatGPT, Open-Ended Questions, and Multiple-Choice Questions (MCQs). MiniGPT4-Video consistently outperforms existing

⁵VideoOFA uses the OFA [Wang et al., 2022] model as its image-text backbone in practice.

state-of-the-art models such as Video-LLaMA [Zhang et al., 2023b] by a large margin on MSVD [Chen and Dolan, 2011], MSRVT [Xu et al., 2016], TGIF [Li et al., 2016], and TVQA [Lei et al., 2018] benchmarks.

5.4 Opportunities in evaluations

While video benchmarks are often similar as image ones, for example captioning, videos also open the door to other types of evaluations. Datasets such as EgoSchema [Mangalam et al., 2024] require the model to answer questions on long videos where interactions between objects/agents must be understood. This enables the evaluation to go beyond describing the scene, which is hard to do on images alone. Similarly, ActivityNet-QA [Yu et al., 2019], MSVD-QA [Xu et al., 2017], and MSRVT-QA [Xu et al., 2017] require to retrieve relevant frames/localize actions to properly answer the questions. However, for a lot of questions looking at a simple frame can be enough to provide accurate answers. For example, showing a football match and asking “Which sport are people playing?” does not require looking beyond a single frame. This raises the question of how much the temporal aspect of the videos is necessary to solve current video benchmarks.

Understanding the semantic aspect of actions in the video is very important, but videos also provide unique opportunities to probe reasoning capabilities or the understanding of the world of the models. To this effect, synthetic data has proven very effective in probing reasoning capabilities of video-based VLMs. In Jassim et al. [2023], videos are generated such that they either follow the laws of physics, or violate them. For example, a ball that suddenly vanishes violates spatio-temporal continuity. Models are then asked if elements in the video, such as the trajectory of a ball, follow the laws of physics. Perhaps surprisingly, models such as VideoLLaMA or PandaGPT [Su et al., 2023] do not exceed random performance, whereas humans achieve more than 80% accuracy. These findings suggest that video VLMs still lack some basic reasoning capabilities that can be probed efficiently thanks to synthetic data.

While the current capabilities of video VLMs are impressive, there are still opportunities to further probe their reasoning capabilities, something possible only by the temporal nature of videos.

5.5 Challenges in leveraging video data

A challenge for video-text pretraining is the current scarcity of (weak) supervision on temporal space, a problem illustrated in VideoPrism [Zhao et al., 2024]. Existing data (e.g., from the Internet) focuses on describing the content of the scenes rather than actions or motion, making a video model downgrade to an image model. CLIP models trained on video can also exhibit a noun bias [Momeni et al., 2023] which makes it harder to model interactions. This yields models that are trained on videos but which are lacking in terms of temporal understanding. Generating paired video-caption data that contain information about the content of the scene as well as temporal aspects is more complex (and costly) than describing the scene in an image. There are possible solutions. For example, a video captioning model can be used to generate more captions, but this requires

an initial high-quality dataset to train this captioner. Another option is to train a video encoder on video alone. This was also exploited for VideoPrism, as it limits the impact of imperfect captions. Beyond data, another challenge is compute. Processing videos is more expensive than images yet it's an even more redundant modality. While an image has a lot of redundant information, two successive frames in a video are even more similar. There is thus a need for more efficient training protocols, for example with masking, a technique that has proven useful on image-based VLMs [Li et al., 2023f]. All of these challenges, whether regarding pretraining data, compute or quality of evaluations, point to promising research directions towards video VLMs with better understanding of the world.

6 Conclusion

Mapping vision to language is still an active research area. From contrastive to generative methods, there are many ways to train VLMs. However, the high compute and data cost is often a barrier for most researchers. This mostly motivates the use of leveraging pre-trained LLMs or image encoders to learn only a mapping between modalities. Whatever the technique to train a VLM might be, there are still general considerations to bear in mind. Large-scale high-quality images and captions are important ingredients to push model performances. Improving model grounding and aligning the model with human preferences are also much needed steps to improve a model's reliability. To assess performances, several benchmarks have been introduced to measure vision-linguistic and reasoning abilities; however, many of them have severe limitations such as being able to be solved only by using language priors. Binding images to text is not the only objective with VLMs; video is also an important modality that can be leveraged to learn representations. However, there are still a lot of challenges to overcome before learning good video representations. Research into VLMs remains very active, as there are still many missing components needed to make these models more reliable.

Acronyms

- ARO** Attribution, Relation, and Order. 36
- BERT** Bidirectional Encoder Representations from Transformers. 6
- BLIP** Bootstrapping Language-Image Pre-training. 19
- CLIP** Contrastive Language-Image Pre-training. 9
- CNN** Convolutional Neural Network. 13
- CoCa** Contrastive Captioner. 11
- DCI** Densely Captioned Images. 37
- EBM** Energy-Based Models. 6
- FFCV** Fast Forward Computer Vision. 22
- FLAVA** Foundational Language And Vision Alignment. 9
- IoU** Intersection over Union. 25
- LLaVA** Large Language-and-Vision Assistant. 19, 26
- LLMs** Large Language Models. 4
- LM** Language Model. 33
- LoRa** Low Rank Adapters. 29
- MCMC** Markov Chain Monte Carlo. 8
- MERLOT** Multimodal Event Representation Learning Over Time. 43
- MIM** Masked Image Modeling. 9
- MLLMs** Multimodal Large Language Models. 28
- MLM** Masked Language Modeling. 9
- NCE** Noise Contrastive Estimation. 8
- OBELICS** Open Bimodal Examples from Large filtered Commoncrawl Snapshots. 20
- OCR** Optical Character Recognition. 28, 33, 35
- OOD** Out-Of Distribution. 36

PEFT Parameter-Efficient Fine-Tuning. 29

PUG Photorealistic Unreal Graphics. 38

RLHF Reinforcement Learning from Human Feedback. 26, 27

ROI Region of Interest. 29

SSL Self-Supervised Learning. 8

STR Scene text recognition. 29

ViT Vision Transformer. 10

VLMs Vision Language Models. 4

VQ-VAE Vector Quantised-Variational AutoEncoder. 13

VQA Visual Question Answering. 31, 33

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. GPT-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 27, 28, 41
- Sandhini Agarwal, Gretchen Krueger, Jack Clark, Alec Radford, Jong Wook Kim, and Miles Brundage. Evaluating CLIP: towards characterization of broader capabilities and downstream implications. *arXiv preprint arXiv:2108.02818*, 2021. 38, 39
- Aishwarya Agrawal, Ivana Kojic, Emanuele Bugliarello, Elnaz Davoodi, Anita Gergely, Phil Blunsom, and Aida Nematzadeh. Reassessing evaluation practices in visual question answering: A case study on out-of-distribution generalization. In *Findings of the Association for Computational Linguistics: EACL 2023*, pages 1171–1196, 2023. 34
01. AI, :, Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, Kaidong Yu, Peng Liu, Qiang Liu, Shawn Yue, Senbin Yang, Shiming Yang, Tao Yu, Wen Xie, Wenhao Huang, Xiaohui Hu, Xiaoyi Ren, Xinyao Niu, Pengcheng Nie, Yuchi Xu, Yudong Liu, Yue Wang, Yuxuan Cai, Zhenyu Gu, Zhiyuan Liu, and Zonghong Dai. Yi: Open foundation models by 01.ai, 2024. 27
- Jean-Baptiste Alayrac, Piotr Bojanowski, Nishant Agrawal, Josef Sivic, Ivan Laptev, and Simon Lacoste-Julien. Unsupervised learning from narrated instruction videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4575–4583, 2016. 43

- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems*, 35:23716–23736, 2022. 15, 20, 21, 43
- Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C. Lawrence Zitnick, and Devi Parikh. VQA: Visual Question Answering. In *International Conference on Computer Vision (ICCV)*, 2015. 31, 34, 40
- Mahmoud Assran, Quentin Duval, Ishan Misra, Piotr Bojanowski, Pascal Vincent, Michael Rabbat, Yann LeCun, and Nicolas Ballas. Self-supervised learning from images with a joint-embedding predictive architecture. In *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 15619–15629, 2023. doi: 10.1109/CVPR52729.2023.01499. 9
- Kirolos Ataallah, Xiaoqian Shen, Eslam Abdelrahman, Essam Sleiman, Deyao Zhu, Jian Ding, and Mohamed Elhoseiny. MiniGPT4-Video: Advancing multimodal llms for video understanding with interleaved visual-textual tokens. *arXiv preprint arXiv:2404.03413*, 2024. 44
- Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, et al. OpenFlamingo: An open-source framework for training large autoregressive vision-language models. *arXiv preprint arXiv:2308.01390*, 2023. 26, 27
- Rabiul Awal, Le Zhang, and Aishwarya Agrawal. Investigating prompting techniques for zero-and few-shot visual question answering. *arXiv preprint arXiv:2306.09996*, 2023. 34
- Shekoofeh Azizi, Simon Kornblith, Chitwan Saharia, Mohammad Norouzi, and David J. Fleet. Synthetic data from diffusion models improves imagenet classification. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL <https://openreview.net/forum?id=DlRsoxjyPm>. 19
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023a. 16
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*, 2023b. 16
- Randall Balestriero and Yann LeCun. Learning by reconstruction produces uninformative features for perception. *arXiv preprint arXiv:2402.11337*, 2024. 24

