

Equivariant Similarity for Vision-Language Foundation Models

Tan Wang¹, Kevin Lin², Linjie Li², Chung-Ching Lin², Zhengyuan Yang²,
Hanwang Zhang¹, Zicheng Liu², Lijuan Wang²

¹Nanyang Technological University ²Microsoft

{TAN317, hanwangzhang}@ntu.edu.sg, {keli, lindsey.li, chungching.lin, zhengyang, zliu, lijuanw}@microsoft.com

Abstract

This study explores the concept of equivariance in vision-language foundation models (VLMs), focusing specifically on the multimodal similarity function that is not only the major training objective but also the core delivery to support downstream tasks. Unlike the existing image-text similarity objective which only categorizes matched pairs as similar and unmatched pairs as dissimilar, equivariance also requires similarity to vary faithfully according to the semantic changes. This allows VLMs to generalize better to nuanced and unseen multimodal compositions. However, modeling equivariance is challenging as the ground truth of semantic change is difficult to collect. For example, given an image-text pair about a dog, it is unclear to what extent the similarity changes when the pixel is changed from dog to cat? To this end, we propose EQSIM, a regularization loss that can be efficiently calculated from any two matched training pairs and easily pluggable into existing image-text retrieval fine-tuning. Meanwhile, to further diagnose the equivariance of VLMs, we present a new challenging benchmark EQBEN. Compared to the existing evaluation sets, EQBEN is the first to focus on “visual-minimal change”. Extensive experiments show the lack of equivariance in current VLMs¹ and validate the effectiveness of EQSIM².

1. Introduction

Vision-language (VL) training is all about learning “good” features for each modality, such that the features should faithfully represent the underlying semantics. Thanks to the large-scale image-text pairs on the Web, we have abundant multimodal supervision for the two features with the same semantic meaning [60, 35, 48, 27]—each matched image-text pair should have “similar” visual and textual features, and each unmatched pair should have “dissimilar” ones. Thus, the image-text similarity plays a crucial role to define the feature quality in training VL founda-

¹We also include results of Multimodal LLM in Appendix A.4.

²Code is available at <https://github.com/Wangt-CN/EqBen>

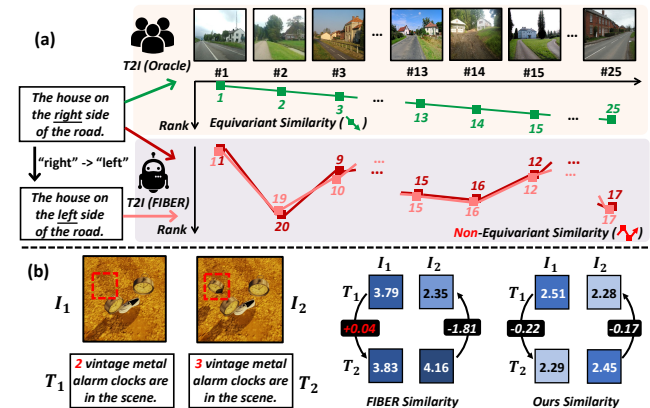


Figure 1: (a) Comparison between oracle and the latest SoTA VLM FIBER [13] similarity measure by ranking the candidate images with the given query texts. Check Appendix for full ranking results. (b) Measuring the similarity score change (number in \blacksquare) of FIBER [13] and our proposed EQSIM by applying a slight text change (“right” \leftrightarrow “left”). Darker color indicates larger similarity.

tion models (VLMs) [68, 74, 35, 48, 27, 14, 13, 34].

Has the prevailing “matched vs. unmatched” similarity fulfilled its duty? Yes and no. On the one hand, recent VLMs [51, 68, 13, 48, 74, 49] have demonstrated impressive results in various downstream VL tasks such as image-text retrieval. However, on the other hand, it is acknowledged by the community that the VLMs still fall short in *nuanced and complex semantic compositions* [49, 9, 44, 62]. In this regard, we present a text-to-image retrieval example on LAION400M [54] with the most recent SOTA VLM FIBER [13]. As shown in Figure 1(a), given the query text “the house on the *right* side of the road”, we first invite 5 graduate students to rank 25 candidate images from most similar to least similar. The continuously decreasing ranking from human judges (\checkmark) is served as the oracle semantic similarity measure. We then compared this ranking with the ones from FIBER [13] (\heartsuit). Although FIBER correctly retrieved the top-1 image (image#1, ranks 1), some semantically incorrect images (e.g., image#25, ranks 17) are falsely ranked higher than the correct ones (e.g., image#2, ranks 20). Furthermore, when modifying the query text with a slight semantic change (“right” \rightarrow “left”), the rankings re-

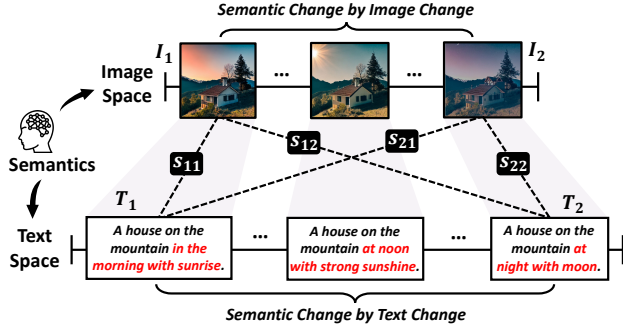


Figure 2: The illustration of the core idea in EQSIM. Besides the two matched pairs $\{I_1, T_1\}$ and $\{I_2, T_2\}$, we don't need extra annotation such as the middle pair.

main almost the same. Clearly, the similarity changes in FIBER do not faithfully reflect the semantic changes in images (#1 \rightarrow #25) or text queries (“right” \rightarrow “left”).

To quantitatively measure the above inconsistency between semantic and similarity score changes, we consider two matched image-text pairs $\{I_1, T_1\}$ and $\{I_2, T_2\}$ that are semantically similar but only different in the number of clocks in Figure 1 (b). With a slight change of clock counts in caption (“2” \rightarrow “3”), FIBER mistakenly assigns a higher similarity score to $\{I_1, T_2\}$ rather than $\{I_1, T_1\}$ (3.83 v.s. 3.79). Furthermore, the changes in similarity scores guided by the semantic change (“2” \leftrightarrow “3”) are highly inconsistent (+0.04 v.s. -1.81). Ideally, an **equivariant** image-text similarity measure should faithfully reflect the semantic change, *i.e.*, the same semantic changes should lead to a similar amount of similarity changes (*e.g.*, -0.22 v.s. -0.17 of ours in Figure 1(b)).

Equivariance Loss. To address this non-equivariance issue, we propose Equivariant Similarity Learning (EQSIM), which imposes additional equivariance regularization on image-text pairs for VLM learning without additional supervision. Figure 2 illustrates the underlying semantics perceived by human, where each matched pair demonstrates the image and text corresponding to the underlying semantic. Given two matched image-text pairs $\{I_1, T_1\}$ as semantic 1 and $\{I_2, T_2\}$ as semantic 2, we can obtain four similarity scores s_{11}, s_{12}, s_{22} , and s_{21} . We define **Equivariant Similarity** to be an image-text similarity function, whose output value should correspond to the underlying semantic change, which can be measured by text or image change.

