

# Curriculum Learning for Data-Efficient Vision-Language Alignment

Tejas Srinivasan      Xiang Ren      Jesse Thomason  
 University of Southern California  
 tejas.srinivasan@usc.edu

## Abstract

Aligning image and text encoders from scratch using contrastive learning requires large amounts of paired image-text data. We alleviate this need by aligning individually pre-trained language and vision representation models using a much smaller amount of paired data, augmented with a curriculum learning algorithm to learn fine-grained vision-language alignments. TOnICS (Training with **O**ntology-**I**nformed **C**ontrastive **S**ampling) initially samples minibatches whose image-text pairs contain a wide variety of objects to learn object-level alignment, and progressively samples minibatches where all image-text pairs contain the same object to learn finer-grained contextual alignment. Aligning pre-trained BERT and ViNVL models to each other using TOnICS outperforms CLIP on downstream zero-shot image retrieval while using less than 1% as much training data.

## 1 Introduction

Aligned representations for language and vision—which encode texts and corresponding images in a common latent space—are necessary to perform effective cross-modal retrieval. CLIP (Radford et al., 2021) and ALIGN (Jia et al., 2021) train individual text and image encoders from scratch to produce aligned image-text representations. Their encoders demonstrate strong cross-modal alignment, evidenced by strong performance on zero-shot retrieval tasks. However, these models were trained on proprietary datasets of 400M and 1B image-text pairs respectively, on hundreds of GPUs and TPUs, which is infeasible for non-industry practitioners.

CLIP and ALIGN align their encoders using the contrastive InfoNCE objective (Oord et al., 2018), which seeks to maximize the mutual information between image and text representations. In the InfoNCE objective, the model must correctly identify the positive image-text pair from among a set of negatives formed by the other minibatch pairs.

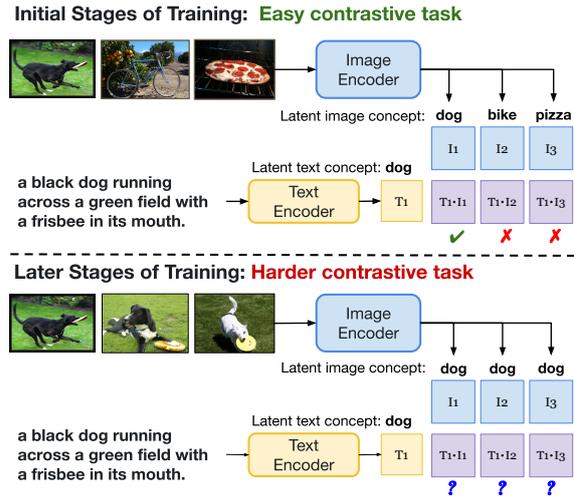


Figure 1: We propose TOnICS, a curriculum learning algorithm for contrastive alignment of language and vision encoders.

Since samples within a minibatch act as negative samples for each other in the InfoNCE objective, the minibatch determines the granularity of alignment that is learned. Minibatches constructed by random sampling contain a large variety of objects in the images and texts (Figure 1, top). To correctly match a *dog*-related caption to its image, it is sufficient to identify that the retrieved image must contain a dog, since the vast majority of randomly sampled negative images will not contain a dog. Thus, random minibatch sampling reduces the contrastive task to object-matching, for which object-level vision-language alignment suffices.

When minibatches are sampled such that the images contain the same objects, object-level alignments no longer suffice (Figure 1, bottom). The contrastive task can no longer be solved by identifying that the retrieved image must contain a dog, since all the negative images will also have a dog. The model must produce language and vision representations that encode shared *context*-level information, resulting in a finer-grained alignment.

In this work, rather than training our image and text encoders from scratch, we leverage rich single-modality pre-trained models—BERT (Devlin et al., 2019) for language, VinVL (Zhang et al., 2021)<sup>1</sup> for vision—and align them to each other using the InfoNCE contrastive objective. We perform the vision-language alignment using TOnICS, a novel ontology-based curriculum learning algorithm. TOnICS initiates training with an easy contrastive task by sampling minibatches randomly and progressively makes the contrastive task harder by constructing minibatches containing the same object class in the image and text inputs. We show that our learned representations have strong cross-modal alignment—outperforming CLIP on zero-shot Flickr30K image retrieval—while using less than 1% as much paired image-text training data.

