

Unsupervised Vision-and-Language Pre-training Without Parallel Images and Captions

Liunian Harold Li[†], Haoxuan You^{*°}, Zhecan Wang^{*°}, Alireza Zareian[°],
Shih-Fu Chang[°] & Kai-Wei Chang[†]

[†]University of California, Los Angeles

[°]Columbia University

liunian.harold.li@cs.ucla.edu,
{hy2612, zw2627, az2407, sc250}@columbia.edu,
kwchang@cs.ucla.edu

Abstract

Pre-trained contextual vision-and-language (V&L) models have achieved impressive performance on various benchmarks. However, existing models require a large amount of parallel image-caption data for pre-training. Such data are costly to collect and require cumbersome curation. Inspired by unsupervised machine translation, we investigate if a strong V&L representation model can be learned through unsupervised pre-training without image-caption corpora. In particular, we propose to conduct “mask-and-predict” pre-training on text-only and image-only corpora and introduce the object tags detected by an object recognition model as anchor points to bridge two modalities. We find that such a simple approach achieves performance close to a model pre-trained with aligned data, on four English V&L benchmarks. Our work challenges the widely held notion that aligned data is necessary for V&L pre-training, while significantly reducing the amount of supervision needed for V&L models.

1 Introduction

Pre-trained contextual vision-and-language (V&L) models (Lu et al., 2019; Tan and Bansal, 2019; Li et al., 2019; Su et al., 2019; Chen et al., 2020c) have achieved high performance on various V&L tasks. However, different from contextual language models, such as BERT (Devlin et al., 2019a), which are trained on easily-accessible unannotated text corpora, existing V&L models are still a step away from self-supervision. They require a massive amount of aligned text-image pairs for “mask-and-predict” pre-training. Such aligned data are costly to collect and hard to scale up. For example, the widely used MS-COCO dataset (Chen et al., 2015) requires extensive

annotation from crowd workers.¹

In this paper, we explore *unsupervised V&L pre-training* with unaligned image and text corpora.² This research direction aligns with the theme of unsupervised and self-supervised learning that moves from heavily-annotated data to unannotated data, e.g. unsupervised machine translation (Lample et al., 2018) and unsupervised image captioning (Feng et al., 2019). Unsupervised V&L pre-training is highly desirable as in many domains, aligned data is scarce (e.g. multimodal hate speech detection (Kiela et al., 2020) and the medical domain (Li et al., 2020c)) and it is easier to collect unaligned text and images. In addition to its practical implication, our endeavour challenges the widely held notion that image-caption corpora is indispensable for pre-training (Lu et al., 2019) and brings valuable insight into the role that aligned data play in V&L pre-training.

We are inspired by works on multi-lingual contextual language models (Pires et al., 2019). If we treat an image as a set of regions and each region as a visual token (Dosovitskiy et al., 2020), V&L models share a similar goal with multi-lingual models as they both learn shared representations across different domains. Although a multi-lingual language model pre-trained on non-parallel corpora such as mBERT (Devlin et al., 2019b) cannot align or translate languages out-of-the-box, its representation spaces for different languages can be easily aligned with a linear probe (Conneau et al., 2020). This property suggests the existence of universal latent symmetries in the unaligned contextual embedding spaces and is believed to contribute to

¹Other datasets also require cumbersome curation. For example, while Conceptual Captions is crawled from the web, the authors report that from 5 billion images gathered over the Internet, only 3 million have paired high-quality captions after filtering (Sharma et al., 2018; Changpinyo et al., 2021).

²Following Lample et al. (2018) and Feng et al. (2019), we use the term “unsupervised” to refer to pre-training with unaligned data, while “supervised” refers to pre-training with aligned text and images.

*The two authors contributed equally.

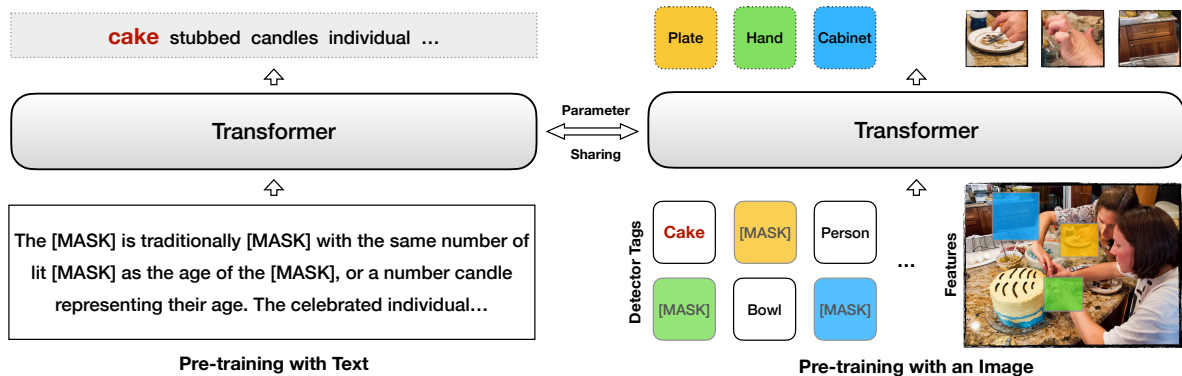


Figure 1: An illustration of pre-training without aligned data. Given text, the model is trained to predict masked words; given an image, the model is trained to predict masked regions and detector tags. The semantic class “cake” appears in both the language modality and the visual modality and is linked through the detector tags. Note that we do not require a text segment with the word *cake* to appear together with the image. Rather, we assume that as long as the text corpora are general enough, the word *cake* will appear in the textual modality eventually. The model can thus learn V&L representations from such weak supervision signals.

mBERT’s cross-lingual transfer ability. Thus we hypothesize that strong V&L representations can be similarly learned by “mask-and-predict” pre-training on unaligned language and vision data.