Definition 1 (Equivariant Similarity) The similarity s between image and text is equivariant if and only if the following equations hold:

$$s_{11} - s_{12} = \underbrace{\sum_{T_1}^{T_2} \mu(T)}_{\text{Semantic Change Measured by Text Change}}, \quad s_{22} - s_{21} = \underbrace{\sum_{T_2}^{T_1} \mu(T)}_{\text{Semantic Change Measured by Text Change}}, \quad (1)$$

$$s_{11} - s_{21} = \underbrace{\sum_{I_1}^{I_2} \mu(I)}_{\text{Semantic Change Measured by Image Change}}, \quad s_{22} - s_{12} = \underbrace{\sum_{I_2}^{I_1} \mu(I)}_{\text{Semantic Change Measured by Image Change}}, \quad (2)$$

where $\mu(I)$ ($\mu(T)$) denotes the measure [52] in image (text) space, *i.e.*, an infinitesimal unit of visual (textual) change. Based on Definition 1, we formally derive EQSIM, an equivariance loss for a hybrid learning strategy on both semantically close and distant training pairs (Section 3). Specifically, EQSIM directly enforces $s_{11} - s_{12} = s_{22} - s_{21}$ and $s_{11} - s_{21} = s_{22} - s_{12}$ for semantically close samples; while for semantically distant samples, we derive a simplified formulation of $s_{12} = s_{21}$. We show that adding EQSIM as a regularization term improves existing similarity training objectives significantly on challenging datasets (*e.g.*, over 4% on Winoground [62]) and tricky tasks (*e.g.*, around 30% on VALSE [44]). EQSIM can also retain or even improve retrieval performance on Flickr30K [46] dataset.

Equivariance Benchmark. To further facilitate the proper evaluation of equivariance in VL community, we present a novel evaluation benchmark dubbed EQBEN (Section 4). Motivated by the examples in Figure 1(b), EQBEN features “slightly” mis-matched pairs with a *minimal semantic drift* from the matched pairs, as opposed to “very different” matched and unmatched pairs that are easily distinguishable by both non-equivariant and equivariant similarities. Unlike recent efforts [44, 62] focusing on minimal semantic changes in captions, EQBEN pivots on diverse *visual-minimal* changes, automatically curated from time-varying visual contents in natural videos and synthetic engines with more precise control. We benchmark a full spectrum of VLMs on EQBEN, and reveal that the non-equivariant similarity in existing VLMs fails easily. On this new test bed, EQSIM can serve as a remedy and bring a large performance gain of $\sim 3\%$ on average.

Our contributions are summarized as follows: (1) We comprehensively study the problem of similarity equivariance in VLMs. We propose EQSIM for equivariant training and EQBEN for diagnostic evaluation; (2) EQSIM is not only theoretically grounded but also simple, effective and easily pluggable; and (3) EQBEN clearly diagnoses that conventional evaluation is not responsive to equivariance. Furthermore, EQSIM can significantly improve VLMs on EQBEN, as well as other challenging benchmarks.

2. Related Work

Pre-training VL Models. Early object detector (OD)-based methods [8, 78, 38, 37, 41, 15, 60] utilized the offline image region features from a pre-trained object detector [50]. More recent methods mainly learn from image pixels directly in an end-to-end manner [71, 27, 69, 58, 34, 65, 74, 1]. Researchers [16] further categorize VLMs into (1) Dual-Encoder (*e.g.*, CLIP [48] and ALIGN [35]) and (2) Fusion-Encoder (*e.g.*, METER [14], FIBER [13], and ALBEF [35]). It is worth noting that our proposed EQSIM is model-agnostic, and can be easily plugged into the image-text alignment objectives such as Image-Text

Matching (ITM) and Image-Text Contrastive (ITC) loss.

Diagnosing VL Models. Years of VL research have spawned a series of VL evaluation kits, from classical VL tasks [77, 64, 76, 46] (e.g., VQA [2] and image captioning [7]), to more complex contexts, such as adversarial examples [36, 6], robustness [21, 5, 66, 61, 28] and counterfactual reasoning [56, 23, 25, 42]. However, these benchmarks require manual annotation and their evaluation relies on task-specific model fine-tuning. Another line of work [62, 44, 79] probe VLMs on similarity measure with minimal *caption* semantic changes while keeping images intact. While our EQBEN tries to test whether the inherent image-text similarity measure in existing VLMs is sensitive to visual semantic changes. The most relevant work is ImageCoDe [32] which leverages video frames toward fine-grained image-text retrieval. However, ImageCoDe requires additional human crowdsourcing and is limited to real-world video sources. In contrast, EQBEN explores both natural and synthetic ways to generate image pairs with minimal semantic change, making the data generation process inclusive, automatic, and extensive.

Equivariance Learning. Unlike the wide usage of invariance in deep neural networks (e.g., shift invariance achieved by convolutional layers), strict group equivariance [11, 10, 72, 4] is hard to apply in practice. However, the equivariance property still plays an important role in various fields, such as self-supervised learning [12, 73, 67, 45, 22], representation learning [47], and language understanding [19]. In this paper, we point out the significance of the equivariant similarity measure in VLMs. Based on this, we further propose a novel loss EQSIM for the regularization of equivariance, as well as a new challenging benchmark EQBEN to diagnose the equivariance of existing VLMs. We notice that the recent CyCLIP [17] delivers a similar idea but with different motivation, implementation and evaluation settings. In Table 6, we compare with CyCLIP-equivalent baseline as EQSIM_{v1}. Check more detailed comparison in Appendix.

3. Improving VLMs with EQSIM

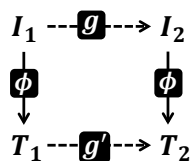


Figure 3: Commutativity in equivariant map.

Recall that VLMs adopt the image-text similarity as the training proxy for learning multimodal feature representation [48]. Therefore, the ultimate goal of pursuing equivariant similarity is to learn an equivariant feature map between image space and text space.

Definition 2 Let \mathcal{I} and \mathcal{T} be two continuous feature spaces. Let \mathcal{G} be a group whose group action on \mathcal{I} is defined by $g : \mathcal{I} \rightarrow \mathcal{I}$, and that on \mathcal{T} is defined by $g' : \mathcal{T} \rightarrow \mathcal{T}$. Then, $\phi : \mathcal{I} \rightarrow \mathcal{T}$ is an equivariant feature map if and only if $g' \cdot \phi(I) = \phi(g \cdot I)$ for all the group actions and $I \in \mathcal{I}$. The

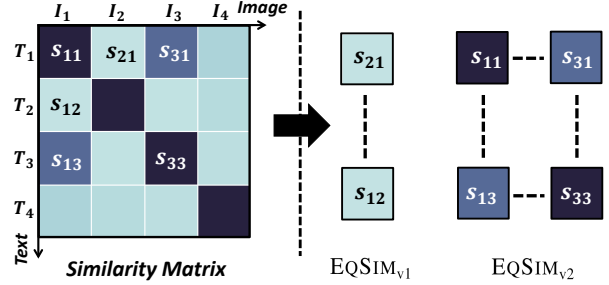


Figure 4: Illustration of EQSIM given a similarity matrix during training. Darker color indicates higher similarity.

commutativity for ϕ , g and g' is shown in Figure 3.