- Hritik Bansal and Aditya Grover. Leaving reality to imagination: Robust classification via generated datasets. *arXiv preprint arXiv:2302.02503*, 2023. 19
- Hugo Berg, Siobhan Hall, Yash Bhalgat, Hannah Kirk, Aleksandar Shtedritski, and Max Bain. A prompt array keeps the bias away: Debiasing vision-language models with adversarial learning. In Yulan He, Heng Ji, Sujian Li, Yang Liu, and Chua-Hui Chang, editors, *Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 806–822, Online only, November 2022. Association for Computational Linguistics. URL <https://aclanthology.org/2022.aacl-main.61>. 38
- Federico Bianchi, Mirac Suzgun, Giuseppe Attanasio, Paul Rottger, Dan Jurafsky, Tatsunori Hashimoto, and James Zou. Safety-tuned LLaMAs: Lessons from improving the safety of large language models that follow instructions. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=gT5hALch9z>. 42
- Fengxiang Bie, Yibo Yang, Zhongzhu Zhou, Adam Ghanem, Minjia Zhang, Zhewei Yao, Xiaoxia Wu, Connor Holmes, Pareesa Golnari, David A Clifton, et al. Renaissance: A survey into ai text-to-image generation in the era of large model. *arXiv preprint arXiv:2309.00810*, 2023. 12
- Ali Furkan Biten, Ruben Tito, Andres Mafla, Lluís Gomez, Marçal Rusinol, Minesh Mathew, CV Jawahar, Ernest Valveny, and Dimosthenis Karatzas. ICDAR 2019 competition on scene text visual question answering. In *2019 International Conference on Document Analysis and Recognition (ICDAR)*, pages 1563–1570. IEEE, 2019. 35
- Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. Man is to computer programmer as woman is to homemaker? Debiasing word embeddings. *Advances in Neural Information Processing Systems*, 29, 2016. 39
- Florian Bordes, Randall Balestriero, and Pascal Vincent. High fidelity visualization of what your self-supervised representation knows about. *Transactions on Machine Learning Research*, 2022. ISSN 2835-8856. URL <https://openreview.net/forum?id=urfWb7VjmL>. 20, 24
- Florian Bordes, Shashank Shekhar, Mark Ibrahim, Diane Bouchacourt, Pascal Vincent, and Ari Morcos. Pug: Photorealistic and semantically controllable synthetic data for representation learning. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 45020–45054. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/8d352fd0f07fde4a74f9476603b3773b-Paper-Datasets_and_Benchmarks.pdf. 32, 38
- Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101 – mining discriminative components with random forests. In David Fleet, Tomas Pajdla, Bernt Schiele, and Tinne Tuytelaars, editors, *Computer Vision – European Conference on Computer Vision*

- 2014, pages 446–461, Cham, 2014. Springer International Publishing. ISBN 978-3-319-10599-4. [35](#)
- Andy Brock, Soham De, Samuel L Smith, and Karen Simonyan. High-performance large-scale image recognition without normalization. In *International Conference on Machine Learning*, pages 1059–1071. PMLR, 2021. [15](#)
- Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh. Video generation models as world simulators, 2024. URL <https://openai.com/research/video-generation-models-as-world-simulators>. [42](#)
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 33:1877–1901, 2020. [12](#)
- Joy Buolamwini and Timnit Gebru. Gender shades: Intersectional accuracy disparities in commercial gender classification. In Sorelle A. Friedler and Christo Wilson, editors, *Proceedings of the 1st Conference on Fairness, Accountability and Transparency*, volume 81 of *Proceedings of Machine Learning Research*, pages 77–91. PMLR, 23–24 Feb 2018. URL <https://proceedings.mlr.press/v81/buolamwini18a.html>. [38](#)
- Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, et al. Extracting training data from large language models. In *30th USENIX Security Symposium (USENIX Security 21)*, pages 2633–2650, 2021. [41](#)
- Nicolas Carlini, Jamie Hayes, Milad Nasr, Matthew Jagielski, Vikash Sehwal, Florian Tramer, Borja Balle, Daphne Ippolito, and Eric Wallace. Extracting training data from diffusion models. In *32nd USENIX Security Symposium (USENIX Security 23)*, pages 5253–5270, 2023. [41](#)
- Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *Advances in Neural Information Processing Systems*, 33:9912–9924, 2020. [20](#)
- Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9650–9660, 2021. [25](#)
- Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12m: Pushing web-scale image-text pre-training to recognize long-tail visual concepts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3558–3568, 2021. [25](#), [28](#)
- David L. Chen and William B. Dolan. Collecting highly parallel data for paraphrase evaluation. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics (ACL-2011)*, Portland, OR, June 2011. [45](#)

- Fei-Long Chen, Du-Zhen Zhang, Ming-Lun Han, Xiu-Yi Chen, Jing Shi, Shuang Xu, and Bo Xu. Vlp: A survey on vision-language pre-training. *Machine Intelligence Research*, 20(1):38–56, January 2023a. ISSN 2731-5398. doi: 10.1007/s11633-022-1369-5. URL <http://dx.doi.org/10.1007/s11633-022-1369-5>. 4
- Jun Chen, Deyao Zhu, Xiaoqian Shen, Xiang Li, Zechun Liu, Pengchuan Zhang, Raghuraman Krishnamoorthi, Vikas Chandra, Yunyang Xiong, and Mohamed Elhoseiny. MiniGPT-v2: large language model as a unified interface for vision-language multi-task learning. *arXiv preprint arXiv:2310.09478*, 2023b. 16
- Junsong Chen, Jincheng YU, Chongjian GE, Lewei Yao, Enze Xie, Zhongdao Wang, James Kwok, Ping Luo, Huchuan Lu, and Zhenguo Li. Pixart- α : Fast training of diffusion transformer for photorealistic text-to-image synthesis. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=eAKmQPe3m1>. 19
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International Conference on Machine Learning*, pages 1597–1607. PMLR, 2020. 8, 19, 25
- Xilun Chen, Lili Yu, Wenhan Xiong, Barlas Oğuz, Yashar Mehdad, and Wen-tau Yih. VideoOFA: Two-stage pre-training for video-to-text generation. *arXiv preprint arXiv:2305.03204*, 2023c. URL <https://arxiv.org/abs/2305.03204>. 44
- Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C Lawrence Zitnick. Microsoft COCO captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325*, 2015. 31
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, March 2023. URL <https://lmsys.org/blog/2023-03-30-vicuna/>. 15
- Jaemin Cho, Abhay Zala, and Mohit Bansal. Visual programming for text-to-image generation and evaluation. *arXiv preprint arXiv:2305.15328*, 2023. 33
- Jaemin Cho, Yushi Hu, Jason Michael Baldrige, Roopal Garg, Peter Anderson, Ranjay Krishna, Mohit Bansal, Jordi Pont-Tuset, and Su Wang. Davidsonian scene graph: Improving reliability in fine-grained evaluation for text-to-image generation. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=ITq4ZRUT4a>. 33
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. Scaling instruction-finetuned language models. *Journal of Machine Learning Research*, 25(70):1–53, 2024. URL <http://jmlr.org/papers/v25/23-0870.html>. 26

- Kevin Clark and Priyank Jaini. Text-to-image diffusion models are zero shot classifiers. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=fxNQJVMwK2>. 14
- Sam Corbett-Davies, Emma Pierson, Avi Feller, Sharad Goel, and Aziz Huq. Algorithmic decision making and the cost of fairness. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '17, page 797–806, New York, NY, USA, 2017. Association for Computing Machinery. ISBN 9781450348874. doi: 10.1145/3097983.3098095. URL <https://doi.org/10.1145/3097983.3098095>. 38
- Xiaoliang Dai, Ji Hou, Chih-Yao Ma, Sam Tsai, Jialiang Wang, Rui Wang, Peizhao Zhang, Simon Vandenhende, Xiaofang Wang, Abhimanyu Dubey, et al. Emu: Enhancing image generation models using photogenic needles in a haystack. *arXiv preprint arXiv:2309.15807*, 2023. 19
- Corentin Dancette, Spencer Whitehead, Rishabh Maheshwary, Ramakrishna Vedantam, Stefan Scherer, Xinlei Chen, Matthieu Cord, and Marcus Rohrbach. Improving selective visual question answering by learning from your peers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 24049–24059, 2023. 34
- Abhishek Das, Satwik Kottur, Khushi Gupta, Avi Singh, Deshraj Yadav, José MF Moura, Devi Parikh, and Dhruv Batra. Visual dialog. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 326–335, 2017. 21, 34
- Terrance de Vries, Ishan Misra, Changhan Wang, and Laurens van der Maaten. Does object recognition work for everyone? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2019. 38
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009. doi: 10.1109/CVPR.2009.5206848. 18, 31, 35
- Leon Derczynski, Hannah Rose Kirk, Vidhisha Balachandran, Sachin Kumar, Yulia Tsvetkov, MR Leiser, and Saif Mohammad. Assessing language model deployment with risk cards. *arXiv preprint arXiv:2303.18190*, 2023. 42
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. QLoRA: Efficient finetuning of quantized LLMs. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=OUIFPHEgJU>. 29
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and Tamar Solorio, editors, *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL <https://aclanthology.org/N19-1423>. 6, 9