We propose unsupervised V&L pre-training with unaligned text and images (see an illustration in Figure 1). Specifically, we take VisualBERT (Li et al., 2019) as a running example and apply unsupervised pre-training, resulting in **Unsupervised VisualBERT** (U-VisualBERT). The model takes the form of a single Transformer that can accept inputs from both modalities. During each step of pre-training, unlike the existing models that observe a batch of text-image pairs, our model observes either a batch of text segments or a batch of images. When provided with text, part of the text is masked and the model is trained to predict the masked words; when provided with an image, part of the image regions are masked and the model is trained to predict properties of the masked regions.

To further encourage cross-modal fusion, we leverage the tags from an object detector as “anchor points” (Li et al., 2020b). For every object, we append its detected tag as a word to the visual input. The mask-and-predict objective is applied to the tags. For instance, for the image in Figure 1, the model can observe “*cake*” appears naturally as a word, a tag, and an image region. The direct typing of image regions and words can be learned and serves as a starting point for further alignment. The function of the detector tags resembles that of the “overlapping vocabulary” in multi-lingual language models, i.e., identical strings that appear in differ-

ent languages with the same meanings (e.g., “DNA” appears in both English and French). As the “overlapping vocabulary” improves cross-lingual transfer (Wu and Dredze, 2019), we argue the detector tags can improve cross-modal grounding.

We first conduct controlled experiments by pre-training on an English image-caption corpus without providing the alignment, following unsupervised machine translation and image captioning (Gu et al., 2019). Results on four English V&L benchmarks (VQA (Goyal et al., 2017), NLVR² (Suhr et al., 2019), Flickr30K Image Retrieval (Plummer et al., 2015), and RefCOCO+ (Yu et al., 2016)) show that U-VisualBERT achieves comparable performance as models with access to text-image pairs (Section 4).

Additionally, our approach is effective in practical settings, 1) when using independently collected images and captions and 2) when using images and general-domain text (BookCorpus (Zhu et al., 2015)) without any captions (Section 5.1). Quantitative and qualitative analysis confirms the anchoring effect of the detector tags (Section 5.2). As a byproduct, we conduct preliminary experiments to show the promise of the approach in a semi-supervised setting, where a hybrid model pre-trained with both aligned and additional unaligned data surpasses a model pre-trained only on aligned data. (Section 6). The above experiments demonstrate the wide applicability of our method. We will open-source the project under <https://github.com/uclanlp/visualbert>.

2 Related Work

Pre-trained V&L Transformers Various V&L models that are pre-trained with a “mask-and-predict” objective on aligned text-image data have been proposed (Lu et al., 2019; Tan and Bansal, 2019; Li et al., 2019; Su et al., 2019; Chen et al., 2020c; Li et al., 2020a; Zhou et al., 2020; Huang et al., 2020; Yu et al., 2020; Gan et al., 2020). Two kinds of designs have been proposed. Two-stream models (Lu et al., 2019; Tan and Bansal, 2019; Yu et al., 2020) utilize separate Transformers (Vaswani et al., 2017) for each modality and a cross-modality module is adopted. Single-stream models (Li et al., 2019; Su et al., 2019; Chen et al., 2020c) directly input the text and visual embeddings into one single Transformer. They have been widely used by downstream tasks (Kiela et al., 2020). Probing tasks (Cao et al., 2020) confirm that they capture useful V&L information after pre-training.

Two studies also try to incorporate “tag” information during pre-training. Oscar (Li et al., 2020b) adds detected tags as additional signals when pre-training with aligned data. We, however, do so for pre-training with unaligned data and show that the tags serve a more important role in unsupervised pre-training (Section 5.2). VIVO (Hu et al., 2020) targets novel object captioning. They use manually annotated image-tag data for pre-training and image-caption data for fine-tuning. We do not use manually annotated data and the tags are noisily generated by a detector.

Self-supervised Representation Learning

Self-supervision involves creating supervision objectives from natural data, often by corrupting the input and training the model to reconstruct the input (Kolesnikov et al., 2019) or contrastive learning (Chen et al., 2020b). Self-supervised training on language (Peters et al., 2018; Devlin et al., 2019a) such as BERT has been proven useful for various NLP tasks (Liu et al., 2019), while self-supervised visual representation learning has been centered around learning low-level visual features, in hope of enhancing the backbone CNN (Doersch et al., 2015; Pathak et al., 2016; Noroozi and Favaro, 2016; Chen et al., 2020b). In this paper, we conduct V&L pre-training by optimizing a reconstructive objective on unlabeled language-only and image-only data. Thus, our proposed model could be regarded as “self-supervised”. Notably, our contextual visual representation is

built on top of a pre-trained detector, operating at a level above local visual features.

Unsupervised Multi-lingual Language Model

This work is inspired by multi-lingual representations trained without parallel corpora (Devlin et al., 2019b). They are effective for cross-lingual transfer, which involves learning a model in one language and applying it to another with no additional training. Studies (Wu and Dredze, 2019; Conneau et al., 2020) have confirmed several design choices that facilitate such transfer, e.g. shared parameters and overlapping vocabularies across languages, and we make similar design choices in U-VisualBERT (Section 3.2). We argue that multi-lingual representations bear resemblance to multi-modal representations as both seek to encode the alignment between two domains (Chen et al., 2020a).

Unsupervised Grounding Learning

Prior works have explored learning grounding with weak or no supervision (Rohrbach et al., 2016; Xiao et al., 2017; Wang et al., 2020). Closest to this paper is unsupervised image captioning (Feng et al., 2019; Laina et al., 2019; Gu et al., 2019), which conducts image captioning with unpaired images and captions. Similar to this work, the detector tags serve as the anchor points for image captioning. However, unsupervised image captioning still requires captions, while our approach works with easy-to-collect general-domain text without any caption text (Section 5.1).

3 Approach

We first take Supervised VisualBERT (S-VisualBERT) as an example and illustrate how a typical V&L model is pre-trained with aligned data. Then we introduce unsupervised V&L pre-training, and the resulting model Unsupervised VisualBERT (U-VisualBERT).

3.1 Background

As mentioned in Section 2, there are several V&L representation learning methods based on BERT. We take Supervised VisualBERT (S-VisualBERT) as an example, which will also be used as a baseline in the experiments. S-VisualBERT is modified from the original VisualBERT (Li et al., 2019) and augmented with the visual objectives from LXMERT (Tan and Bansal, 2019) and detector tags similar to Oscar (Li et al., 2020b) (discussed in detail in Section 3.2).