In Definition 1, the measure μ can be considered as the semantic group acting on an infinitesimal region in image or text space. Thus, by applying the commutativity of Definition 2 in Definition 1 to change the sum from image space to text space. Without loss of generality, we only show the results of sum $1 \rightarrow 2$:

$$\sum_{T_1}^{T_2} \mu(T) = \sum_{\phi(I_1)}^{\phi(I_2)} \mu(\phi(I)) = \sum_{\phi(I_1)}^{\phi(I_2)} \phi(\mu(I)). \quad (3)$$

This implies that the equivariant map establishes an isometry for the measure $\mu(I)$ in image space and $\phi(\mu(I))$ in text space. Thus, they only differ by a constant scale $C > 0$, i.e., $\phi(\mu(I)) = C\mu(I)$:

$$\sum_{\phi(I_1)}^{\phi(I_2)} \phi(\mu(I)) = C \sum_{I_1}^{I_2} \mu(I). \quad (4)$$

By combining Eq. (1), (2), (3), and (4), we have the following ratio equality as our EQSIM constraint:

$$\frac{s_{11} - s_{12}}{s_{11} - s_{21}} = \frac{s_{22} - s_{21}}{s_{22} - s_{12}} = C = 1. \quad (5)$$

Note that $C = 1$ can be derived by using the fact that $s_{11} > s_{12}$ and $s_{22} > s_{21}$. By simplifying Eq. (5) further, we have the following two regularizations:

$$\begin{aligned} \text{EQSIM}_{v1} : s_{12} &= s_{21} \\ \text{EQSIM}_{v2} : s_{11} - s_{12} &= s_{22} - s_{21}, \quad s_{11} - s_{21} = s_{22} - s_{12}. \end{aligned} \quad (6)$$

Note that the viable space of EQSIM_{v2} is a subset of EQSIM_{v1}, because EQSIM_{v1} is exactly equivalent to Eq. (5) while EQSIM_{v2} further requires $s_{11} = s_{22}$. Empirically, we find that EQSIM_{v2} is more suitable to the semantically close pairs (I_1, T_1) and (I_2, T_2) ; and EQSIM_{v1} to distant pairs. Figure 4 illustrates such hybrid training loss within a training batch. Semantically “close” and “distant” are determined by the similarity score s , where we regard samples with top- k s as “close” samples. For dual encoder VLMs with ITC loss, s is the cosine similarity between image and text features. For fusion encoder VLMs with ITM, s is the scoring output from the ITM head.

In our implementation, we adopt Mean Square Error (MSE) loss to regularize the equation of similarities. In





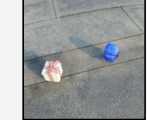







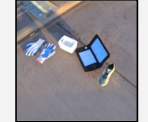

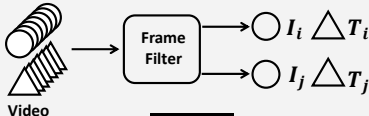
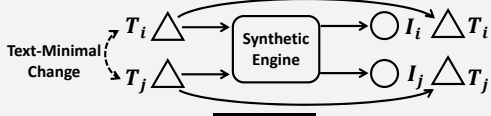
Name	EQ-YouCook2	EQ-GEBC	EQ-AG	Attribute	EQ-Kubric Counting	Location	EQ-SD
Example	 Put chicken and cheese on top of bread.	 Human hands pressing the cardboard with thumb.	 The person is sitting on the sofa which is behind him/her.	 The blue plant container is on the left side of the brown bull.	 1 white gift box with red straps is on the platform.	 The blue gloves is in front of the fabric basket.	 A photo of beach.
	 Spread butter over a slice of bread.	 Human hands adjust the cardboard on the paper.	 The person is lying on the sofa which is beneath him/her.	 The red plant container is on the left hand of the brown bull.	 There are 4 white gift box with red straps on the platform.	 The blue gloves is on the left hand of the fabric basket.	 A photo of beach at sunset.
# Sample	45849	1378	195872	2000	2000	2000	1513
Instrument	Act./Obj.	Act./Obj.	Act.	Attr.	Count	Location	Attr./Obj.
Source Data	YouCook2	GEBC	Action Genome	Google Scanned Objects (GSO)			N/A
Pipeline							

Figure 5: Overview of the proposed benchmark EQBEN, which consists of 5 sub-datasets and can be categorized to natural and synthetic. Act., Obj., Attr. denote action, object and attribute, respectively.

addition, motivated by the hinge loss [3, 63], we utilize a margin parameter α to control the strength of regularization. EQSIM_{v1} can be written as $[\|s_{12} - s_{21}\|_2^2 - \alpha]_+$, where $[x]_+ = \max(x, 0)$ and $\|\cdot\|_2$ denotes the L2 norm. EQSIM_{v2} can be implemented similarly. In practice, given a retrieval fine-tuning objective \mathcal{L}_{Ret} , the final loss can be written as: $\mathcal{L} = \mathcal{L}_{Ret} + \beta\mathcal{L}_{EQ}$, where \mathcal{L}_{EQ} is EQSIM_{v1} (EQSIM_{v2}) for semantically distant (close) samples, β is the balancing factor. Experiments in Section 5.4 validate that the hybrid training is better than only using EQSIM_{v1}.

4. Diagnosing VLMs with EQBEN

We argue that standard VL evaluation kits [46, 39] are too coarse to evaluate the equivariance of VLMs similarity. Existing VLMs can easily distinguish most samples in conventional retrieval benchmarks, *e.g.*, images of a group of people against those with cars, given a caption of “people standing on the street”. Therefore, we propose EQBEN to focus on visual minimal semantic changes to check whether VLMs can faithfully respond, *i.e.*, the equivariance of the similarity measure in VLMs. Specifically, EQBEN contains 5 sub-datasets, covering diverse image domains, from real-life scenarios to synthetic well-controlled scenes. And it is designed to stress test VLMs with accurate semantic changes in action, location, and attribution (*e.g.*, color, count and size). Figure 5 presents an overview of EQBEN.

Next, we introduce our design principle for constructing EQBEN. Each sample in EQBEN consists of a pair of

images (I_1, I_2) and a pair of captions (T_1, T_2). A valid EQBEN sample must satisfy: (1) T_i is preferred to be used as the description for I_i ; (2) I_1 and I_2 are visually-minimally different. The former one requires that $\{I_1, T_1\}$ and $\{I_2, T_2\}$ should be semantically distinguishable without confusion, while the latter one limits the extent of the distinction – “visual-minimal change”. Previous work [62] defines “minimal” semantic change in the caption space as the same words but in a different order. However, due to the continuity and entanglement of image pixels, “minimal” semantic change in visual space is hard to determine. In this paper, we roughly define it as changes in the foreground (*e.g.*, attribute, action, location, *etc.*) while sharing the same scene and background.

In practice, we source image pairs with “visual-minimal change” in two ways: (1) from natural videos and (2) from synthetic engines, where we adopt different construction pipelines, as shown at the bottom of Figure 5. For the former one, we directly leverage the continuity of scene changes along the temporal dimension in *natural videos*, which can provide massive image pairs with minimal visual changes. Specifically, we leverage the existing video-language datasets [80, 26, 70] to construct EQBEN samples. To more precisely control the varying component in images, we further explore the photo-realistic scene generator (Kubric [20]) and the open-source diffusion model (Stable Diffusion [51, 24]) to synthetically generate pairs of images by providing two captions that are minimally different from

each other. In what follows, we introduce the construction pipeline for each sub-dataset in detail.

4.1. Construction from Natural Video

Let’s define a video with caption annotations as $\mathcal{V} = \{I_i, T_i\}_{i=1}^N$, where N is the number of sampled frames. We assume the “visual-minimal change” is naturally guaranteed between any two frames I_i and I_j ideally, where $I_i, I_j \in \mathcal{V}$, $i \neq j$, as we limit the source video to be either short segment [70, 26] or capturing a fixed scene [80]. However, we find that it is hard to ensure the validity of all the video frame pairs in practice. Therefore, we utilize a frame filter to filter out invalid samples automatically for different video sources. Below, we briefly introduce the dataset construction process and delay the details to Appendix.