- Anuj Diwan, Layne Berry, Eunsol Choi, David Harwath, and Kyle Mahowald. Why is winoground hard? investigating failures in visuolinguistic compositionality. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang, editors, *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 2236–2250, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.143. URL <https://aclanthology.org/2022.emnlp-main.143>. 32, 36
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=YicbFdNTTy>. 10
- Yifan Du, Zikang Liu, Junyi Li, and Wayne Xin Zhao. A survey of vision-language pre-trained models, 2022. 4
- Yann Dubois, Benjamin Bloem-Reddy, Karen Ullrich, and Chris J Maddison. Lossy compression for lossless prediction. *Advances in Neural Information Processing Systems*, 34:14014–14028, 2021. 10
- Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12873–12883, 2021. 13
- M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The pascal visual object classes (voc) challenge. *International Journal of Computer Vision*, 88(2): 303–338, June 2010. 35
- Haoqi Fan, Bo Xiong, Karttikeya Mangalam, Yanghao Li, Zhicheng Yan, Jitendra Malik, and Christoph Feichtenhofer. Multiscale vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6824–6835, 2021. 43
- Marco Federici, Anjan Dutta, Patrick Forré, Nate Kushman, and Zeynep Akata. Learning robust representations via multi-view information bottleneck. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=B1xwcyHFDr>. 10
- Christoph Feichtenhofer, Haoqi Fan, Yanghao Li, and Kaiming He. Masked autoencoders as spatiotemporal learners. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=UaXD4Al3mdb>. 43
- Enrico Fini, Pietro Astolfi, Adriana Romero-Soriano, Jakob Verbeek, and Michal Drozdal. Improved baselines for vision-language pre-training. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL <https://openreview.net/forum?id=a7nvXxNmdV>. Featured Certification. 20
- Ronald A Fisher. The use of multiple measurements in taxonomic problems. *Annals of Eugenics*, 7(2):179–188, 1936. 13

- Samir Yitzhak Gadre, Gabriel Ilharco, Alex Fang, Jonathan Hayase, Georgios Smyrnis, Thao Nguyen, Ryan Marten, Mitchell Wortsman, Dhruva Ghosh, Jieyu Zhang, Eyal Orgad, Rahim Entezari, Giannis Daras, Sarah M Pratt, Vivek Ramanujan, Yonatan Bitton, Kalyani Marathe, Stephen Mussmann, Richard Vencu, Mehdi Cherti, Ranjay Krishna, Pang Wei Koh, Olga Saukh, Alexander Ratner, Shuran Song, Hannaneh Hajishirzi, Ali Farhadi, Romain Beaumont, Sewoong Oh, Alex Dimakis, Jenia Jitsev, Yair Carmon, Vaishaal Shankar, and Ludwig Schmidt. Datacomp: In search of the next generation of multimodal datasets. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023. URL <https://openreview.net/forum?id=dVaWCDMBof>. 17, 18
- Oran Gafni, Adam Polyak, Oron Ashual, Shelly Sheynin, Devi Parikh, and Yaniv Taigman. Make-a-scene: Scene-based text-to-image generation with human priors. In *European Conference on Computer Vision*, pages 89–106. Springer, 2022. 11
- Peng Gao, Jiaming Han, Renrui Zhang, Ziyi Lin, Shijie Geng, Aojun Zhou, Wei Zhang, Pan Lu, Conghui He, Xiangyu Yue, et al. Llama-adapter v2: Parameter-efficient visual instruction model. *arXiv preprint arXiv:2304.15010*, 2023. 30
- Peng Gao, Shijie Geng, Renrui Zhang, Teli Ma, Rongyao Fang, Yongfeng Zhang, Hongsheng Li, and Yu Qiao. Clip-adapter: Better vision-language models with feature adapters. *International Journal of Computer Vision*, 132(2):581–595, 2024. 21, 30
- Noa Garcia, Yusuke Hirota, Yankun Wu, and Yuta Nakashima. Uncurated image-text datasets: Shedding light on demographic bias. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6957–6966, 2023. 40
- Akash Ghosh, Arkadeep Acharya, Sriparna Saha, Vinija Jain, and Aman Chadha. Exploring the frontier of vision-language models: A survey of current methodologies and future directions, 2024. 4
- Shashank Goel, Hritik Bansal, Sumit Bhatia, Ryan Rossi, Vishwa Vinay, and Aditya Grover. Cyclip: Cyclic contrastive language-image pretraining. *Advances in Neural Information Processing Systems*, 35:6704–6719, 2022. 26
- Priya Goyal, Adriana Romero Soriano, Caner Hazirbas, Levent Sagun, and Nicolas Usunier. Fairness indicators for systematic assessments of visual feature extractors. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, pages 70–88, 2022. 38
- Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. 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*, pages 6904–6913, 2017. 34, 40
- Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent—a new approach to self-supervised learning. *Advances in Neural Information Processing Systems*, 33:21271–21284, 2020. 20, 25

- Anisha Gunjal, Jihan Yin, and Erhan Bas. Detecting and preventing hallucinations in large vision language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 18135–18143, 2024. 41
- Danna Gurari, Qing Li, Abigale J Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo, and Jeffrey P Bigham. Vizwiz grand challenge: Answering visual questions from blind people. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3608–3617, 2018. 34
- Laura Gustafson, Chloe Rolland, Nikhila Ravi, Quentin Duval, Aaron Adcock, Cheng-Yang Fu, Melissa Hall, and Candace Ross. Facet: Fairness in computer vision evaluation benchmark. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 20370–20382, October 2023. 38
- Michael Gutmann and Aapo Hyvärinen. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pages 297–304. JMLR Workshop and Conference Proceedings, March 2010. 8
- Melissa Hall, Bobbie Chern, Laura Gustafson, Denisse Ventura, Harshad Kulkarni, Candace Ross, and Nicolas Usunier. Towards reliable assessments of demographic disparities in multi-label image classifiers, 2023a. 39
- Melissa Hall, Laura Gustafson, Aaron Adcock, Ishan Misra, and Candace Ross. Vision-language models performing zero-shot tasks exhibit disparities between gender groups. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) Workshops*, pages 2778–2785, October 2023b. 39
- Kimia Hamidieh, Haoran Zhang, Thomas Hartvigsen, and Marzyeh Ghassemi. Identifying implicit social biases in vision-language models, 2023. 40
- Hasan Abed Al Kader Hammoud, Hani Itani, Fabio Pizzati, Philip Torr, Adel Bibi, and Bernard Ghanem. Synthclip: Are we ready for a fully synthetic clip training? *arXiv preprint arXiv:2402.01832*, 2024. 19
- Caner Hazirbas, Alicia Sun, Yonathan Efroni, and Mark Ibrahim. The bias of harmful label associations in vision-language models. *arXiv preprint arXiv: 2402.07329*, 2024. 39
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *arXiv preprint arXiv:1512.03385*, 2015. 9
- Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9729–9738, 2020. 25
- Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16000–16009, 2022. 9

- Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth. Eurosat: A novel dataset and deep learning benchmark for land use and land cover classification. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 12(7): 2217–2226, 2019. **35**
- Reyhane Askari Hemmat, Mohammad Pezeshki, Florian Bordes, Michal Drozdal, and Adriana Romero-Soriano. Feedback-guided data synthesis for imbalanced classification. *arXiv preprint arXiv:2310.00158*, 2023. **19**
- Tom Henighan, Jared Kaplan, Mor Katz, Mark Chen, Christopher Hesse, Jacob Jackson, Heewoo Jun, Tom B. Brown, Prafulla Dhariwal, Scott Gray, Chris Hallacy, Benjamin Mann, Alec Radford, Aditya Ramesh, Nick Ryder, Daniel M. Ziegler, John Schulman, Dario Amodei, and Sam McCandlish. Scaling laws for autoregressive generative modeling. *ArXiv*, abs/2010.14701, 2020a. URL <https://api.semanticscholar.org/CorpusID:225094178>. **16**
- Tom Henighan, Jared Kaplan, Mor Katz, Mark Chen, Christopher Hesse, Jacob Jackson, Heewoo Jun, Tom B Brown, Prafulla Dhariwal, Scott Gray, et al. Scaling laws for autoregressive generative modeling. *arXiv preprint arXiv:2010.14701*, 2020b. **16**
- Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. CLIPScore: A reference-free evaluation metric for image captioning. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih, editors, *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7514–7528, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.595. URL <https://aclanthology.org/2021.emnlp-main.595>. **18, 33**
- Wenyi Hong, Weihang Wang, Qingsong Lv, Jiazheng Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan Wang, Yuxiao Dong, Ming Ding, and Jie Tang. Cogagent: A visual language model for gui agents, 2023. **27**
- Matthew Honnibal and Ines Montani. spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing, 2017. **25**
- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Larousilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for NLP. In *International Conference on Machine Learning*, pages 2790–2799. PMLR, 2019. **30**
- Cheng-Yu Hsieh, Jieyu Zhang, Zixian Ma, Aniruddha Kembhavi, and Ranjay Krishna. Sugarcrepe: Fixing hackable benchmarks for vision-language compositionality. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023. URL <https://openreview.net/forum?id=Jsc7WSCzd4>. **37, 40**
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=nZeVKeeFYf9>. **29**