## 2 Contrastive Vision-Language Alignment

We align language representations from BERT (Devlin et al., 2019) and visual representations from a VinVL object detector (Zhang et al., 2021). Our BERT-VinVL Aligner model is similar to the phrase grounding model from Gupta et al. (2020).

At every training step, the input to the model is a minibatch of  $N_B$  triplets, where each triplet  $X_i = \{t^i, v^i, w\}$  comes from an image-text pair. Each image caption  $t^i$  is encoded using BERT. The caption contains a noun  $w$ , whose word representation is denoted as  $h^i$ . For the corresponding image,  $v^i$  is a set of region features extracted from a frozen pre-trained VinVL object detector.<sup>2</sup> We add a learnable linear projection atop these region features.

In the cross-modal interaction, we employ a single Transformer (Vaswani et al., 2017) layer that uses  $i$ -th noun representation  $h^i$  as the query and  $j$ -th image features  $v^j$  as the keys and values. This layer outputs a visual representation  $v_{att}(i, j)$ , which is an attended representation of the  $j$ -th image, conditioned on the noun from the  $i$ -th caption. We then compute a dot product between the  $i$ -th noun representation  $h^i$  and the attended representation of  $j$ -th image  $v_{att}(i, j)$  to get an image-text score  $s(i, j) = \phi(h^i, v_{att}(i, j))$  (Figure 2).

To align the noun representation  $h^i$  to its corresponding image  $v^i$ , we use the InfoNCE loss (Oord et al., 2018) which maximizes a lower bound of

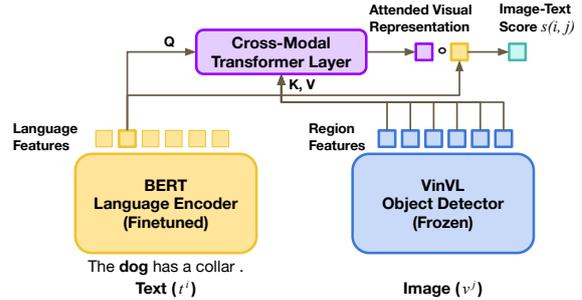


Figure 2: Our BERT-VinVL Aligner model scores every image-text combination  $(t^i, v^j)$  in the minibatch.

the mutual information between  $h^i$  and  $v_{att}(i, i)$ . InfoNCE minimizes the cross-entropy of correctly retrieving an image  $v^i$  from the set of all minibatch images, given the query noun representation  $h^i$ , with other instances in the minibatch acting as negative samples. We refer to the objective in this setup as the image retrieval loss,  $\mathcal{L}_{IR}$ :

$$\mathcal{L}_{IR}(i) = -\log \frac{\exp(s(i, i))}{\sum_{j=1}^{N_B} \exp(s(i, j))}$$

The training loss  $\mathcal{L}_{IR}$  is the mean loss  $\mathcal{L}_{IR}(i)$  over all images  $i = \{1 \dots N_B\}$  in the minibatch  $\mathcal{B}$ . We also similarly define a text retrieval loss,  $\mathcal{L}_{TR}$ , where the image  $v^i$  is used to retrieve the correct noun representation  $h^i$ :

$$\mathcal{L}_{TR}(i) = -\log \frac{\exp(s(i, i))}{\sum_{j=1}^{N_B} \exp(s(j, i))}$$

We experiment with training our model using just the image retrieval loss  $\mathcal{L}_{IR}$ , as well as the sum of the two losses  $\mathcal{L}_{IR} + \mathcal{L}_{TR}$ .

## 3 TOnICS: Training with Ontology Informed Contrastive Sampling

As noted above, negative samples for the contrastive learning objective come from other pairs in the minibatch. Therefore, the minibatch sampling itself influences the alignment learned by the model. We hypothesize that sampling minibatches randomly will yield object-level alignments, while sampling harder minibatches containing the same object in the image may result in fine-grained contextual alignments.

We introduce TOnICS, **T**rainig with **O**ntology-**I**nformed **C**ontrastive **S**ampling (Figure 3), a curriculum learning algorithm that initially seeks to align vision and language representations at the object level, and later learns contextual alignments.

<sup>1</sup>We use VinVL to refer to their pre-trained object detector.