We construct three sub-datasets based on real images from natural videos, including EQ-AG, EQ-GEBC and EQ-YOUCOOK2. We construct **EQ-AG** by leveraging the scene graph annotations from Action Genome (AG) [26], which capture detailed changes between objects and their pairwise relationships while action occurs. We first use a slot-filling template to translate scene graphs to captions. As videos usually come with redundant frames, we avoid the nearly duplicated frames by sampling frames I_i and I_j if and only if at least 2 of 3 pairwise relationships are different. Furthermore, if I_j is chosen for previous samples, we empirically skip the subsequent 2 frames to I_{j+2} . **EQ-GEBC** is built on GEBC [70] that contains captions describing the event before and after an event boundary. We adopt the frames before and after the boundary as our visual minimally different images. Similarly, we avoid temporal redundancy via sparse sampling across multiple boundaries. We construct **EQ-YOUCOOK2** based on YouCook2 [80], which is sparsely annotated with captions for each cooking step. We construct the dataset by sampling the middle frame of a short video segment as I_i with its annotated caption as T_i . We then apply off-the-shelf object detectors to filter scene changes. Please note that EQBEN can be easily extended to other video-language datasets by applying the same construction pipeline.

4.2. Construction from Synthetic Engine

Synthetic engine may provide more precise and controllable visual changes in the generated images, to allow more accurate diagnosis in terms of model failure when evaluating with EQBEN. We assume that a synthetic engine can faithfully generate images based on a text prompt describing the image content. Based on this assumption, given a pair of semantic-minimally different captions T_i and T_j , we expect the generated I_i and I_j to be correspondingly visual-minimally different. In the following, we briefly introduce the utilized engine and how to construct semantic-minimally different captions for each sub-dataset and leave

Dataset	# Testing Samples	Visual Semantic Change	Pairwise	Domain Diversity	Scalability
Flickr30K [46]	1K	✗	✗	✗	✓
COCO [7]	1/5K	✗	✗	✗	✓
VALSE [44]	6795	✗	✓	✓	✓
Winoground [62]	400	✗	✓	✗	✗
EQBEN (Ours)	250K	✓	✓	✓	✓

Table 1: Comparison between EQBEN and related benchmarks.

more details to Appendix.

EQ-KUBRIC takes advantages of Kubric [20], an open-source graphics engine to generate photo-realistic scenes. Here we adopt Google Scanned Objects (GSO) for scene construction and categorize caption change into three aspects: *attribute*, *counting*, and *location*. For each aspect, we construct 2000 image-text pairs by intervening corresponding phrases of sentences while leaving other words unchanged. **EQ-SD** is inspired by the recent advances in diffusion models for text-to-image generation [49, 53, 75]. We utilize the open-source checkpoint v1.4 of Stable Diffusion (SD) with prompt-to-prompt image editing framework [24] to translate two semantic-minimally different captions to a pair of images. Specifically, we elaborately design a set of textual semantic-minimal editing: 1) object change (e.g., “dog”→“cat”); 2) scene change (e.g., + “in the winter”); 3) attribute change (e.g., + “with a sunglasses”). Finally, we perform a human evaluation to filter out poor-quality generations. Notably, we can adopt more rendered objects (e.g., rendered animals) and various synthetic engine (e.g., better generative models) to further extend our EQBEN following the proposed pipeline.

4.3. Comparisons with Other Datasets

In Table 1, we conduct a direct comparison of EQBEN against two widely adopted retrieval benchmarks (Flickr30K [46] and COCO [7]) and two recent datasets with textual-minimal change (VALSE [44] and Winoground [62]) from four aspects. 1) On the dataset characteristics, to the best of our knowledge, EQBEN is the first diagnosing benchmark to examine the equivariance of VLMs in terms of **minimal visual semantic change**. 2) For evaluation setting, **pairwise** setting asks VLMs to select the correct counterpart within a pair of slightly different samples rather than thousands of very different samples in conventional retrieval datasets. The minimal semantic drift between the pair of samples makes the evaluation of equivariant similarity measure more effective. 3) For **domain diversity**, our EQBEN contains rich visual contents collected from different video domains as well as synthetic domains, as opposed to the common image-text datasets which existing diagnosing kits are built upon. 4) In terms of **scalability**, EQBEN is highly scalable as our automatic pipeline can be easily applied to other video-language datasets and synthetic engines, as opposed to manual annotation for building

Method	Winoground			VALSE		F30K Text-to-image Ret.			F30K Image-to-text Ret.		
	Text	Image	Group	$\min(p_c, p_f)$	acc	R@1	R@5	R@10	R@1	R@5	R@10
METER [14]	39.25	15.75	11.99	25.43	54.01	79.60	94.96	97.28	90.90	98.30	99.50
+ FT (F30K) [14]	43.50	20.75	14.75	22.41	53.06	82.22	96.34	98.36	94.30	99.60	99.90
+ EQSIM	44.99	22.75	18.75	30.08	53.38	82.16	94.70	96.64	95.30	99.60	99.90
FIBER [13]	46.25	25.75	22.24	11.42	51.76	79.26	95.70	97.92	91.60	99.50	99.80
+ FT (F30K) [13]	51.24	26.49	23.00	20.78	54.90	81.44	96.72	98.48	92.90	99.50	99.90
+ EQSIM	51.49	31.49	27.50	52.42	58.06	83.56	96.78	98.28	96.00	99.60	99.90

Table 2: Results of EQSIM on the challenging Winoground [62], VALSE [44] benchmark and Flickr-30K (F30K) [46] test split for image-text retrieval. FT and Ret. are short for fine-tuning and retrieval.

traditional retrieval datasets. While VALSE only focuses on linguistic editing of the captions, EQBEN can be further scale up with more diverse visual contents.

5. Experiments

We first introduce our experimental setting in Section 5.1, followed by evaluation of EQSIM on existing benchmarks in Section 5.2. Section 5.3 benchmarks SOTA VLMs on EQBEN to show their insensitivity to minimal visual semantic changes, and we further validate EQSIM on EQBEN. Section 5.4 presents additional ablation studies to examine the design of EQSIM.

5.1. Experimental Setting

Training Details. Recent efforts on diagnosing benchmarks [62, 44] only provide testing data and directly evaluate models after VL pre-training. The low performance reported on these benchmarks can mainly be attributed to two factors: 1) the inherent weaknesses of VLMs, *e.g.*, non-equivariant similarity measure; and 2) the domain gap between training and testing. To better validate the effectiveness of our method, we fine-tune the VLMs on limited image-text pairs from conventional retrieval dataset Flickr30K [46] with or without the regularization term of EQSIM, and then test the fine-tuned VLMs on the challenging Winoground [62], VALSE [44] and our EQBEN. Under a fair comparison, we argue that the absolute performance improvements from EQSIM thus would suggest that the gain is entirely from the remedy of model weaknesses.

To validate the effectiveness and generalizability of our proposed method, we apply EQSIM to two SOTA end-to-end methods with different architectures and retrieval losses. Specifically, FIBER [13] supports the dual encoder with ITC loss for fast retrieval, which computes similarities for N^2 image-text pairs with only $O(N)$ forwarding. In contrast, the SOTA fusion-encoder model METER [14], optimized with ITM task during pre-training, computes the similarity by forwarding the concatenation of each pair of image and text, resulting in $O(N^2)$ time complexity. Fine-tuning details for each model can be found in Appendix.

Evaluation Metric. On **Winoground** [62], given two image-text pairs $\{I_1, T_1\}$ and $\{I_2, T_2\}$, a VLM mea-

sures similarity s_{ij} between image I_i and text T_j ($i, j \in \{0, 1\}, i \neq j$). Three metrics are computed based on s_{ij} : 1) *Text score* measures whether the model can select the correct text for a given image. The model wins one point if $s_{ii} > s_{ij}$. 2) *Image score* evaluates if VLMs can select the correct image for a given text and the model wins one point when $s_{ii} > s_{ji}$. 3) *Group score* combines the previous two, such that the VLMs win one point if and only if both text score and image score are 1, meaning the following condition must be satisfied: $s_{ii} > s_{ij}$ and $s_{ii} > s_{ji}$. On **VALSE** [44], the two image-text pairs share a common image, *i.e.*, $\{I_1, T_1\}$ (correct) and $\{I_1, T_2\}$ (foil). We follow [44] to report the following metrics: 1) *acc* is the overall accuracy on both correct and foil image-text pairs; and 2) $\min(p_c, p_f)$ is the minimum of precision p_c and foil precision p_f , where p_c (p_f) measures how well models identify the correct (foil) pair. We also report performance on the conventional image-text retrieval task, where recall R@K (K=1,5,10) is used as the evaluation metric.