- Yushi Hu, Benlin Liu, Jungo Kasai, Yizhong Wang, Mari Ostendorf, Ranjay Krishna, and Noah A Smith. Tifa: Accurate and interpretable text-to-image faithfulness evaluation with question answering. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 20406–20417, 2023. 33
- Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *arXiv preprint arXiv:2311.05232*, 2023. 40
- Zheng Huang, Kai Chen, Jianhua He, Xiang Bai, Dimosthenis Karatzas, Shijian Lu, and C. V. Jawahar. ICDAR 2019 competition on scanned receipt ocr and information extraction. In *2019 International Conference on Document Analysis and Recognition (ICDAR)*. IEEE, September 2019. doi: 10.1109/icdar.2019.00244. URL <http://dx.doi.org/10.1109/ICDAR.2019.00244>. 35
- Drew A Hudson and Christopher D Manning. GQA: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6700–6709, 2019. 34
- Aapo Hyvärinen. Estimation of non-normalized statistical models by score matching. *Journal of Machine Learning Research*, 6(Apr):695–709, 2005. 8
- Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali Farhadi, and Ludwig Schmidt. Openclip, July 2021. URL <https://doi.org/10.5281/zenodo.5143773>. If you use this software, please cite it as below. 16, 22
- Priyank Jaini, Kevin Clark, and Robert Geirhos. Intriguing properties of generative classifiers. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=rmg0qMKYRQ>. 14
- Serwan Jassim, Mario Holubar, Annika Richter, Cornelius Wolff, Xenia Ohmer, and Elia Bruni. Grasp: A novel benchmark for evaluating language grounding and situated physics understanding in multimodal language models. *arXiv preprint arXiv:2311.09048*, 2023. 45
- Guillaume Jaume, Hazim Kemal Ekenel, and Jean-Philippe Thiran. Funsd: A dataset for form understanding in noisy scanned documents. In *2019 International Conference on Document Analysis and Recognition Workshops (ICDARW)*, volume 2, pages 1–6. IEEE, 2019. 35
- Bargav Jayaraman, Chuan Guo, and Kamalika Chaudhuri. Déjà vu memorization in vision-language models. *arXiv preprint arXiv:2402.02103*, 2024. 41
- Menglin Jia, Luming Tang, Bor-Chun Chen, Claire Cardie, Serge Belongie, Bharath Hariharan, and Ser-Nam Lim. Visual prompt tuning. In *European Conference on Computer Vision*, pages 709–727. Springer, 2022. 30
- Jeff Johnson, Matthijs Douze, and Hervé Jégou. Billion-scale similarity search with gpus. *IEEE Transactions on Big Data*, 7(3):535–547, 2019. 18

- Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. Bag of tricks for efficient text classification. In Mirella Lapata, Phil Blunsom, and Alexander Koller, editors, *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 427–431, Valencia, Spain, April 2017. Association for Computational Linguistics. URL <https://aclanthology.org/E17-2068>. 18
- Dimosthenis Karatzas, Faisal Shafait, Seiichi Uchida, M. Iwamura, Lluís Gómez i Bigorda, Sergi Robles Mestre, Joan Mas Romeu, David Fernández Mota, Jon Almazán, and Lluís-Pere de las Heras. ICDAR 2013 robust reading competition. In *2013 12th International Conference on Document Analysis and Recognition*, pages 1484–1493, 2013. 35
- Kimmo Karkkainen and Jungseock Joo. Fairface: Face attribute dataset for balanced race, gender, and age for bias measurement and mitigation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1548–1558, 2021. 39
- Junjie Ke, Keren Ye, Jiahui Yu, Yonghui Wu, Peyman Milanfar, and Feng Yang. VILA: Learning image aesthetics from user comments with vision-language pretraining. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10041–10051, 2023. 21
- Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollar, and Ross Girshick. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 4015–4026, October 2023. 21, 37
- Dawid Jan Kopiczko, Tijmen Blankevoort, and Yuki M Asano. VeRA: Vector-based random matrix adaptation. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=NjNfLdxr3A>. 29
- Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *Proceedings - 2013 IEEE International Conference on Computer Vision Workshops, ICCVW 2013*, Proceedings of the IEEE International Conference on Computer Vision, pages 554–561, United States, 2013. Institute of Electrical and Electronics Engineers Inc. ISBN 9781479930227. doi: 10.1109/ICCVW.2013.77. 2013 14th IEEE International Conference on Computer Vision Workshops, ICCVW 2013 ; Conference date: 01-12-2013 Through 08-12-2013. 35
- Jonathan Krause, Justin Johnson, Ranjay Krishna, and Li Fei-Fei. A hierarchical approach for generating descriptive image paragraphs. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 317–325, 2017. 21
- Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International Journal of Computer Vision*, 123:32–73, 2017. 21, 25, 34
- Alex Krizhevsky. Learning multiple layers of features from tiny images, 2009. URL <https://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf>. 35

- Jianfeng Kuang, Wei Hua, Dingkang Liang, Mingkun Yang, Deqiang Jiang, Bo Ren, and Xiang Bai. Visual information extraction in the wild: Practical dataset and end-to-end solution. In Gernot A. Fink, Rajiv Jain, Koichi Kise, and Richard Zanibbi, editors, *Document Analysis and Recognition – ICDAR 2023*, pages 36–53, Cham, 2023. Springer Nature Switzerland. ISBN 978-3-031-41731-3. 35
- Zhanghui Kuang, Hongbin Sun, Zhizhong Li, Xiaoyu Yue, Tsui Hin Lin, Jianyong Chen, Huaqiang Wei, Yiqin Zhu, Tong Gao, Wenwei Zhang, et al. MMOCR: a comprehensive toolbox for text detection, recognition and understanding. In *Proceedings of the 29th ACM International Conference on Multimedia*, pages 3791–3794, 2021. 18
- Gukyeong Kwon, Zhaowei Cai, Avinash Ravichandran, Erhan Bas, Rahul Bhotika, and Stefano Soatto. Masked vision and language modeling for multi-modal representation learning. In *The Eleventh International Conference on Learning Representations, 2023*. URL <https://openreview.net/forum?id=ZhuXksSJYWn>. 9, 10
- Hugo Laurençon, Lucile Saulnier, Leo Tronchon, Stas Bekman, Amanpreet Singh, Anton Lozhkov, Thomas Wang, Siddharth Karamcheti, Alexander M Rush, Douwe Kiela, Matthieu Cord, and Victor Sanh. OBELICS: An open web-scale filtered dataset of interleaved image-text documents. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track, 2023*. URL <https://openreview.net/forum?id=SKN2hf1BIZ>. 20, 21
- Samuel Lavoie, Polina Kirichenko, Mark Ibrahim, Mahmoud Assran, Andrew Gordon Wildon, Aaron Courville, and Nicolas Ballas. Modeling caption diversity in contrastive vision-language pretraining. *arXiv preprint arXiv:2405.00740*, 2024. 9
- Guillaume Leclerc, Andrew Ilyas, Logan Engstrom, Sung Min Park, Hadi Salman, and Aleksander Mądry. FFCV: Accelerating training by removing data bottlenecks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12011–12020, 2023. 22
- Yann LeCun and Yoshua Bengio. *Convolutional Networks for Images, Speech, and Time Series*, page 255–258. MIT Press, Cambridge, MA, USA, 1998. ISBN 0262511029. 13
- Yann LeCun, Sumit Chopra, Raia Hadsell, Marc Aurelio Ranzato, and Fu Jie Huang. A tutorial on energy-based learning, 2006. 6
- Tony Lee, Michihiro Yasunaga, Chenlin Meng, Yifan Mai, Joon Sung Park, Agrim Gupta, Yunzhi Zhang, Deepak Narayanan, Hannah Teufel, Marco Bellagente, et al. Holistic evaluation of text-to-image models. *Advances in Neural Information Processing Systems*, 36, 2024. 42
- Benjamin Lefaudeux, Francisco Massa, Diana Liskovich, Wenhan Xiong, Vittorio Caggiano, Sean Naren, Min Xu, Jieru Hu, Marta Tintore, Susan Zhang, Patrick Labatut, Daniel Haziza, Luca Wehrstedt, Jeremy Reizenstein, and Grigory Sizov. xformers: A modular and hackable transformer modelling library. <https://github.com/facebookresearch/xformers>, 2022. 22