Every input to S-VisualBERT contains a text segment T and an image I . The text and the image are first mapped into embedding vectors respectively. Text embeddings \mathbf{T} is a matrix in which each column vector represents the embedding of a subword in the text sequence, i.e. $\mathbf{T} = [\mathbf{w}_{1:n}]$. Following BERT, each subword embedding \mathbf{w}_i is the sum of its token, position, and segment embedding. Image embeddings \mathbf{I} include both the image region embeddings $\mathbf{r}_{1:m}$ and the detector tag embeddings $\mathbf{d}_{1:l}$ (see Section 3.2 for details). Each region embedding \mathbf{r}_i is the sum of a visual feature vector from the detector and a spatial box coordinate embedding (Tan and Bansal, 2019). The text and visual embeddings are then passed through a Transformer to built contextual representations.

The model is pre-trained with a mask-and-predict objective. Given a text-image pair $[T, I]$ from the aligned dataset D , we randomly mask out some words w_i , some regions r_j , and some tags d_k to obtain masked $[\tilde{T}, \tilde{I}]$. The model is trained to predict the masked words, the properties of the masked regions, and the masked tags given $[\tilde{T}, \tilde{I}]$. The pre-training objective can be summarized as:

$$\min_{\theta} \sum_{[T, I] \in D} L_{T+I+M} \left(f_{\theta}([\tilde{T}, \tilde{I}]), [T, I] \right).$$

f_{θ} represents the embedding layer and the multi-layer Transformer. L_{T+I+M} is the sum of 1) the masked language model loss L_T , 2) the image reconstruction loss L_I , and 3) an “text-image match” objective L_M . Specifically, L_I includes a *tag reconstruction* loss L_I^{tag} (more details in Section 3.2) and the two visual losses as in LXMERT (Tan and Bansal, 2019): the *region feature regression* loss L_I^{ref} , which forces the model to regress to the visual vector, and the *noisy label classification* loss L_I^{cls} , which predicts the detected labels of masked objects with the cross-entropy loss. With a probability of 0.5, we provide the model with a mismatched text-image pair instead of a matched pair, and L_M asks the model to predict whether the image matches the text. After the model is pre-trained, it can be fine-tuned for V&L tasks similar to how BERT is fine-tuned for NLP tasks.

3.2 Unsupervised Pre-training

We introduce the two core design choices of unsupervised pre-training: mask-and-predict pre-training with unaligned data and the detector tags.

Mask-and-Predict Pre-training with Unaligned

Data We assume access to a text corpus D_T and an image corpus D_I for pre-training. During every pre-training step, we randomly sample either a batch of text from D_T or a batch of images from D_I . No alignment between text and images is provided to the model. When pre-training with a text segment T , the model is trained to reconstruct T given the masked \tilde{T} .³ When pre-training with an image I , the model is trained to reconstruct I given the masked \tilde{I} . A single Transformer is used throughout two modalities (i.e. θ shared across modalities). The pre-training objective can be summarized as:

$$\min_{\theta} \sum_{T \in D_T} L_T(f_{\theta}(\tilde{T}), T) + \sum_{I \in D_I} L_I(f_{\theta}(\tilde{I}), I).$$

After pre-training, the model is fine-tuned on downstream tasks just as its supervised counterpart, with the input being a text-image pair.

Detector Tags While mask-and-predict pre-training with unaligned data in itself achieves non-trivial performance (Section 5.2), we find it beneficial to provide noisy alignment signals in the form of the detector tags. When modeling an image I , for each region detected, we append the tag outputted by the object detector to the input. The detector (Ren et al., 2015) is pre-trained on a general object detection dataset (Krishna et al., 2017; Anderson et al., 2018) and the tags are essentially a bag of words that provide some noisy grounding signals to the model. During pre-training, we apply the mask-and-predict objective to the tags, which further encourages grounding.

We process the detector tags as a subword sequence $d_{1:l}$ with spatial coordinates.⁴ Every tag subword is embedded as the sum of its token embedding and a spatial coordinate embedding. The token embedding is the same as the token embedding used in text modeling, while the spatial coordinate embedding is the same as the coordinate embedding of the corresponding region. The coordinate embedding allows the model to distinguish tags from different regions.⁵ With the de-

³We adopt the next sentence prediction task in BERT when long documents are available.

⁴Each tag corresponds to a region. A tag could be split into multiple subwords, so the total length of the tag subword sequence l is equal to or larger than the number of regions m .

⁵This design differs from that of Oscar (Li et al., 2020b). Oscar does not add the coordinate embeddings to tags to encourage the fusion of tag and visual representations.

tector tags added, the image I is embedded as a sequence of image region features $r_{1:m}$ followed by a sequence of detector tag embeddings $d_{1:l}$, i.e. $I = [r_{1:m}; d_{1:l}]$. The tags are added during both pre-training and fine-tuning. Further, during pre-training, certain tag subwords are masked and the *tag reconstruction loss* L_I^{tag} supervises the model to predict the masked tags. The tags are predicted just as masked subwords are predicted in text modeling. The prediction softmax layer is shared between the tag and text subwords.

The parameters involved in modeling tags include the token embedding, the coordinate embedding, and the subword softmax embedding. These embedding parameters are shared across modalities and encourage the model to project text, visual, and tag representations into the same space (see Section 5.2 for an example). This resembles the design in multi-lingual language models, which use shared BPE embeddings and softmax weights across languages (Wu and Dredze, 2019).

4 Experiment

As the domain and quality of data may affect the model performance, the conventional practice in unsupervised learning is to use aligned corpora without providing alignments, allowing for controlled comparison with a supervised model. For example, unsupervised machine translation creates unaligned corpora by splitting up parallel corpora (Lample et al., 2018) while unsupervised image captioning (Gu et al., 2019) create unaligned corpus by shuffling images and captions from MSCOCO (Chen et al., 2015). Following prior work, we first conduct experiments by using Conceptual Captions (CC) (Sharma et al., 2018) as the source of images and text for both the supervised and unsupervised model. Later in Section 5.1, we show that our method is effective when the images and captions are collected independently and when no caption text is used.