5.2. Evaluation of EQSIM

In Table 2, we compare model performance under three settings: (i) direct evaluation after pre-training (the first rows of each block); (ii) standard fine-tuning (FT) on Flickr30K training data (the second rows of each block); and (iii) fine-tuning with EQSIM regularization (the third rows of each block). We observe that standard fine-tuning can somewhat bring a little performance improvement on both Winoground and VALSE benchmarks, indicating that some domain overlap between Flickr30K training data and testing samples. It is difficult to entirely rule out the domain influence, but comparing fine-tuning with EQSIM against standard fine-tuning, our method brings consistent and significant performance improvements on both of the challenging Winoground and VALSE across METER and FIBER models. Specifically, EQSIM improves the group score over standard fine-tuning by 4% for METER and 4.5% for FIBER on Winoground, respectively. While for VALSE, the performance improvement on $\min(p_c, p_f)$ is as large as 31.6%, further validating the effectiveness of our EQSIM. In addition, we observe that the equivariance regularization from EQSIM does not sacrifice retrieval performance. On Flickr30K, EQSIM can mostly retain the retrieval perfor-

Method	Natural Subsets									Synthetic Subsets						Avg
	EQ-YouCook2			EQ-GEBC			EQ-AG			EQ-KUBRIC			EQ-SD			
	Text	Image	Group	Text	Image	Group	Text	Image	Group	Text	Image	Group	Text	Image	Group	
LXMERT [60]	13.96	11.98	4.55	13.56	12.73	4.19	18.17	9.02	4.46	18.50	15.35	7.26	11.16	6.15	1.98	10.20
ViLBERT [41]	14.78	12.75	5.18	14.67	12.64	4.82	17.43	8.36	3.89	17.55	18.44	8.13	12.37	7.37	2.78	10.74
CLIP (RN-50) [48]	47.72	47.99	34.05	10.80	18.03	3.97	14.52	10.44	3.50	21.33	21.93	9.75	90.09	85.92	79.11	33.28
CLIP (ViT-B/32) [48]	49.48	51.10	36.50	12.57	20.12	4.47	13.91	8.72	3.32	20.56	21.29	9.66	89.16	86.05	78.98	33.73
FLAVA [59]	51.66	54.78	39.68	12.24	16.81	5.07	6.59	13.47	2.15	28.88	28.18	15.90	79.64	84.47	71.10	34.04
ViLT [30]	44.61	46.69	31.74	14.72	16.70	5.62	15.37	9.89	3.45	31.23	27.00	17.90	80.37	79.04	68.93	32.88
ALBEF [35]	57.01	58.04	44.90	13.56	19.63	5.89	11.28	15.17	3.93	29.87	30.18	18.58	88.96	90.41	83.07	38.03
BLIP [34]	59.22	58.36	46.31	15.87	19.79	7.27	19.76	13.87	6.31	29.38	32.25	18.73	85.39	85.52	77.13	38.34
METER [14]	52.18	49.42	36.81	20.95	18.19	6.95	28.70	15.80	7.88	44.28	35.20	27.26	89.62	84.93	79.44	39.84
+ FT (F30K) [14]	52.68	48.31	36.52	18.08	19.85	7.33	29.50	16.30	8.12	41.11	34.59	24.33	86.64	84.46	77.46	39.02
+ EQSIM	54.12	53.12	40.29	24.20	26.02	11.69	28.85	20.09	10.76	<u>43.68</u>	39.08	28.42	88.04	84.07	<u>77.79</u>	42.28
FIBER [13]	52.04	50.84	38.32	25.19	22.66	11.08	32.49	24.05	13.70	47.94	45.60	33.53	86.05	88.63	79.97	44.86
+ FT (F30K) [13]	57.70	56.46	44.33	18.24	21.33	8.54	26.99	18.69	9.24	50.31	46.06	34.66	90.48	86.64	81.29	43.40
+ EQSIM	58.26	57.10	45.10	<u>21.55</u>	26.07	<u>10.58</u>	<u>29.93</u>	<u>23.42</u>	<u>12.64</u>	51.90	48.40	37.38	90.81	85.98	80.70	45.32

Table 3: Results on EQBEN. Rows highlighted in gray are results with our EQSIM. Numbers in bold, underline respectively represent the best results and the inferior results compared to pre-training but better than fine-tuning baseline. Avg denotes the average over all scores.

Method	EQ-KUBRIC		
	Location	Counting	Attribute
LXMERT	2.15	1.89	17.90
ViLBERT [41]	1.98	2.03	20.40
CLIP (RN-50) [48]	1.25	5.20	22.80
CLIP (ViT-B/32) [48]	0.75	5.80	22.44
FLAVA [59]	1.00	7.35	39.35
ViLT [30]	1.95	6.19	45.55
ALBEF [35]	1.25	9.49	44.99
BLIP [34]	1.15	10.60	44.44
METER [14]	3.59	15.29	62.90
+ FT (F30K) [14]	2.40	17.49	53.10
+ EQSIM	3.80	23.85	57.60
FIBER [13]	11.34	19.65	69.59
+ FT (F30K) [13]	8.95	28.49	66.54
+ EQSIM	<u>11.05</u>	30.90	70.20

Table 4: Detailed results on EQ-KUBRIC. We report group score on three splits of EQ-KUBRIC, capturing the visual semantic changes in Location, Counting and Attribute.

performance, and sometimes even yield performance gain, *e.g.*, 3.1% on R@1 for image-to-text retrieval with FIBER.

5.3. Benchmarking VLMs with EQBEN

We evaluate a wide range of VLMs with different configurations on EQBEN in a zero-shot manner, to examine the equivariance of their similarity measures for distinguishing visually-minimal different samples. We consider representative VLMs, including (i) LXMERT [60], ViLBERT [41] for OD-Based models; and (ii) CLIP [48] variants, FLAVA [58], ViLT [30], ALBEF [35], BLIP [34] and METER [14] and FIBER [13] as prominent examples of end-to-end SOTA methods. Full results on more VLMs can be found in Appendix A.6. For evaluation metrics, we adopt text score, image score and group score to compare model performance, similar to Winoground [62].

Table 3 presents the evaluation results of existing VLMs on EQBEN and we summarize our observations below.

- Regardless of the subsets, end-to-end VLMs generally achieve better performance as it is not constrained by the fixed visual representation from a pre-trained object detector [50], as in OD-based methods.
- Among all subsets, VLMs obtain evidently higher performance on EQ-SD. The stable diffusion model [51] is pre-trained on similar VL corpus to these VLMs. Hence, the generated images can be biased towards the same underlying data distribution, much easier for VLMs to tell the differences. Besides, the generated images maybe visually minimally different to human eyes, but it is unclear whether in the pixel space, they are minimally different w.r.t. the model input. It is worth noting that LXMERT and ViLBERT are the exception due to the totally different distribution with the off-the-shelf object detector.
- Interestingly, a larger pre-training corpus (*e.g.*, CLIP [48] and FLAVA [58]) does not always guarantee better results. This implies training loss may be more critical in learning equivariant similarity measure.
- In Table 4, we further conduct a fine-grained examination with the synthetic subset EQ-KURIC, where we focus on specific visual changes in location, counting and attribute. VLMs fail substantially in terms of location and counting, while being sensitive to attribute changes. Similar findings are also observed by [62, 44] from the text side.