- Jie Lei, Licheng Yu, Mohit Bansal, and Tamara L Berg. Tvqa: Localized, compositional video question answering. In *EMNLP*, 2018. 45
- Alexander C. Li, Mihir Prabhudesai, Shivam Duggal, Ellis Brown, and Deepak Pathak. Your diffusion model is secretly a zero-shot classifier. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 2206–2217, October 2023a. 14
- Baiqi Li, Zhiqiu Lin, Deepak Pathak, Jiayao Li, Yixin Fei, Kewen Wu, Xide Xia, Pengchuan Zhang, Graham Neubig, and Deva Ramanan. Evaluating and improving compositional text-to-visual generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2024a. 33
- Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Fanyi Pu, Jingkang Yang, Chunyuan Li, and Ziwei Liu. MIMIC-IT: Multi-modal in-context instruction tuning. *arXiv preprint arXiv:2306.05425*, 2023b. 27
- Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Jingkang Yang, and Ziwei Liu. Otter: A multi-modal model with in-context instruction tuning. *arXiv preprint arXiv:2305.03726*, 2023c. 27
- Chen Li, Yixiao Ge, Dian Li, and Ying Shan. Vision-language instruction tuning: A review and analysis. *arXiv preprint arXiv: 2311.08172*, 2023d. 26
- Fei-Fei Li, Marco Andreeto, Marc’Aurelio Ranzato, and Pietro Perona. Caltech 101, avr 2022a. 35
- Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. BLIP: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning*, pages 12888–12900. PMLR, 2022b. 19, 40
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. BLIP-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett, editors, *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 19730–19742. PMLR, 23–29 Jul 2023e. URL <https://proceedings.mlr.press/v202/li23q.html>. 15, 16, 19, 34
- Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. VisualBERT: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*, 2019. 6
- Liunian Harold Li, Pengchuan Zhang, Haotian Zhang, Jianwei Yang, Chunyuan Li, Yiwu Zhong, Lijuan Wang, Lu Yuan, Lei Zhang, Jenq-Neng Hwang, et al. Grounded language-image pre-training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10965–10975, 2022c. 25
- Mukai Li, Lei Li, Yuwei Yin, Masood Ahmed, Zhenguang Liu, and Qi Liu. Red teaming visual language models. *arXiv preprint arXiv:2401.12915*, 2024b. 42

- Yangguang Li, Feng Liang, Lichen Zhao, Yufeng Cui, Wanli Ouyang, Jing Shao, Fengwei Yu, and Junjie Yan. Supervision exists everywhere: A data efficient contrastive language-image pre-training paradigm. *arXiv preprint arXiv:2110.05208*, 2021. 26
- Yanghao Li, Haoqi Fan, Ronghang Hu, Christoph Feichtenhofer, and Kaiming He. Scaling language-image pre-training via masking. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 23390–23400, 2023f. 22, 46
- Yi Li and Nuno Vasconcelos. Debias your VLM with counterfactuals: A unified approach, 2024. URL <https://openreview.net/forum?id=xx05gm7oQw>. 39
- Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Xin Zhao, and Ji-Rong Wen. Evaluating object hallucination in large vision-language models. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 292–305, Singapore, December 2023g. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.20. URL <https://aclanthology.org/2023.emnlp-main.20>. 41
- Yuncheng Li, Yale Song, Liangliang Cao, Joel Tetreault, Larry Goldberg, Alejandro Jaimes, and Jiebo Luo. TGIF: A New Dataset and Benchmark on Animated GIF Description. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016. 45
- Zhang Li, Biao Yang, Qiang Liu, Zhiyin Ma, Shuo Zhang, Jingxu Yang, Yabo Sun, Yuliang Liu, and Xiang Bai. Monkey: Image resolution and text label are important things for large multi-modal models. *arXiv preprint arXiv:2311.06607*, 2023h. 28
- Paul Pu Liang, Amir Zadeh, and Louis-Philippe Morency. Foundations & trends in multimodal machine learning: Principles, challenges, and open questions. *ACM Comput. Surv.*, apr 2024. ISSN 0360-0300. doi: 10.1145/3656580. URL <https://doi.org/10.1145/3656580>. Just Accepted. 4
- Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, et al. Holistic evaluation of language models. *arXiv preprint arXiv:2211.09110*, 2022. 42
- Bin Lin, Yang Ye, Bin Zhu, Jiayi Cui, Munan Ning, Peng Jin, and Li Yuan. Video-LLaVA: Learning united visual representation by alignment before projection. *arXiv preprint arXiv:2311.10122*, 2023. 44
- Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain, July 2004. Association for Computational Linguistics. URL <https://aclanthology.org/W04-1013>. 33
- 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, 2014. 21, 25, 28, 40
- Zhiqiu Lin, Xinyue Chen, Deepak Pathak, Pengchuan Zhang, and Deva Ramanan. Revisiting the role of language priors in vision-language models. *arXiv preprint arXiv:2306.01879*, 2024a. 40

- Zhiqiu Lin, Deepak Pathak, Baiqi Li, Jiayao Li, Xide Xia, Graham Neubig, Pengchuan Zhang, and Deva Ramanan. Evaluating text-to-visual generation with image-to-text generation. *arXiv preprint arXiv:2404.01291*, 2024b. **33**
- Fangyu Liu, Guy Emerson, and Nigel Collier. Visual spatial reasoning. *Transactions of the Association for Computational Linguistics*, 11:635–651, 2023a. doi: 10.1162/tacl_a_00566. URL <https://aclanthology.org/2023.tacl-1.37>. **21**
- Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. Aligning large multi-modal model with robust instruction tuning. *arXiv preprint arXiv:2306.14565*, 2023b. **41**
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. *arXiv preprint arXiv: 2310.03744*, 2023c. **19, 25, 27**
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 34892–34916. Curran Associates, Inc., 2023d. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/6dcf277ea32ce3288914faf369fe6de0-Paper-Conference.pdf. **19, 26, 27, 28**
- Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next: Improved reasoning, ocr, and world knowledge, January 2024a. URL <https://llava-vl.github.io/blog/2024-01-30-llava-next/>. **19, 27**
- Shih-Yang Liu, Chien-Yi Wang, Hongxu Yin, Pavlo Molchanov, Yu-Chiang Frank Wang, Kwang-Ting Cheng, and Min-Hung Chen. DoRA: Weight-decomposed low-rank adaptation. *arXiv preprint arXiv:2402.09353*, 2024b. **29**
- Yuliang Liu, Zhang Li, Hongliang Li, Wenwen Yu, Mingxin Huang, Dezhi Peng, Mingyu Liu, Mingrui Chen, Chunyuan Li, Lianwen Jin, et al. On the hidden mystery of OCR in large multimodal models. *arXiv preprint arXiv:2305.07895*, 2023e. **28**
- Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper_files/paper/2019/file/c74d97b01eae257e44aa9d5bade97baf-Paper.pdf. **6**
- Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for science question answering. *Advances in Neural Information Processing Systems*, 35:2507–2521, 2022. **26, 34**
- Zixian Ma, Jerry Hong, Mustafa Omer Gul, Mona Gandhi, Irena Gao, and Ranjay Krishna. CREPE: Can vision-language foundation models reason compositionally? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10910–10921, 2023. **37, 40**

- Anas Mahmoud, Mostafa Elhoushi, Amro Abbas, Yu Yang, Newsha Ardalani, Hugh Leather, and Ari Morcos. Sieve: Multimodal dataset pruning using image captioning models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2024. 18
- Pratyush Maini, Sachin Goyal, Zachary C Lipton, J Zico Kolter, and Aditi Raghunathan. T-MARS: Improving visual representations by circumventing text feature learning. *arXiv preprint arXiv:2307.03132*, 2023. 18
- Subhransu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. Fine-grained visual classification of aircraft. *arXiv preprint arXiv:1306.5151*, 2013. 35
- Oscar Mañas, Pau Rodriguez Lopez, Saba Ahmadi, Aida Nematzadeh, Yash Goyal, and Aishwarya Agrawal. MAPL: Parameter-efficient adaptation of unimodal pre-trained models for vision-language few-shot prompting. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 2523–2548, 2023. 30, 31
- Oscar Mañas, Benno Krojer, and Aishwarya Agrawal. Improving automatic VQA evaluation using large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 4171–4179, 2024. 34
- Karttikeya Mangalam, Raiymbek Akshulakov, and Jitendra Malik. Egoschema: A diagnostic benchmark for very long-form video language understanding. *Advances in Neural Information Processing Systems*, 36, 2024. 45
- Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. OK-VQA: A visual question answering benchmark requiring external knowledge. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3195–3204, 2019. 21, 34
- Ahmed Masry, Xuan Long Do, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. ChartQA: A benchmark for question answering about charts with visual and logical reasoning. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio, editors, *Findings of the Association for Computational Linguistics: ACL 2022*, pages 2263–2279, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.findings-acl.177. URL <https://aclanthology.org/2022.findings-acl.177>. 35
- Minesh Mathew, Dimosthenis Karatzas, and C. V. Jawahar. DocVQA: A dataset for vqa on document images. In *2021 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 2199–2208, 2021. doi: 10.1109/WACV48630.2021.00225. 35
- Minesh Mathew, Viraj Bagal, Rubèn Tito, Dimosthenis Karatzas, Ernest Valveny, and CV Jawahar. InfographicVQA. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1697–1706, 2022. 35
- Nitika Mathur, Timothy Baldwin, and Trevor Cohn. Tangled up in BLEU: Reevaluating the evaluation of automatic machine translation evaluation metrics. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault, editors, *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4984–4997, Online, July