U-VisualBERT The model is pre-trained with shuffled captions and images. At each training step, we sample either a batch of images or a batch of text. Following VL-BERT (Su et al., 2019), we find it beneficial to include BookCorpus (Zhu et al., 2015), a general-domain text corpus, during pre-training. In sum, U-VisualBERT is trained on 3M images from CC, 3M captions from CC, and 2.5M

text segments from BookCorpus⁶.

S-VisualBERT We introduce a Supervised VisualBERT (S-VisualBERT) trained with aligned data as introduced in Section 3.1. S-VisualBERT is pre-trained on 3M caption-image pairs from CC and 2.5M text segments from BookCorpus.

Compared Models Additionally, we list the performance of a **Base** VisualBERT that is initialized from BERT and does not undergo further pre-training. Previously reported supervised models that are trained on CC are also listed, including ViLBERT, VL-BERT, and UNITER. For UNITER, we include the version that is trained only on CC (UNITER_{cc})⁷. Although their network architectures differ from ours and cannot be directly compared, they jointly paint the picture of the performance we should expect by pre-training on CC. Models developed before BERT are listed as **Pre-BERT** (Gao et al. (2019) for VQA, Suhr et al. (2019) for NLVR², Lee et al. (2018) for Flickr30K, and Yu et al. (2018) for RefCOCO+).

Setup For all the VisualBERT variants introduced in the paper, we initialize them from BERT_{base} and pre-train for 10 epochs on their respective pre-training datasets with a batch size of 144. All models can be trained within 3 days on 4 V100s each with 16GB of memory. We use the Adam optimizer (Kingma and Ba, 2015) with a linear-decayed learning-rate schedule (Devlin et al., 2019a) and a peak learning rate at 6×10^{-5} . We conduct evaluations by fine-tuning on four downstream tasks: Visual Question Answering (VQA 2.0) (Goyal et al., 2017), Natural Language for Visual Reasoning (NLVR²) (Suhr et al., 2019), Image Retrieval (Flickr 30K) (Plummer et al., 2015), and Referring Expression (RefCOCO+) (Yu et al., 2016). We use a Faster R-CNN pre-trained on the Visual Genome dataset to extract region features (Anderson et al., 2018). For each task, we follow the recommended setting in previous works. For details, please refer to the appendix.

Results Table 1 summarizes the results. For each model, we list the type and amount of data used

⁶Our version of BookCorpus contains around 5M text segments with 64 words per segment. For computational reasons, we downsample the dataset such that during each epoch, the model observes only half of the text segments from BookCorpus. This downsampling is also done for the other VisualBERT variants.

⁷The results are from Appendix A.6 of Chen et al. (2020c).

Model	Aligned	Unaligned		VQA Test-Dev	NLVR ²		Flickr30K			RefCOCO+		
		Image	Text		Dev	Test-P	R@1	R@5	R@10	Dev	TestA	TestB
Pre-BERT	-	-	-	70.22	54.1	54.8	48.60	77.70	85.20	65.33	71.62	56.02
ViLBERT	3M	0	0	70.55	-	-	58.78	85.60	91.42	72.34	78.52	62.61
VL-BERT	3M	0	~50M	71.16	-	-	-	-	-	71.60	77.72	60.99
UNITER _{cc}	3M	0	0	71.22	-	-	-	-	-	72.49	79.36	63.65
S-VisualBERT	3M	0	2.5M	70.87 \pm .02	73.44 \pm .51	73.93 \pm .51	61.19 \pm .06	86.32 \pm .12	91.90 \pm .02	73.65 \pm .11	79.48 \pm .36	64.49 \pm .22
Base	0	0	0	69.26	68.40	68.65	42.86	73.62	83.28	70.66	77.06	61.43
U-VisualBERT	0	3M	5.5M	70.74 \pm .06	71.74 \pm .24	71.02 \pm .47	55.37 \pm .49	82.93 \pm .07	89.84 \pm .21	72.42 \pm .06	79.11 \pm .08	64.19 \pm .54

Table 1: Evaluation results on four V&L benchmarks. Our unsupervised model trained with unaligned data (U-VisualBERT) achieves close performance with a supervised model trained with aligned data (S-VisualBERT). U-VisualBERT also rivals with several supervised models such as ViLBERT on most metrics.

Model	Text		VQA Test-Dev	NLVR ²		Flickr30K			RefCOCO+		
	Caption	General		Dev	Test-P	R@1	R@5	R@10	Dev	TestA	TestB
Base	-	-	69.26	68.40	68.65	42.86	73.62	83.28	70.66	77.06	61.43
U-VisualBERT	CC	BC	70.74	71.74	71.02	55.37	82.93	89.84	72.42	79.11	64.19
U-VisualBERT _{SBU}	SBU	BC	70.70	71.97	72.11	56.12	82.82	90.12	73.05	79.48	64.19
U-VisualBERT _{NC}	-	BC	70.47	71.47	71.19	54.36	82.22	89.24	72.96	79.30	64.25

Table 2: Unsupervised pre-training is applicable when images and captions are collected independently (U-VisualBERT_{SBU}) or when no caption text is provided (U-VisualBERT_{NC}).

during pre-training.⁸ To control for randomness, we report the means and standard deviations of U-VisualBERT and S-VisualBERT across three runs.

U-VisualBERT outperforms the Base model on all benchmarks, while only lagging behind S-VisualBERT slightly on VQA, NLVR², and RefCOCO+. U-VisualBERT even surpasses or rivals with some supervised models (e.g., ViLBERT on VQA and RefCOCO+, VL-BERT on RefCOCO+, and UNITER_{cc} on RefCOCO+). This shows that a model through unsupervised pre-training can perform comparably with supervised models.