We again equip the two strong baseline models (METER and FIBER) with EQSIM and fine-tune on Flickr30K. As EQBEN covers diverse domains, standard fine-tuning on Flickr30K can hardly improve or even hurt model performance, compared with direct evaluation after pre-training (with -0.62% and -1.46% performance drop for METER and FIBER, respectively). However, by enforcing equivariant constraint with EQSIM, we observe significant performance improvements than standard fine-tuning, with an absolute gain of 3.26% for METER and 1.92% for FIBER.

FT Data	Method	EQ-AG	EQ-Y.	EQ-G.	Wino.	Avg
F30K [46]	FT	9.24	44.33	8.54	23.00	21.28
	+ EQSIM	12.64	45.10	10.58	27.50	23.96
COCO [7]	FT	10.14	42.90	8.93	21.50	20.87
	+ EQSIM	12.52	45.68	9.37	25.75	23.33
F30K + COCO	FT	9.96	43.81	8.93	22.75	21.36
	+ EQSIM	11.98	45.80	10.47	26.50	23.69
4M [†]	FT	10.49	40.95	7.49	20.99	19.98
	+ EQSIM	12.78	40.96	9.81	21.25	21.20

Table 5: Group accuracy (%) of fine-tuning (FT) and EQSIM based on FIBER on different FT data corpus. 4M data denotes the commonly used pre-training data with about 4M images, including COCO, Visual Genome [31], Conceptual Captions [55] and SBU [43]. Y., G., Wino. are the short for YouCook2, GEBC and Winoground. † The model is fine-tuned for 10K steps.

5.4. Ablation Study

In this section, we conduct ablation studies to validate the scalability, design and effectiveness of EQSIM in terms of enforcing equivariant similarity.

Scalability of EQSIM. Table 5 evaluates the scalability of EQSIM and standard fine-tuning baseline on the natural subsets of EQBEN and Winoground by gradually including more training data. Under the same fine-tuning data, EQSIM achieves consistent and significant improvements (2% - 3%) over the baseline. Interestingly, there is no remarkable correlation between the corpus size and model performance. This may be due to the distribution of standard VL data is far away from that of EQBEN and Winoground. Note that for the 4M experiment, we fine-tune the models for 10K steps due to computational constraints. Our results demonstrate the potential of EQSIM to benefit VL pre-training on large-scale data. Additionally, we validate the generalizability of EQSIM in other relevant downstream tasks. Further details are provided in Appendix A.7. **Ablation on EQSIM design.** Table 6 compares EQSIM against the four ablated instances on EQ-KUBRIC and Winoground [62], including 1) fine-tuning with hard negative sampling (HardNeg); 2) applying EQSIM_{v1} to all samples in the training batch (EQSIM_{v1}-all); 3) applying EQSIM_{v2} to all samples in the training batch (EQSIM_{v2}-all); and 4) applying EQSIM_{v2} for only semantically close samples (EQSIM_{v2}-close). The final EQSIM is equivalent to EQSIM_{v1}-all + EQSIM_{v2}-close, which achieves the best performance. Notably, enforcing EQSIM_{v2} on all (EQSIM_{v2}-all) even degrades the performance by -0.58% on average, compared to applying only to semantically close samples (EQSIM_{v2}-close). This validates our claim in Section 3 that EQSIM_{v2} is better suited for semantically close samples.

Validation of equivariance via EQSIM. Given the similarity scores s calculated by a VLM, we can define the equivariance score as the derivation of EQSIM_{v2} (headline of Figure 6) to measure the degree of equivariance (the smaller, the better). In Figure 6, we plot the distribution of EQSIM_{v2}

Method	EQ-KUBRIC			Wino.	Avg
	Location	Counting	Attribute		
FT (F30K)	8.95	28.49	66.54	22.24	31.55
+ HardNeg	10.89	29.49	67.69	27.00	33.77
+ EQSIM _{v1} -all	9.79	29.94	68.75	26.49	33.74
+ EQSIM _{v2} -all	10.25	29.05	68.30	25.49	33.27
+ EQSIM _{v2} -close	11.15	29.25	69.25	25.75	33.85
+ EQSIM	11.05	30.90	70.20	27.50	34.91

Table 6: Ablation studies of the loss design for our EQSIM on EQ-KUBRIC and Winoground (Wino.) using group score (%).

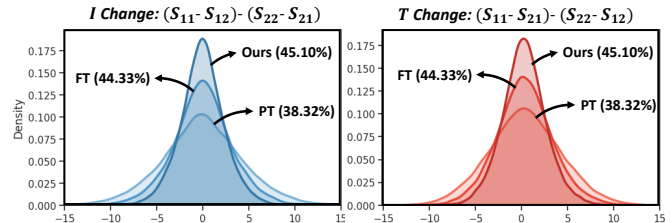


Figure 6: The equivariant score of three FIBER [13] variant models: Pre-trained (PT), Fine-tuned (FT) and Ours (EQSIM) on EQ-YOUCOOK2. The equivariant score is defined by the derivation of EQSIM_{v2}. The tighter distribution curve, the better equivariant similarity measure.

values across all samples in EQ-YOUCOOK2 dataset for FIBER [13] and its variants, attached with their group scores. A tighter curve indicates smaller derivation, hence better equivariance similarity measure. Full results on other EQBEN subset are presented in Appendix A.8. Compared with pre-training only (PT), fine-tuning on Flickr30K (FT) can improve the group score while being more equivariant in the similarity measure. Adding EQSIM (Ours) obtains additional improvements on both similarity equivariance and group score, indicating EQSIM indeed enforces equivariant similarity measure. Additionally, due to the space limitation, we leave more visualizations in Appendix A.9.

6. Conclusion

In this study, we investigated the non-equivariant similarity issue in VLMs, hidden behind their excellent performances on standard evaluation benchmarks. To address this issue, we proposed Equivariance Similarity Learning (EQSIM), an elegant and effective regularization method that can be easily integrated into the fine-tuning process of existing VLMs. Meanwhile, to better diagnose the equivariance of VLMs, we further introduced a new challenging benchmark EQBEN, the first to focus on “visual-minimal change”. Our proposed EQSIM is backed by the strong results on both challenging benchmarks (*e.g.*, Winoground, VALSE, EQBEN) and the conventional Flickr30K dataset. In future work, we plan to explore the application of EQSIM in VL pre-training and instruction tuning. **Acknowledgement.** This work is partly supported by AISG, A*STAR under its AME YIRG Grant (Project No.A20E6c0101).