2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.448. URL <https://aclanthology.org/2020.acl-main.448>. 33
- Brandon McKinzie, Zhe Gan, Jean-Philippe Fauconnier, Sam Dodge, Bowen Zhang, Philipp Dufter, Dhruvi Shah, Xianzhi Du, Futang Peng, Floris Weers, et al. MM1: Methods, analysis & insights from multimodal llm pre-training. *arXiv preprint arXiv:2403.09611*, 2024. 20, 21, 23
- Sachit Menon and Carl Vondrick. Visual classification via description from large language models. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=j1AjNL8z5cs>. 36
- Jack Merullo, Louis Castricato, Carsten Eickhoff, and Ellie Pavlick. Linearly mapping from image to text space. In *The Eleventh International Conference on Learning Representations*, 2022. 30, 31
- Anand Mishra, KartEEK Alahari, and C. V. Jawahar. Top-down and bottom-up cues for scene text recognition. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pages 2687–2694, 2012. doi: 10.1109/CVPR.2012.6247990. 35
- Anand Mishra, Shashank Shekhar, Ajeet Kumar Singh, and Anirban Chakraborty. OCR-VQA: Visual question answering by reading text in images. In *2019 International Conference on Document Analysis and Recognition (ICDAR)*, pages 947–952. IEEE, 2019. 35
- Liliane Momeni, Mathilde Caron, Arsha Nagrani, Andrew Zisserman, and Cordelia Schmid. Verbs in action: Improving verb understanding in video-language models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15579–15591, 2023. 45
- Norman Mu, Alexander Kirillov, David Wagner, and Saining Xie. SLIP: Self-supervision meets language-image pre-training. In *European Conference on Computer Vision*, pages 529–544. Springer, 2022. 19
- Andrew Ng and Michael Jordan. On discriminative vs. generative classifiers: A comparison of logistic regression and naive bayes. *Advances in Neural Information Processing Systems*, 14, 2001. 13
- Thao Nguyen, Samir Yitzhak Gadre, Gabriel Ilharco, Sewoong Oh, and Ludwig Schmidt. Improving multimodal datasets with image captioning. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023. URL <https://openreview.net/forum?id=VIRKdeFJIg>. 19
- Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of classes. *2008 Sixth Indian Conference on Computer Vision, Graphics & Image Processing*, pages 722–729, 2008. URL <https://api.semanticscholar.org/CorpusID:15193013>. 35
- Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. 8

- Vicente Ordonez, Girish Kulkarni, and Tamara Berg. Im2text: Describing images using 1 million captioned photographs. In J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 24. Curran Associates, Inc., 2011. URL https://proceedings.neurips.cc/paper_files/paper/2011/file/5dd9db5e033da9c6fb5ba83c7a7e9-Paper.pdf. 15
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022. 42
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting on association for computational linguistics*, pages 311–318. Association for Computational Linguistics, 2002. 33
- Shubham Parashar, Zhiqiu Lin, Yanan Li, and Shu Kong. Prompting scientific names for zero-shot species recognition. *arXiv preprint arXiv:2310.09929*, 2023. 36
- Shubham Parashar, Zhiqiu Lin, Tian Liu, Xiangjue Dong, Yanan Li, Deva Ramanan, James Caverlee, and Shu Kong. The neglected tails of vision-language models. *arXiv preprint arXiv:2401.12425*, 2024. 19, 36
- Omkar M. Parkhi, Andrea Vedaldi, Andrew Zisserman, and C. V. Jawahar. Cats and dogs. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2012. 35
- Deepak Pathak, Philipp Krähenbühl, Jeff Donahue, Trevor Darrell, and Alexei A. Efros. Context encoders: Feature learning by inpainting. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2536–2544, 2016. doi: 10.1109/CVPR.2016.278. 9
- Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, Qixiang Ye, and Furu Wei. Grounding multimodal large language models to the world. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=1LmqxkfSIw>. 25
- Ethan Perez, Saffron Huang, Francis Song, Trevor Cai, Roman Ring, John Aslanides, Amelia Glaese, Nat McAleese, and Geoffrey Irving. Red teaming language models with language models. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang, editors, *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 3419–3448, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.225. URL <https://aclanthology.org/2022.emnlp-main.225>. 42
- Mihir Prabhudesai, Tsung-Wei Ke, Alexander C Li, Deepak Pathak, and Katerina Fragkiadaki. Test-time adaptation of discriminative models via diffusion generative feedback. *arXiv preprint arXiv:2311.16102*, 2023. 14
- Sarah Pratt, Ian Covert, Rosanne Liu, and Ali Farhadi. What does a platypus look like? generating customized prompts for zero-shot image classification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15691–15701, 2023. 36

- Filip Radenovic, Abhimanyu Dubey, Abhishek Kadian, Todor Mihaylov, Simon Vandenhende, Yash Patel, Yi Wen, Vignesh Ramanathan, and Dhruv Mahajan. Filtering, distillation, and hard negatives for vision-language pre-training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6967–6977, 2023a. 18
- Filip Radenovic, Abhimanyu Dubey, Abhishek Kadian, Todor Mihaylov, Simon Vandenhende, Yash Patel, Yi Wen, Vignesh Ramanathan, and Dhruv Mahajan. Filtering, distillation, and hard negatives for vision-language pre-training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6967–6977, 2023b. 26
- 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 Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 8748–8763. PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/radford21a.html>. 9, 12, 16, 18, 19, 21, 22, 31, 35, 36
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140): 1–67, 2020. 15
- Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024. 27
- Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns, Trevor Darrell, and Kate Saenko. Object hallucination in image captioning. In Ellen Riloff, David Chiang, Julia Hockenmaier, and Jun'ichi Tsujii, editors, *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4035–4045, Brussels, Belgium, October–November 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1437. URL <https://aclanthology.org/D18-1437>. 40
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2021. 15
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022. 11, 12, 19
- Candace Ross, Boris Katz, and Andrei Barbu. Measuring social biases in grounded vision and language embeddings. *arXiv preprint arXiv:2002.08911*, 2020. 39

- Y Dan Rubinstein, Trevor Hastie, et al. Discriminative vs informative learning. In *KDD*, volume 5, pages 49–53, 1997. 13
- Paul Röttger, Fabio Pernisi, Bertie Vidgen, and Dirk Hovy. Safetyprompts: a systematic review of open datasets for evaluating and improving large language model safety, 2024. 42
- Noveen Sachdeva, Benjamin Coleman, Wang-Cheng Kang, Jianmo Ni, Lichan Hong, Ed H Chi, James Caverlee, Julian McAuley, and Derek Zhiyuan Cheng. How to train data-efficient LLMs. *arXiv preprint arXiv:2402.09668*, 2024. 21
- Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems*, 35:36479–36494, 2022. 11, 12
- Jameel Hassan Abdul Samadh, Hanan Gani, Noor Hazim Hussein, Muhammad Uzair Khattak, Muzammal Naseer, Fahad Khan, and Salman Khan. Align your prompts: Test-time prompting with distribution alignment for zero-shot generalization. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=CusNOTRkQw>. 36
- Shibani Santurkar, Yann Dubois, Rohan Taori, Percy Liang, and Tatsunori Hashimoto. Is a caption worth a thousand images? A controlled study for representation learning. *arXiv preprint arXiv:2207.07635*, 2022. 19
- Morgan Klaus Scheuerman, Katy Weathington, Tarun Mugunthan, Emily Denton, and Casey Fiesler. From human to data to dataset: Mapping the traceability of human subjects in computer vision datasets. *Proceedings of the ACM on Human-Computer Interaction*, 7(CSCW1):1–33, 2023. 39
- Christoph Schuhmann. Laion-aesthetics. <https://laion.ai/blog/laion-aesthetics/>, 2023. 21
- Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of CLIP-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114*, 2021. 15, 18
- Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. Laion-5b: An open large-scale dataset for training next generation image-text models, 2022. 28
- Dustin Schwenk, Apoorv Khandelwal, Christopher Clark, Kenneth Marino, and Roozbeh Mottaghi. A-okvqa: A benchmark for visual question answering using world knowledge. In *European Conference on Computer Vision*, pages 146–162. Springer, 2022. 21

- Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2556–2565, Melbourne, Australia, July 2018a. Association for Computational Linguistics. doi: 10.18653/v1/P18-1238. URL <https://aclanthology.org/P18-1238>. 18
- Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2556–2565, 2018b. 15
- Vasu Sharma, Karthik Padthe, Newsha Ardalani, Kushal Tirumala, Russell Howes, Hu Xu, Po-Yao Huang, Shang-Wen Li, Armen Aghajanyan, and Gargi Ghosh. Text quality-based pruning for efficient training of language models. *arXiv preprint arXiv:2405.01582*, 2024. 21
- Ashish Shenoy, Yichao Lu, Srihari Jayakumar, Debojeet Chatterjee, Mohsen Moslehpour, Pierce Chuang, Abhay Harpale, Vikas Bhardwaj, Di Xu, Shicong Zhao, et al. Lumos: Empowering multimodal llms with scene text recognition. *arXiv preprint arXiv:2402.08017*, 2024. 29
- Cunzhao Shi, Chunheng Wang, Baihua Xiao, Song Gao, and Jinlong Hu. End-to-end scene text recognition using tree-structured models. *Pattern Recognition*, 47:2853–2866, 2014. 35
- Ravid Shwartz Ziv and Yann LeCun. To compress or not to compress—self-supervised learning and information theory: A review. *Entropy*, 26(3):252, 2024. 10
- Chenglei Si, Weijia Shi, Chen Zhao, Luke Zettlemoyer, and Jordan Lee Boyd-Graber. Getting MoRE out of Mixture of language model Reasoning Experts. *Findings of Empirical Methods in Natural Language Processing*, 2023. 34
- Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. Towards VQA models that can read. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8317–8326, 2019. 34, 35
- Amanpreet Singh, Ronghang Hu, Vedanuj Goswami, Guillaume Couairon, Wojciech Galuba, Marcus Rohrbach, and Douwe Kiela. Flava: A foundational language and vision alignment model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15638–15650, 2022. 9
- Brandon Smith, Miguel Farinha, Siobhan Mackenzie Hall, Hannah Rose Kirk, Aleksandar Shtedritski, and Max Bain. Balancing the picture: Debiasing vision-language datasets with synthetic contrast sets. *arXiv preprint arXiv:2305.15407*, 2023. 39
- Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Diffusion art or digital forgery? Investigating data replication in diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6048–6058, 2023. 41

- Ben Sorscher, Robert Geirhos, Shashank Shekhar, Surya Ganguli, and Ari S. Morcos. Beyond neural scaling laws: beating power law scaling via data pruning. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=UmvSlP-PyV>. 16
- Tejas Srinivasan and Yonatan Bisk. Worst of both worlds: Biases compound in pre-trained vision-and-language models. *arXiv preprint arXiv:2104.08666*, 2021. 39
- Yixuan Su, Tian Lan, Huayang Li, Jialu Xu, Yan Wang, and Deng Cai. Pandagpt: One model to instruction-follow them all. *arXiv preprint arXiv:2305.16355*, 2023. 45
- Chen Sun, Austin Myers, Carl Vondrick, Kevin Murphy, and Cordelia Schmid. Videobert: A joint model for video and language representation learning. In *ICCV*, 2019. 43
- Lichao Sun, Yue Huang, Haoran Wang, Siyuan Wu, Qihui Zhang, Chujie Gao, Yixin Huang, Wenhan Lyu, Yixuan Zhang, Xiner Li, et al. TrustLLM: Trustworthiness in large language models. *arXiv preprint arXiv:2401.05561*, 2024. 42
- Zhiqing Sun, Sheng Shen, Shengcao Cao, Haotian Liu, Chunyuan Li, Yikang Shen, Chuang Gan, Liang-Yan Gui, Yu-Xiong Wang, Yiming Yang, et al. Aligning large multimodal models with factually augmented rlhf. *arXiv preprint arXiv:2309.14525*, 2023. 26, 27, 41
- Yi-Lin Sung, Jaemin Cho, and Mohit Bansal. VI-adapter: Parameter-efficient transfer learning for vision-and-language tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5227–5237, 2022. 30
- Makarand Tapaswi, Yukun Zhu, Rainer Stiefelwagen, Antonio Torralba, Raquel Urtasun, and Sanja Fidler. MovieQA: Understanding stories in movies through question-answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4631–4640, 2016. 42
- Chameleon Team. Chameleon: Mixed-modal early-fusion foundation models. *arXiv preprint arXiv:2405.09818*, 2024. 11, 12
- Tristan Thrush, Ryan Jiang, Max Bartolo, Amanpreet Singh, Adina Williams, Douwe Kiela, and Candace Ross. Winoground: Probing vision and language models for visio-linguistic compositionality. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5238–5248, 2022. 14, 31, 36, 40
- Yonglong Tian, Lijie Fan, Kaifeng Chen, Dina Katabi, Dilip Krishnan, and Phillip Isola. Learning vision from models rivals learning vision from data. *arXiv preprint arXiv:2312.17742*, 2023a. 19
- Yonglong Tian, Lijie Fan, Phillip Isola, Huiwen Chang, and Dilip Krishnan. Stablerep: Synthetic images from text-to-image models make strong visual representation learners. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023b. URL <https://openreview.net/forum?id=xpjsOQtKqx>. 19

- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models, 2023. 6
- Maria Tsimpoukelli, Jacob L Menick, Serkan Cabi, SM Eslami, Oriol Vinyals, and Felix Hill. Multimodal few-shot learning with frozen language models. *Advances in Neural Information Processing Systems*, 34:200–212, 2021. 15
- Vishaal Udandarao, Ameya Prabhu, Adhiraj Ghosh, Yash Sharma, Philip H. S. Torr, Adel Bibi, Samuel Albanie, and Matthias Bethge. No “zero-shot” without exponential data: Pretraining concept frequency determines multimodal model performance. *arXiv preprint arXiv:2404.04125*, 2024. 19, 36, 40
- Shagun Uppal, Sarthak Bhagat, Devamanyu Hazarika, Navonil Majumder, Soujanya Poria, Roger Zimmermann, and Amir Zadeh. Multimodal research in vision and language: A review of current and emerging trends. *Information Fusion*, 77:149–171, 2022. 4
- Jack Urbanek, Florian Bordes, Pietro Astolfi, Mary Williamson, Vasu Sharma, and Adriana Romero-Soriano. A picture is worth more than 77 text tokens: Evaluating clip-style models on dense captions. *arXiv preprint arXiv:2312.08578*, 2023. 21, 32, 37
- Théophane Vallaey, Mustafa Shukor, Matthieu Cord, and Jakob Verbeek. Improved baselines for data-efficient perceptual augmentation of LLMs. *arXiv preprint arXiv:2403.13499*, 2024. 31
- Aaron Van Den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learning. *Advances in Neural Information Processing Systems*, 30, 2017. 13
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf. 5, 6, 9, 10
- Andreas Veit, Tomas Matera, Lukas Neumann, Jiri Matas, and Serge Belongie. COCO-Text: Dataset and benchmark for text detection and recognition in natural images. *arXiv preprint arXiv:1601.07140*, 2016. 35
- Bertie Vidgen, Hannah Rose Kirk, Rebecca Qian, Nino Scherrer, Anand Kannappan, Scott A Hale, and Paul Röttger. Simple safety tests: a test suite for identifying critical safety risks in large language models. *arXiv preprint arXiv:2311.08370*, 2023. 42
- Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural Computation*, 23(7):1661–1674, 2011. 8
- Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. Extracting and composing robust features with denoising autoencoders. In *Proceedings of the 25th International Conference on Machine Learning, ICML '08*, page 1096–1103, New York,

- NY, USA, 2008. Association for Computing Machinery. ISBN 9781605582054. doi: 10.1145/1390156.1390294. URL <https://doi.org/10.1145/1390156.1390294>. 9
- C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. Caltech-ucsd birds-200-2011. Technical Report CNS-TR-2011-001, California Institute of Technology, 2011. 35
- Peng Wang, An Yang, Rui Men, Junyang Lin, Shuai Bai, Zhikang Li, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. OFA: Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 23318–23340. PMLR, 17–23 Jul 2022. URL <https://proceedings.mlr.press/v162/wang22a1.html>. 44
- Tan Wang, Kevin Lin, Linjie Li, Chung-Ching Lin, Zhengyuan Yang, Hanwang Zhang, Zicheng Liu, and Lijuan Wang. Equivariant similarity for vision-language foundation models. *arXiv preprint arXiv:2303.14465*, 2023a. 40
- Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, Jiazheng Xu, Bin Xu, Juanzi Li, Yuxiao Dong, Ming Ding, and Jie Tang. Cogvlm: Visual expert for pretrained language models, 2023b. 27
- Xinyu Wang, Yuliang Liu, Chunhua Shen, Chun Chet Ng, Canjie Luo, Lianwen Jin, Chee Seng Chan, Anton van den Hengel, and Liangwei Wang. On the general value of evidence, and bilingual scene-text visual question answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10126–10135, 2020. 35
- Laura Weidinger, Jonathan Uesato, Maribeth Rauh, Conor Griffin, Po-Sen Huang, John Mellor, Amelia Glaese, Myra Cheng, Borja Balle, Atoosa Kasirzadeh, et al. Taxonomy of risks posed by language models. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, pages 214–229, 2022. 42
- Alexander Wettig, Aatmik Gupta, Saumya Malik, and Danqi Chen. Qurating: Selecting high-quality data for training language models. *arXiv preprint arXiv:2402.09739*, 2024. 21
- Spencer Whitehead, Suzanne Petryk, Vedaad Shakib, Joseph Gonzalez, Trevor Darrell, Anna Rohrbach, and Marcus Rohrbach. Reliable visual question answering: Abstain rather than answer incorrectly. In *European Conference on Computer Vision*, pages 148–166. Springer, 2022. 34
- Olivia Wiles, Isabela Albuquerque, and Sven Gowal. Discovering bugs in vision models using off-the-shelf image generation and captioning. *arXiv preprint arXiv:2208.08831*, 2022. 39
- Robert Wolfe and Aylin Caliskan. American == white in multimodal language-and-image ai. In *Proceedings of the 2022 AAAI/ACM Conference on AI, Ethics, and Society*, pages 800–812, 2022. 40