On Flickr30K Image Retrieval, the difference between U-VisualBERT and S-VisualBERT is more evident. The task focuses on identifying if an image and a text segment are coherent. S-VisualBERT is provided with explicit signals for such a task with the “text-image match” objective L_M during pre-training (Section 3.1). While U-VisualBERT is not provided with such explicit signals, it still performs better than the Base model. Further, if we were to remove the explicit signal (i.e. the “text-image match” objective) when pre-training on aligned data, S-VisualBERT without L_M achieves only 57.98 on R@1, much closer to U-VisualBERT

⁸For models initialized from BERT, we do not count the BERT pre-training data. VL-BERT uses both BookCorpus and Wikipedia during V&L pre-training. We estimate that the two corpora roughly have 50M segments with 64 words per segment. With a different pre-processing style (e.g. longer segments), the number of segments may change.

5 Analysis

In this section, we analyze the effect of the text data and the role of the detector tags.

5.1 The Effect of Text Data

The assumption behind unsupervised pre-training is that the detector tags should appear both in the images and text corpus, serving as the grounding anchor points. When the images and captions come from the same corpus, such an assumption clearly holds, and unsupervised pre-training works well (Section 4). However, we are curious if such an assumption still holds 1) if images and captions come from independently collected corpora (U-VisualBERT_{SBU}) and 2) if no caption text but general-domain text is provided (U-VisualBERT_{NC}).

The latter setting bears great practical value. Conceptually, collecting caption-style text could be as hard as collecting image-caption data as images and captions seldom appear separately. It is desirable to explore training V&L representations without caption-style text. Thus we experiment pre-training with general-domain text, which could be easier to collect.

U-VisualBERT_{SBU} We use 3M images from CC and 1M captions from SBU captions (Ordonez et al., 2011). To compensate for the different amounts of text between CC and SBU, we upsam-

Model	VQA	NLVR ²		R@1	Flickr30K		RefCOCO+		
	Test-Dev	Dev	Test-P		R@5	R@10	Dev	TestA	TestB
Base _{NT}	69.06	51.98	52.73	48.40	78.20	87.18	70.15	76.91	61.72
U-VisualBERT _{NT}	69.87	67.90	68.92	50.56	80.22	88.32	71.94	77.79	62.38
U-VisualBERT	70.74	71.74	71.02	55.37	82.93	89.84	72.42	79.11	64.19
S-VisualBERT _{NT}	70.49	72.56	73.53	60.26	85.58	91.64	72.70	77.93	62.99
S-VisualBERT	70.87	73.44	73.93	61.19	86.32	91.90	73.65	79.48	64.49
H-VisualBERT	71.05 \pm .02	73.80 \pm .26	74.82 \pm .25	60.28 \pm .60	86.30 \pm .35	92.06 \pm .28	74.01 \pm .25	80.18 \pm .23	64.89 \pm .24

Table 3: Detector tags show a larger impact in the unsupervised setting (U-VisualBERT_{NT} vs. U-VisualBERT) than in the supervised setting (S-VisualBERT_{NT} vs. S-VisualBERT). Semi-supervised pre-training (H-VisualBERT) shows marginal improvement over supervised pre-training (S-VisualBERT).

ple the BookCorpus so that the amount of text data used by U-VisualBERT_{SBU} is roughly the same as U-VisualBERT.

U-VisualBERT_{NC} The model is trained on images from CC and text from BookCorpus, a general-domain corpus.

Results Unsupervised pre-training is effective in both scenarios (Table 1). When pre-training images and text are collected independently, U-VisualBERT_{SBU} achieves similar performance as U-VisualBERT, with the latter higher on VQA, and the former higher on the other three tasks.

When no caption text is used, the performance on NLVR² and RefCOCO+ remains unaffected while the performance on VQA and Flickr30K drops slightly, potentially because the language style of VQA and Flickr30K is similar to captions, benefiting U-VisualBERT. Such results are not surprising. In general-domain corpora like Wikipedia, grounded words take up a decent portion (>25%) (Tan and Bansal, 2020). Thus the tags appear in pre-training text corpora with a non-trivial frequency and U-VisualBERT_{NC} learns from such signals. The above results suggest the applicability of unsupervised pre-training to many language-only and image-only datasets, which are easier to collect than image-caption datasets (Trinh and Le, 2018; Sun et al., 2017).

5.2 The Detector Tags as Anchor Points

We study the effect of the detector tags in unsupervised and supervised pre-training, respectively.

W-VisualBERT_{NT} U-VisualBERT_{NT} observes no tags and only dense region features for image embeddings during pre-training and fine-tuning. For comparison, a base model without tags is introduced (Base_{NT}), which is initialized from BERT

and does undergo further pre-training.

S-VisualBERT_{NT} To study the effect of the detector tags when aligned data are present, we introduce S-VisualBERT_{NT} which is trained on aligned data but observes no tags for image embeddings.

Result We first find that even without tags, unsupervised pre-training benefits downstream tasks (Table 3). U-VisualBERT_{NT} outperforms Base_{NT} on all metrics with a large margin. We attribute this to the (unaligned) contextual V&L representation learned through pre-training. This bears resemblance to the observation in multi-lingual language models that the shared vocabulary across languages (i.e. anchor points) is not necessary for cross-lingual transfer (Conneau et al., 2020).

Further, while the detector tags are beneficial for both supervised and unsupervised pre-training, the performance improvement is more evident for the latter. For example, performance difference on VQA between U-VisualBERT and U-VisualBERT_{NT} is 0.95 (70.82 vs. 69.87) while the difference between S-VisualBERT and S-VisualBERT_{NT} is 0.41 (70.90 vs. 70.49). The results are expected. When aligned data are present, object tags serve as additional signals while in unsupervised pre-training, they serve as the only source from which grounding is learned.