References

- [1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *arXiv preprint*, 2022. [2](#)
- [2] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. VQA: Visual question answering. In *ICCV*, 2015. [3](#)
- [3] Bernhard E Boser, Isabelle M Guyon, and Vladimir N Vapnik. A training algorithm for optimal margin classifiers. In *Proceedings of the fifth annual workshop on Computational learning theory*, pages 144–152, 1992. [4](#)
- [4] Michael M Bronstein, Joan Bruna, Taco Cohen, and Petar Veličković. Geometric deep learning: Grids, groups, graphs, geodesics, and gauges. *arXiv preprint arXiv:2104.13478*, 2021. [3](#)
- [5] Remi Cadene, Corentin Dancette, Matthieu Cord, Devi Parikh, et al. Rubi: Reducing unimodal biases for visual question answering. *NeurIPS*, 32, 2019. [3](#)
- [6] Hongge Chen, Huan Zhang, Pin-Yu Chen, Jinfeng Yi, and Cho-Jui Hsieh. Attacking visual language grounding with adversarial examples: A case study on neural image captioning. *arXiv preprint arXiv:1712.02051*, 2017. [3](#)
- [7] 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*, 2015. [3](#), [5](#), [8](#), [12](#), [14](#)
- [8] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. UNITER: Universal image-text representation learning. In *ECCV*, 2020. [2](#)
- [9] Jaemin Cho, Abhay Zala, and Mohit Bansal. Dall-eval: Probing the reasoning skills and social biases of text-to-image generative transformers. *arXiv preprint arXiv:2202.04053*, 2022. [1](#)
- [10] Taco Cohen and Max Welling. Group equivariant convolutional networks. In *ICML*, pages 2990–2999. PMLR, 2016. [3](#)
- [11] Taco S Cohen, Mario Geiger, Jonas Köhler, and Max Welling. Spherical cnns. *arXiv preprint arXiv:1801.10130*, 2018. [3](#)
- [12] Rumen Dangovski, Li Jing, Charlotte Loh, Seungwook Han, Akash Srivastava, Brian Cheung, Pulkit Agrawal, and Marin Soljačić. Equivariant contrastive learning. *arXiv preprint arXiv:2111.00899*, 2021. [3](#)
- [13] Zi-Yi Dou, Aishwarya Kamath, Zhe Gan, Pengchuan Zhang, Jianfeng Wang, Linjie Li, Zicheng Liu, Ce Liu, Yann LeCun, Nanyun Peng, et al. Coarse-to-fine vision-language pre-training with fusion in the backbone. *arXiv preprint arXiv:2206.07643*, 2022. [1](#), [2](#), [6](#), [7](#), [8](#), [12](#), [15](#), [17](#), [20](#)
- [14] Zi-Yi Dou, Yichong Xu, Zhe Gan, Jianfeng Wang, Shuohang Wang, Lijuan Wang, Chenguang Zhu, Pengchuan Zhang, Lu Yuan, Nanyun Peng, Zicheng Liu, and Michael Zeng. An empirical study of training end-to-end vision-and-language transformers. In *CVPR*, 2022. [1](#), [2](#), [6](#), [7](#), [15](#), [20](#)
- [15] Zhe Gan, Yen-Chun Chen, Linjie Li, Chen Zhu, Yu Cheng, and Jingjing Liu. Large-scale adversarial training for vision-and-language representation learning. In *NeurIPS*, 2020. [2](#)
- [16] Zhe Gan, Linjie Li, Chunyuan Li, Lijuan Wang, Zicheng Liu, and Jianfeng Gao. Vision-language pre-training: Basics, recent advances, and future trends. *arXiv preprint arXiv:2210.09263*, 2022. [2](#)
- [17] Shashank Goel, Hritik Bansal, Sumit Bhatia, Ryan A Rossi, Vishwa Vinay, and Aditya Grover. Cyclip: Cyclic contrastive language-image pretraining. *arXiv preprint arXiv:2205.14459*, 2022. [3](#), [14](#)
- [18] Tao Gong, Chengqi Lyu, Shilong Zhang, Yudong Wang, Miao Zheng, Qian Zhao, Kuikun Liu, Wenwei Zhang, Ping Luo, and Kai Chen. Multimodal-gpt: A vision and language model for dialogue with humans. *arXiv preprint arXiv:2305.04790*, 2023. [12](#)
- [19] Jonathan Gordon, David Lopez-Paz, Marco Baroni, and Diane Bouchacourt. Permutation equivariant models for compositional generalization in language. In *ICLR*, 2019. [3](#)
- [20] Klaus Greff, Francois Belletti, Lucas Beyer, Carl Doersch, Yilun Du, Daniel Duckworth, David J Fleet, Dan Gnanaprasgam, Florian Golemo, Charles Herrmann, Thomas Kipf, Abhijit Kundu, Dmitry Lagun, Issam Laradji, Hsueh-Ti (Derek) Liu, Henning Meyer, Yishu Miao, Derek Nowrouzezahrai, Cengiz Oztireli, Etienne Pot, Noha Radwan, Daniel Rebain, Sara Sabour, Mehdi S. M. Sajjadi, Matan Sela, Vincent Sitzmann, Austin Stone, Deqing Sun, Suhani Vora, Ziyu Wang, Tianhao Wu, Kwang Moo Yi, Fangcheng Zhong, and Andrea Tagliasacchi. Kubric: a scalable dataset generator. 2022. [4](#), [5](#), [18](#)
- [21] Tanmay Gupta, Ryan Marten, Aniruddha Kembhavi, and Derek Hoiem. Grit: General robust image task benchmark. *arXiv preprint arXiv:2204.13653*, 2022. [3](#)
- [22] Tengda Han, Weidi Xie, and Andrew Zisserman. Self-supervised co-training for video representation learning. *Advances in Neural Information Processing Systems*, 33:5679–5690, 2020. [3](#)
- [23] Lisa Anne Hendricks and Aida Nematzadeh. Probing image-language transformers for verb understanding. *arXiv preprint arXiv:2106.09141*, 2021. [3](#)
- [24] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626*, 2022. [4](#), [5](#)
- [25] Hexiang Hu, Ishan Misra, and Laurens van der Maaten. Evaluating text-to-image matching using binary image selection (bison). In *ICCV Workshops*, pages 0–0, 2019. [3](#)
- [26] Jingwei Ji, Ranjay Krishna, Li Fei-Fei, and Juan Carlos Niebles. Action genome: Actions as compositions of spatio-temporal scene graphs. In *CVPR*, pages 10236–10247, 2020. [4](#), [5](#), [16](#)
- [27] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc V Le, Yunhsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *ICML*, 2021. [1](#), [2](#)
- [28] Carlos E Jimenez, Olga Russakovsky, and Karthik Narasimhan. Carets: A consistency and robustness evalu-