- Robert Wolfe, Yiwei Yang, Bill Howe, and Aylin Caliskan. Contrastive language-vision ai models pretrained on web-scraped multimodal data exhibit sexual objectification bias. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*, pages 1174–1185, 2023. 40
- Zhirong Wu, Yuanjun Xiong, Stella X Yu, and Dahua Lin. Unsupervised feature learning via non-parametric instance discrimination. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3733–3742, 2018. 8
- Dejing Xu, Zhou Zhao, Jun Xiao, Fei Wu, Hanwang Zhang, Xiangnan He, and Yueting Zhuang. Video question answering via gradually refined attention over appearance and motion. In *Proceedings of the 25th ACM International Conference on Multimedia*, pages 1645–1653, 2017. 45
- Hu Xu, Saining Xie, Xiaoqing Ellen Tan, Po-Yao Huang, Russell Howes, Vasu Sharma, Shang-Wen Li, Gargi Ghosh, Luke Zettlemoyer, and Christoph Feichtenhofer. Demystifying clip data. In *International Conference on Learning Representations*, 2024. URL <https://openreview.net/pdf?id=5BCFlnfE1g>. 18, 21, 23
- Jun Xu, Tao Mei, Ting Yao, and Yong Rui. MSR-VTT: A large video description dataset for bridging video and language. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5288–5296, 2016. 45
- Ran Xu, Caiming Xiong, Wei Chen, and Jason Corso. Jointly modeling deep video and compositional text to bridge vision and language in a unified framework. In *Proceedings of the AAAI conference on artificial intelligence*, volume 29, 2015. 42
- Jingkang Yang, Wenxuan Peng, Xiangtai Li, Zujin Guo, Liangyu Chen, Bo Li, Zheng Ma, Kaiyang Zhou, Wayne Zhang, Chen Change Loy, et al. Panoptic video scene graph generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18675–18685, 2023. 27
- Jiahui Yu, Xin Li, Jing Yu Koh, Han Zhang, Ruoming Pang, James Qin, Alexander Ku, Yuanzhong Xu, Jason Baldridge, and Yonghui Wu. Vector-quantized image modeling with improved VQGAN. In *International Conference on Learning Representations*, 2022a. URL <https://openreview.net/forum?id=pfNyExj7z2>. 14
- Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu. Coca: Contrastive captioners are image-text foundation models. *Transactions on Machine Learning Research*, 2022b. ISSN 2835-8856. URL <https://openreview.net/forum?id=Ee277P3AYC>. 11
- Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive models for content-rich text-to-image generation. *arXiv preprint arXiv:2206.10789*, 2(3): 5, 2022c. 11, 12
- Lili Yu, Bowen Shi, Ramakanth Pasunuru, Benjamin Muller, Olga Golovneva, Tianlu Wang, Arun Babu, Binh Tang, Brian Karrer, Shelly Sheynin, et al. Scaling autoregressive multi-modal models: Pretraining and instruction tuning. *arXiv preprint arXiv:2309.02591*, 2023. 11

- Zhou Yu, Dejing Xu, Jun Yu, Ting Yu, Zhou Zhao, Yueting Zhuang, and Dacheng Tao. Activitynet-qa: A dataset for understanding complex web videos via question answering. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 9127–9134, 2019. 45
- Ye Yuan, Xiao Liu, Wondimu Dikubab, Hui Liu, Zhilong Ji, Zhongqin Wu, and Xiang Bai. Syntax-aware network for handwritten mathematical expression recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4553–4562, June 2022. 35
- Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. MMMU: A massive multi-discipline multimodal understanding and reasoning benchmark for expert AGI. *arXiv preprint arXiv:2311.16502*, 2023. 34
- Alan Yuille and Daniel Kersten. Vision as bayesian inference: analysis by synthesis? *Trends in Cognitive Sciences*, 10(7):301–308, 2006. 13
- Mert Yuksekgonul, Federico Bianchi, Pratyusha Kalluri, Dan Jurafsky, and James Zou. When and why vision-language models behave like bags-of-words, and what to do about it? In *International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=KRLUvvh8uaX>. 26, 32, 36, 37, 40
- Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stéphane Deny. Barlow twins: Self-supervised learning via redundancy reduction. In *International Conference on Machine Learning*, pages 12310–12320. PMLR, 2021. 20
- Rowan Zellers, Ximing Lu, Jack Hessel, Youngjae Yu, Jae Sung Park, Jize Cao, Ali Farhadi, and Yejin Choi. Merlot: Multimodal neural script knowledge models. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 23634–23651. Curran Associates, Inc., 2021. URL https://proceedings.neurips.cc/paper_files/paper/2021/file/c6d4eb15f1e84a36eff58eca3627c82e-Paper.pdf. 43
- Yan Zeng, Xinsong Zhang, and Hang Li. Multi-grained vision language pre-training: Aligning texts with visual concepts. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato, editors, *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pages 25994–26009. PMLR, 2022. URL <https://proceedings.mlr.press/v162/zeng22c.html>. 25
- Bohan Zhai, Shijia Yang, Xiangchen Zhao, Chenfeng Xu, Sheng Shen, Dongdi Zhao, Kurt Keutzer, Manling Li, Tan Yan, and Xiangjun Fan. Halle-switch: Rethinking and controlling object existence hallucinations in large vision language models for detailed caption. *arXiv preprint arXiv:2310.01779*, 2023a. 41
- Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *2023 IEEE/CVF International Conference on Computer Vision*

- (*ICCV*), pages 11941–11952, Los Alamitos, CA, USA, oct 2023b. IEEE Computer Society. doi: 10.1109/ICCV51070.2023.01100. URL <https://doi.ieeecomputersociety.org/10.1109/ICCV51070.2023.01100>. 9
- Chenshuang Zhang, Chaoning Zhang, Mengchun Zhang, and In So Kweon. Text-to-image diffusion model in generative AI: A survey. *arXiv preprint arXiv:2303.07909*, 2023a. 12
- Hang Zhang, Xin Li, and Lidong Bing. Video-LLaMA: An instruction-tuned audio-visual language model for video understanding. In Yansong Feng and Els Lefever, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 543–553, Singapore, December 2023b. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-demo.49. URL <https://aclanthology.org/2023.emnlp-demo.49>. 44, 45
- Jingyi Zhang, Jiaxing Huang, Sheng Jin, and Shijian Lu. Vision-language models for vision tasks: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024a. 4
- Kai Zhang, Lingbo Mo, Wenhui Chen, Huan Sun, and Yu Su. MagicBrush: A manually annotated dataset for instruction-guided image editing. *Advances in Neural Information Processing Systems*, 36, 2024b. 21
- Muhan Zhang, Shali Jiang, Zhicheng Cui, Roman Garnett, and Yixin Chen. D-VAE: A variational autoencoder for directed acyclic graphs. In *Advances in Neural Information Processing Systems*, pages 1586–1598, 2019. 10
- Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. OPT: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*, 2022. 11, 12
- Yanzhe Zhang, Ruiyi Zhang, Jiuxiang Gu, Yufan Zhou, Nedim Lipka, Diyi Yang, and Tong Sun. LLaVAR: Enhanced visual instruction tuning for text-rich image understanding. *arXiv preprint arXiv:2306.17107*, 2023c. 28
- Youcai Zhang, Xinyu Huang, Jinyu Ma, Zhaoyang Li, Zhaochuan Luo, Yanchun Xie, Yuzhuo Qin, Tong Luo, Yaqian Li, Shilong Liu, et al. Recognize anything: A strong image tagging model. *arXiv preprint arXiv:2306.03514*, 2023d. 40
- Long Zhao, Nitesh B Gundavarapu, Liangzhe Yuan, Hao Zhou, Shen Yan, Jennifer J Sun, Luke Friedman, Rui Qian, Tobias Weyand, Yue Zhao, et al. Videoprism: A foundational visual encoder for video understanding. *arXiv preprint arXiv:2402.13217*, 2024. 45
- Tiancheng Zhao, Tianqi Zhang, Mingwei Zhu, Haozhan Shen, Kyusong Lee, Xiaopeng Lu, and Jianwei Yin. V1-checklist: Evaluating pre-trained vision-language models with objects, attributes and relations. *arXiv preprint arXiv:2207.00221*, 2022. 37, 40
- Kaizhi Zheng, Xuehai He, and Xin Eric Wang. MiniGPT-5: Interleaved vision-and-language generation via generative tokens. *arXiv preprint arXiv:2310.02239*, 2023. 15

- Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision-language models. *International Journal of Computer Vision*, 130(9):2337–2348, 2022. 30
- Yutong Zhou and Nobutaka Shimada. Vision + language applications: A survey. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, pages 826–842, June 2023. 4
- Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. MiniGPT-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023a. 15
- Wanrong Zhu, Jack Hessel, Anas Awadalla, Samir Yitzhak Gadre, Jesse Dodge, Alex Fang, Youngjae Yu, Ludwig Schmidt, William Yang Wang, and Yejin Choi. Multimodal c4: An open, billion-scale corpus of images interleaved with text. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 8958–8974. Curran Associates, Inc., 2023b. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/1c6bed78d3813886d3d72595dbecb80b-Paper-Datasets_and_Benchmarks.pdf. 20