Visualization To gain a direct sense of how the detector tags help bridge the modalities, we visualize the contextual representation spaces of S-VisualBERT, U-VisualBERT, and U-VisualBERT_{NT} in Figure 2. For each of the most frequent 15 object classes in the COCO dataset (Chen et al., 2015), we randomly sample at most 50 instances and take the last-layer contextual representations of the words, the objects, and the tags (when available) and visualize them with t-SNE

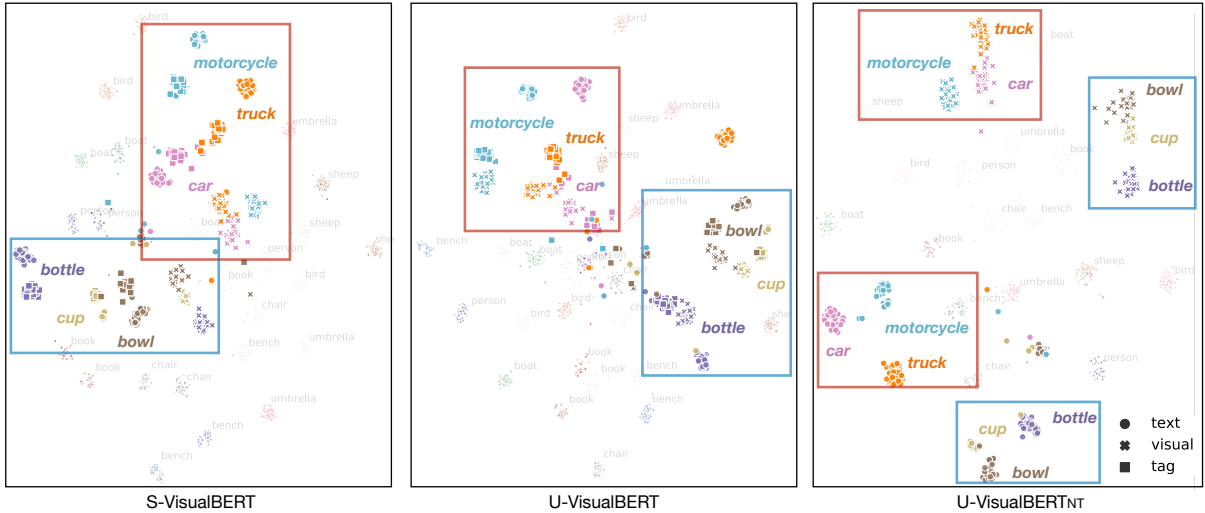


Figure 2: Visualization of the contextual representations of S-VisualBERT, U-VisualBERT, and U-VisualBERT_{NT}. The tags help to fuse text and visual representations for S-VisualBERT and U-VisualBERT. In U-VisualBERT_{NT}, common structures emerge in the text and visual representation spaces even though they are not aligned.

(Maaten and Hinton, 2008). We highlight the representations of six selected classes.

Though trained without aligned data, U-VisualBERT can group text, tag, and visual representations by their semantic classes. Similar phenomena can be observed in S-VisualBERT. U-VisualBERT_{NT}, lacking any signal to align the two spaces, does not show signs of such behaviour. In U-VisualBERT_{NT}, text and visual representations are almost completely separated (e.g., the two disjoint red rectangles in the figure on the right). However, some common structures emerge in both modalities. For instance, representations for “car”, “truck”, and “motorcycle”, the three semantically-related classes, are close to each other, in both the textual and visual modality (the red rectangles); representations for “cup”, “bottle”, and “bowl” are close (the blue rectangles). This also holds for the other two models and resembles what is observed in Li et al. (2020b) and Ilharco et al. (2020).

6 Semi-Supervised Pre-Training

Unsupervised pre-training in itself has great practical and research value in many domains where aligned data is scarce. As a byproduct, we wonder if the approach could find its use in a semi-supervised setting, where we pre-train a model with both aligned data and unaligned data.

H-VisualBERT We introduce a *hybrid* model that is trained on the 3M aligned data from Conceptual Captions (CC) and additional unaligned 1.7M images from Open Images (OI) (Kuznetsova et al.,

2020). When a training sample comes from CC, we provide the model with a text-image pair, and when the training sample comes from OI, we provide only the image. We do not use any manually annotated visual labels provided in OI.

Result We control for randomness by running H-VisualBERT for three times and report the means and stand deviations. We observe that H-VisualBERT brings consistent improvement upon S-VisualBERT on most tasks (Table 3) except Flickr30K⁹. This preliminary result is promising as the dataset scale in this experiment is relatively small (million-scale). Meanwhile, unannotated data generally could not improve upon a model trained with annotated data significantly, unless drastically scaled up (He et al., 2020). We leave large-scale experiments to future work.

7 Conclusion

In this paper, we explore unsupervised pre-training with unaligned data. We conduct mask-and-predict pre-training on textual data and visual data and the detector tags are used as anchor points to bridge the two modalities. Experiments show that unsupervised pre-training can achieve performance similar to supervised pre-training.

⁹On Flickr30K, the performance between H-VisualBERT and S-VisualBERT is similar, potentially because the “image-text match” objective is the dominant contributor and additional image-only data during pre-training have limited benefit (Section 4).

Ethical Considerations

One caveat of the proposed method is that data collected from the web may contain biases (Zhao et al., 2017), toxic contents (Schmidt and Wiegand, 2017), and other ethical issues. This problem is common to ML models and we stress that de-biasing (Zhao et al., 2019) and a rigorous examination are needed before deploying the system.

Acknowledgement

We would like to thank Hao Tan, members of UCLA NLP, and members of UCLA PlusLab for their helpful comments. We also thank the reviewers for the valuable reviews. This work was supported in part by DARPA MCS program under Cooperative Agreement N66001-19-2-4032. The views expressed are those of the authors and do not reflect the official policy or position of the Department of Defense or the U.S. Government.