- ative test suite for vqa. *arXiv preprint arXiv:2203.07613*, 2022. 3
- [29] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, et al. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950*, 2017. 16
- [30] Wonjae Kim, Bokyung Son, and Ildoo Kim. Vilt: Vision-and-language transformer without convolution or region supervision. In *ICML*, pages 5583–5594. PMLR, 2021. 7, 15
- [31] 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. *IJCV*, 2017. 8
- [32] Benno Krojer, Vaibhav Adlakha, Vibhav Vineet, Yash Goyal, Edoardo Ponti, and Siva Reddy. Image retrieval from contextual descriptions. In *ACL*, pages 3426–3440, 2022. 3
- [33] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*, 2023. 15
- [34] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. BLIP: Bootstrapping language-image pre-training for unified vision-language understanding and generation. *arXiv preprint*, 2022. 1, 2, 7, 15
- [35] Junnan Li, Ramprasaath R Selvaraju, Akhilesh Deepak Gotmare, Shafiq Joty, Caiming Xiong, and Steven Hoi. Align before fuse: Vision and language representation learning with momentum distillation. In *NeurIPS*, 2021. 1, 2, 7, 15
- [36] Linjie Li, Jie Lei, Zhe Gan, and Jingjing Liu. Adversarial vqa: A new benchmark for evaluating the robustness of vqa models. In *ICCV*, pages 2042–2051, 2021. 3
- [37] 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*, 2019. 2
- [38] Xiujun Li, Xi Yin, Chunyuan Li, Pengchuan Zhang, Xiaowei Hu, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, et al. Oscar: Object-semantics aligned pre-training for vision-language tasks. In *ECCV*, 2020. 2
- [39] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: Common objects in context. In *ECCV*, 2014. 4
- [40] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*, 2023. 12
- [41] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. ViLBERT: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In *NeurIPS*, 2019. 2, 7, 15
- [42] Yulei Niu, Kaihua Tang, Hanwang Zhang, Zhiwu Lu, Xian-Sheng Hua, and Ji-Rong Wen. Counterfactual vqa: A cause-effect look at language bias. In *CVPR*, pages 12700–12710, 2021. 3
- [43] Vicente Ordonez, Girish Kulkarni, and Tamara Berg. Im2Text: Describing images using 1 million captioned photographs. In *NeurIPS*, 2011. 8
- [44] Letitia Parcalabescu, Michele Cafagna, Lilitta Muradjan, Anette Frank, Iacer Calixto, and Albert Gatt. Valse: A task-independent benchmark for vision and language models centered on linguistic phenomena. *arXiv preprint arXiv:2112.07566*, 2021. 1, 2, 3, 5, 6, 7
- [45] Mandela Patrick, Po-Yao Huang, Yuki Asano, Florian Metze, Alexander Hauptmann, Joao Henriques, and Andrea Vedaldi. Support-set bottlenecks for video-text representation learning. *arXiv preprint arXiv:2010.02824*, 2020. 3
- [46] Bryan A Plummer, Liwei Wang, Chris M Cervantes, Juan C Caicedo, Julia Hockenmaier, and Svetlana Lazebnik. Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. In *ICCV*, pages 2641–2649, 2015. 2, 3, 4, 5, 6, 8, 14
- [47] Guo-Jun Qi, Liheng Zhang, Feng Lin, and Xiao Wang. Learning generalized transformation equivariant representations via autoencoding transformations. *TPAMI*, 2020. 3
- [48] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, pages 8748–8763. PMLR, 2021. 1, 2, 3, 7, 15
- [49] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022. 1, 5
- [50] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. *TPAMI*, 2016. 2, 7
- [51] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, pages 10684–10695, 2022. 1, 4, 7
- [52] Halsey Lawrence Royden and Patrick Fitzpatrick. *Real analysis*, volume 32. Macmillan New York, 1988. 2
- [53] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint arXiv:2205.11487*, 2022. 5
- [54] 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. 1
- [55] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *ACL*, 2018. 8
- [56] Ravi Shekhar, Sandro Pezzelle, Yauhen Klimovich, Aurélie Herbelot, Moin Nabi, Enver Sangineto, and Raffaella Bernardi. Foil it! find one mismatch between image and language caption. *arXiv preprint arXiv:1705.01359*, 2017. 3
- [57] Mike Zheng Shou, Stan Weixian Lei, Weiyao Wang, Deepti Ghadiyaram, and Matt Feiszli. Generic event boundary de-

- tection: A benchmark for event segmentation. In *ICCV*, pages 8075–8084, 2021. 13
- [58] Amanpreet Singh, Ronghang Hu, Vedanuj Goswami, Guillaume Couairon, Wojciech Galuba, Marcus Rohrbach, and Douwe Kiela. FLAVA: A foundational language and vision alignment model. In *CVPR*, 2022. 2, 7
- [59] Amanpreet Singh, Ronghang Hu, Vedanuj Goswami, Guillaume Couairon, Wojciech Galuba, Marcus Rohrbach, and Douwe Kiela. Flava: A foundational language and vision alignment model. In *CVPR*, pages 15638–15650, 2022. 7, 15
- [60] Hao Tan and Mohit Bansal. LXMERT: Learning cross-modality encoder representations from transformers. In *EMNLP*, 2019. 1, 2, 7, 15
- [61] Kaihua Tang, Yulei Niu, Jianqiang Huang, Jiaxin Shi, and Hanwang Zhang. Unbiased scene graph generation from biased training. In *CVPR*, pages 3716–3725, 2020. 3
- [62] Tristan Thrush, Ryan Jiang, Max Bartolo, Amanpreet Singh, Adina Williams, Douwe Kiela, and Candace Ross. Winoground: Probing vision and language models for visiolinguistic compositionality. In *CVPR*, pages 5238–5248, 2022. 1, 2, 3, 4, 5, 6, 7, 8
- [63] Vladimir Vapnik. *The nature of statistical learning theory*. Springer science & business media, 1999. 4
- [64] Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. Show and tell: A neural image caption generator. In *CVPR*, pages 3156–3164, 2015. 3
- [65] Peng Wang, An Yang, Rui Men, Junyang Lin, Shuai Bai, Zhikang Li, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework. *arXiv preprint*, 2022. 2
- [66] Tan Wang, Jianqiang Huang, Hanwang Zhang, and Qianru Sun. Visual commonsense r-cnn. In *CVPR*, pages 10760–10770, 2020. 3
- [67] Tan Wang, Zhongqi Yue, Jianqiang Huang, Qianru Sun, and Hanwang Zhang. Self-supervised learning disentangled group representation as feature. *NeurIPS*, 34:18225–18240, 2021. 3
- [68] Wenhui Wang, Hangbo Bao, Li Dong, Johan Bjorck, Zhiliang Peng, Qiang Liu, Kriti Aggarwal, Owais Khan Mohammed, Saksham Singhal, Subhojit Som, et al. Image as a foreign language: Beit pretraining for all vision and vision-language tasks. *arXiv preprint arXiv:2208.10442*, 2022. 1
- [69] Wenhui Wang, Hangbo Bao, Li Dong, and Furu Wei. VLMo: Unified vision-language pre-training with mixture-of-modality-experts. *arXiv preprint*, 2021. 2
- [70] Yuxuan Wang, Difei Gao, Licheng Yu, Weixian Lei, Matt Feiszli, and Mike Zheng Shou. Geb+: A benchmark for generic event boundary captioning, grounding and retrieval. In *ECCV*, pages 709–725. Springer, 2022. 4, 5, 13, 16
- [71] Zirui Wang, Jiahui Yu, Adams Wei Yu, Zihang Dai, Yulia Tsvetkov, and Yuan Cao. SimVLM: Simple visual language model pretraining with weak supervision. In *ICLR*, 2022. 2
- [72] Maurice Weiler and Gabriele Cesa. General e (2)-equivariant steerable cnns. *NeurIPS*, 32, 2019. 3
- [73] Yuyang Xie, Jianhong Wen, Kin Wai Lau, Yasar Abbas Ur Rehman, and Jiajun Shen. What should be equivariant in self-supervised learning. In *CVPR*, pages 4111–4120, 2022. 3
- [74] Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mostafa Seyedhosseini, and Yonghui Wu. CoCa: Contrastive captioners are image-text foundation models. *arXiv preprint*, 2022. 1, 2
- [75] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Guntjan 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*, 2022. 5
- [76] Licheng Yu, Patrick Poirson, Shan Yang, Alexander C. Berg, and Tamara L. Berg. Modeling context in referring expressions. In *ECCV*, 2016. 3
- [77] Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin Choi. From recognition to cognition: Visual commonsense reasoning. In *CVPR*, pages 6720–6731, 2019. 3
- [78] Pengchuan Zhang, Xiujuan Li, Xiaowei Hu, Jianwei Yang, Lei Zhang, Lijuan Wang, Yejin Choi, and Jianfeng Gao. VinVL: Revisiting visual representations in vision-language models. In *CVPR*, 2021. 2
- [79] Tiancheng Zhao, Tianqi Zhang, Mingwei Zhu, Haozhan Shen, Kyusong Lee, Xiaopeng Lu, and Jianwei Yin. VI-checklist: Evaluating pre-trained vision-language models with objects, attributes and relations. *arXiv preprint arXiv:2207.00221*, 2022. 3
- [80] Luowei Zhou, Chenliang Xu, and Jason J Corso. Towards automatic learning of procedures from web instructional videos. In *AAAI*, 2018. 4, 5, 16
- [81] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023. 12