References

- Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. 2018. Bottom-up and top-down attention for image captioning and visual question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Jize Cao, Zhe Gan, Yu Cheng, Licheng Yu, Yen-Chun Chen, and Jingjing Liu. 2020. Behind the scene: Revealing the secrets of pre-trained vision-and-language models. *Proceedings of the European Conference on Computer Vision (ECCV)*.
- Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. 2021. Conceptual 12m: Pushing web-scale image-text pre-training to recognize long-tail visual concepts. *arXiv preprint arXiv:2102.08981*.
- Liquan Chen, Zhe Gan, Yu Cheng, Linjie Li, Lawrence Carin, and Jingjing Liu. 2020a. Graph optimal transport for cross-domain alignment. *Proceedings of the International Conference on Machine Learning (ICML)*.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020b. A simple framework for contrastive learning of visual representations. *arXiv preprint arXiv:2002.05709*.
- Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C Lawrence Zitnick. 2015. Microsoft COCO captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325*.
- Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. 2020c. UNITER: Universal image-text representation learning. *ECCV*.
- Alexis Conneau, Shijie Wu, Haoran Li, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Emerging cross-lingual structure in pretrained language models. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019a. BERT: Pre-training of deep bidirectional transformers for language understanding. In *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)*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019b. Multilingual BERT readme document. <https://github.com/google-research/bert/blob/master/multilingual.md>.
- Carl Doersch, Abhinav Gupta, and Alexei A Efros. 2015. Unsupervised visual representation learning by context prediction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- 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. 2020. An image is worth 16x16 words: Transformers for image recognition at scale.
- Yang Feng, Lin Ma, Wei Liu, and Jiebo Luo. 2019. Unsupervised image captioning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Zhe Gan, Yen-Chun Chen, Linjie Li, Chen Zhu, Yu Cheng, and Jingjing Liu. 2020. Large-scale adversarial training for vision-and-language representation learning. *arXiv preprint arXiv:2006.06195*.
- Peng Gao, Zhengkai Jiang, Haoxuan You, Pan Lu, Steven CH Hoi, Xiaogang Wang, and Hongsheng Li. 2019. Dynamic fusion with intra-and inter-modality attention flow for visual question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. 2017. 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 (CVPR)*.
- Jiuxiang Gu, Shafiq Joty, Jianfei Cai, Handong Zhao, Xu Yang, and Gang Wang. 2019. Unpaired image captioning via scene graph alignments. In *Proceedings of the International Conference on Computer Vision (ICCV)*.

- Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. 2020. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Xiaowei Hu, Xi Yin, Kevin Lin, Lijuan Wang, Lei Zhang, Jianfeng Gao, and Zicheng Liu. 2020. Vivo: Surpassing human performance in novel object captioning with visual vocabulary pre-training. *arXiv preprint arXiv:2009.13682*.
- Zhicheng Huang, Zhaoyang Zeng, Bei Liu, Dongmei Fu, and Jianlong Fu. 2020. Pixel-bert: Aligning image pixels with text by deep multi-modal transformers. *arXiv preprint arXiv:2004.00849*.
- Gabriel Ilharco, Rowan Zellers, Ali Farhadi, and Hananeh Hajishirzi. 2020. Probing text models for common ground with visual representations. *arXiv preprint arXiv:2005.00619*.
- Andrej Karpathy and Li Fei-Fei. 2015. Deep visual-semantic alignments for generating image descriptions. In *CVPR*.
- Douwe Kiela, Hamed Firooz, Aravind Mohan, Vedanuj Goswami, Amanpreet Singh, Pratik Ringshia, and Davide Testuggine. 2020. The hateful memes challenge: Detecting hate speech in multimodal memes. *arXiv preprint arXiv:2005.04790*.
- Diederik P Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. *International Conference on Learning Representations (ICLR)*.
- Alexander Kolesnikov, Xiaohua Zhai, and Lucas Beyer. 2019. Revisiting self-supervised visual representation learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. 2017. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International Journal of Computer Vision (IJCV)*.
- Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Mallocci, Alexander Kolesnikov, Tom Duerig, and Vittorio Ferrari. 2020. The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. In *International Journal of Computer Vision (IJCV)*.
- Iro Laina, Christian Rupprecht, and Nassir Navab. 2019. Towards unsupervised image captioning with shared multimodal embeddings. In *Proceedings of the International Conference on Computer Vision (ICCV)*.
- Guillaume Lample, Alexis Conneau, Ludovic Denoyer, and Marc’ Aurelio Ranzato. 2018. Unsupervised machine translation using monolingual corpora only. In *International Conference on Learning Representations (ICLR)*.
- Kuang-Huei Lee, Xi Chen, Gang Hua, Houdong Hu, and Xiaodong He. 2018. Stacked cross attention for image-text matching. In *Proceedings of the European Conference on Computer Vision (ECCV)*.
- Gen Li, N. Duan, Yuejian Fang, Daxin Jiang, and M. Zhou. 2020a. Unicoder-vl: A universal encoder for vision and language by cross-modal pre-training. In *Proceedings of the National Conference on Artificial Intelligence (AAAI)*.
- Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. 2019. Visualbert: A simple and performant baseline for vision and language.
- Xiujun Li, Xi Yin, Chunyuan Li, Pengchuan Zhang, Xiaowei Hu, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, et al. 2020b. Oscar: Object-semantic aligned pre-training for vision-language tasks. In *European Conference on Computer Vision*, pages 121–137. Springer.
- Yikuan Li, Hanyin Wang, and Yuan Luo. 2020c. A comparison of pre-trained vision-and-language models for multimodal representation learning across medical images and reports. In *2020 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*.
- Nelson F. Liu, Matt Gardner, Yonatan Belinkov, Matthew E. Peters, and Noah A. Smith. 2019. Linguistic knowledge and transferability of contextual representations. In *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)*.
- Jiasen Lu, Batra Dhruv, Parikh Devi, and Lee Lee. 2019. ViLBERT: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. *arXiv preprint arXiv:1908.02265*.
- Laurens van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-sne. *Journal of machine learning research*.
- Mehdi Noroozi and Paolo Favaro. 2016. Unsupervised learning of visual representations by solving jigsaw puzzles. In *Proceedings of the European Conference on Computer Vision (ECCV)*.
- Vicente Ordonez, Girish Kulkarni, and Tamara L. Berg. 2011. Im2text: Describing images using 1 million captioned photographs. In *Proceedings of the Conference on Advances in Neural Information Processing Systems (NeurIPS)*.

- Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A Efros. 2016. Context encoders: Feature learning by inpainting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Matthew Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*.
- Telmo Pires, Eva Schlinger, and Dan Garrette. 2019. How multilingual is multilingual BERT? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*.
- Bryan A Plummer, Liwei Wang, Chris M Cervantes, Juan C Caicedo, Julia Hockenmaier, and Svetlana Lazebnik. 2015. Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. In *ICCV*.
- Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. 2015. Faster R-CNN: Towards real-time object detection with region proposal networks. In *NeurIPS*.
- Anna Rohrbach, Marcus Rohrbach, Ronghang Hu, Trevor Darrell, and Bernt Schiele. 2016. Grounding of textual phrases in images by reconstruction. In *Proceedings of the European Conference on Computer Vision (ECCV)*.
- Anna Schmidt and Michael Wiegand. 2017. A survey on hate speech detection using natural language processing. In *Proceedings of the Fifth International Workshop on Natural Language Processing for Social Media*.
- Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. 2018. 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)*.
- Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. 2019. VL-BERT: Pre-training of generic visual-linguistic representations. *arXiv preprint arXiv:1908.08530*.
- Alane Suhr, Stephanie Zhou, Iris Zhang, Huajun Bai, and Yoav Artzi. 2019. A corpus for reasoning about natural language grounded in photographs. *ACL*.
- Chen Sun, Abhinav Shrivastava, Saurabh Singh, and Abhinav Gupta. 2017. Revisiting unreasonable effectiveness of data in deep learning era. In *Proceedings of the International Conference on Computer Vision (ICCV)*.
- Hao Tan and Mohit Bansal. 2019. LXMERT: Learning cross-modality encoder representations from transformers. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*.
- Hao Tan and Mohit Bansal. 2020. Vokenization: Improving language understanding with contextualized, visual-grounded supervision. In *Proc. of the Conference on Empirical Methods for Natural Language Processing (EMNLP)*.
- Trieu H Trinh and Quoc V Le. 2018. A simple method for commonsense reasoning. *arXiv preprint arXiv:1806.02847*.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *NeurIPS*.
- Qinxin Wang, Hao Tan, Sheng Shen, Michael Mahoney, and Zhewei Yao. 2020. MAF: Multimodal alignment framework for weakly-supervised phrase grounding. In *Proc. of the Conference on Empirical Methods for Natural Language Processing (EMNLP)*.
- Shijie Wu and Mark Dredze. 2019. Beto, bentz, becas: The surprising cross-lingual effectiveness of BERT. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*.
- Fanyi Xiao, Leonid Sigal, and Yong Jae Lee. 2017. Weakly-supervised visual grounding of phrases with linguistic structures. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Fei Yu, Jiji Tang, Weichong Yin, Yu Sun, Hao Tian, Hua Wu, and Haifeng Wang. 2020. Ernievil: Knowledge enhanced vision-language representations through scene graph. *arXiv preprint arXiv:2006.16934*.
- Licheng Yu, Zhe Lin, Xiaohui Shen, Jimei Yang, Xin Lu, Mohit Bansal, and Tamara L Berg. 2018. Mattnet: Modular attention network for referring expression comprehension. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Licheng Yu, Patrick Poirson, Shan Yang, Alexander C Berg, and Tamara L Berg. 2016. Modeling context in referring expressions. In *ECCV*.
- Jieyu Zhao, Tianlu Wang, Mark Yatskar, Ryan Cotterell, Vicente Ordonez, and Kai-Wei Chang. 2019. Gender bias in contextualized word embeddings. In *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)*.

Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2017. Men also like shopping: Reducing gender bias amplification using corpus-level constraints. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*.

Luowei Zhou, Hamid Palangi, Lei Zhang, Houdong Hu, Jason J Corso, and Jianfeng Gao. 2020. Unified vision-language pre-training for image captioning and vqa. In *Proceedings of the National Conference on Artificial Intelligence (AAAI)*.

Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. In *Proceedings of the International Conference on Computer Vision (ICCV)*.

A Fine-Tuning on Downstream Tasks

We describe the details of fine-tuning on the four downstream tasks: Visual Question Answering (VQA 2.0) (Goyal et al., 2017), Natural Language for Visual Reasoning (NLVR²) (Suhr et al., 2019), Image Retrieval (Flickr 30K) (Plummer et al., 2015), and Referring Expression (RefCOCO+) (Yu et al., 2016).

VQA Given an image and a question, the task is to correctly answer the question. We use the VQA 2.0 and use the Karpathy split for training and validation (Karpathy and Fei-Fei, 2015). We fine-tune with a binary cross-entropy loss. The model is trained with a batch size of 32 and a peak learning rate of 5×10^{-5} over 8 epochs.

NLVR² NLVR² involves determining whether a natural language caption is true about a pair of images. While more sophisticated fine-tuning strategy exists (Chen et al., 2020c), we follow LXMERT (Tan and Bansal, 2019) to pair the caption with each image, concatenate the “[CLS]” representation of the two pairs, and build a classifier on top. We find it beneficial to conduct a moderate amount of “task-specific pre-training” where we use the data from the dataset to conduct mask-and-predict pre-training as suggested by VisualBERT (Li et al., 2019). We conduct task-specific pre-training for at most 5 epochs and fine-tune from the epoch with the best validation LM loss. Fine-tuning is conducted for 8 epochs with a batch size of 32 and a peak learning rate of 2×10^{-5} .

Flickr30K The task of image retrieval involves finding the corresponding image from a collection of images given a caption. We follow the split of

Lee et al. (2018) and use 1,000 images for validation and test each and train on the rest of the dataset. During fine-tuning, we follow UNITER (Chen et al., 2020c) and sample two negative text-image pairs along with a positive sample. We train for 5K steps with a batch size of 8 and a peak learning rate of 5×10^{-5} .

RefCOCO+ The referring expression task involves locating an image region given a natural language phrase. We follow ViLBERT (Lu et al., 2019) and conduct evaluation on the RefCOCO+ dataset. We use the bounding box proposals provided by Yu et al. (2018). For each box proposal, the model is trained to classify if it matches the reference phrase or not. A proposal box is considered correct if it has an IoU with the gold box larger than 0.5. We train for 12 epochs with a batch size of 32 and a peak learning rate of 5×10^{-5} .

B Data Accessibility

The version of BookCorpus we used is downloaded from https://github.com/jackroos/VL-BERT/blob/master/data/PREPARE_DATA.md. The other datasets we used including Conceptual Captions, Open Images, VQA, NLVR², Flickr30K, and RefCOCO+ are publicly available.