<sup>2</sup>Region features provided at <https://github.com/pzzhang/VinVL/blob/main/DOWNLOAD.md>

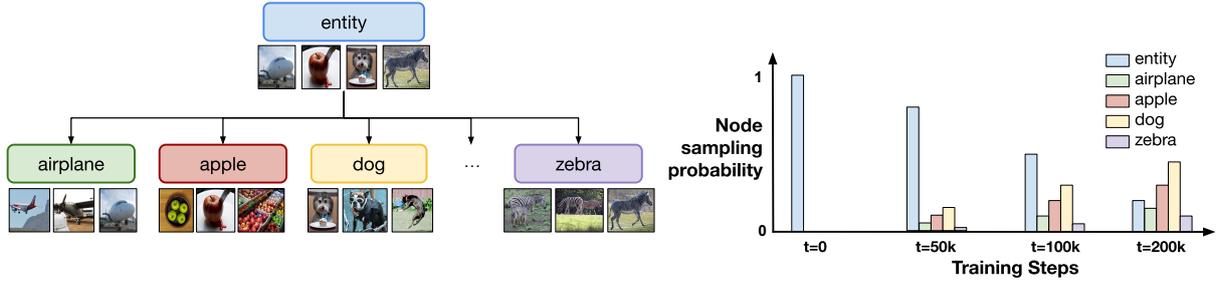


Figure 3: TOnICS selects image-text pairs for the minibatch by first sampling a node  $\eta$  from an ontology, according to a distribution  $P_S(\eta)$ . Sampling the root *entity* node yields easy minibatches containing pairs with a variety of objects, whereas sampling one of its children *object nodes* yields harder minibatches containing pairs sharing a common object, such as *apple* or *dog*, in a variety of contexts (left). TOnICS performs curriculum learning by moving node sampling mass away from the root *entity* node to the object nodes as training progresses (right).

TOnICS initiates the training by generating minibatches with randomly sampled image-text-noun triplets. As training progresses, TOnICS samples harder minibatches whose instances share the same object class in the image.

**Ontology Construction** We begin by extracting object detections from our training images using the pre-trained VinVL model. We next map each noun in the training data to an object class, wherever possible, resulting in a set of object classes  $\Theta$ . Every object class  $o \in \Theta$  has a corresponding set of nouns  $w(o)$ . For instance, the object class *dog*’s noun set  $w(o) = \{dog, dogs, puppy\}$ .

We construct the ontology (Figure 3, left), which contains an *entity* root node and its children *object nodes*  $\eta_o$ , each corresponding to an object class  $o$ . Every object node  $\eta_o$  has a corresponding set of triplet instances  $X(\eta_o)$ , a subset of the full training dataset whose triplet instances all contain the same object class  $o$  in the image, and all containing a noun from the noun set  $w(o)$  in the caption.

**TOnICS Minibatch Sampling** At every training step, TOnICS proceeds in two stages. First, a node  $\eta$  is sampled from the ontology, according to a sampling probability distribution  $P_S(\eta)$ . Second, we sample a minibatch according to the node that was just sampled. If we sample the entity node  $\eta_e$ , we sample the minibatch by sampling  $N_B$  instances from the full training data at random. If we sample an object node  $\eta_o$ , we sample  $N_B$  instances from the corresponding set  $X(\eta_o)$ , ensuring the minibatch contains images with the same object.

**TOnICS Curriculum Refresh** The curriculum is formed by varying the nodes’ sampling probability distribution throughout training. We initialize training by setting  $P_S(\eta_e) = 1$  and  $P_S(\eta_o) = 0$  for

all object nodes. After every fixed number of training steps, we evaluate the model’s image retrieval performance on a set of 100 held-out instances. If the held-out retrieval accuracy is greater than a certain threshold, we say that the model has learned the object-level alignment task, and we can start introducing harder minibatches in the training by *refreshing* the curriculum. The refresh step is performed by multiplying the entity node’s current sampling probability  $P_S(\eta_e)$  by a factor  $\alpha$ ;  $\alpha < 1$ . The remaining probability mass  $(1 - \alpha) \times P_S(\eta_e)$  is distributed among the object nodes. For each object node  $\eta_o$ , we update its sampling probability:

$$P_S(\eta_o) = P_S(\eta_o) + (1 - \alpha)P_S(\eta_e) \times \frac{|X(\eta_o)|}{\sum |X(\eta_o)|}$$

Object classes that are more common in the training data have more sampling probability mass distributed to their object node  $\eta_o$ , by weighting mass according to the size of the node’s instance set,  $|X(\eta_o)|$ . With each curriculum refresh, sampling mass is pushed down from the entity node to the object nodes, as long as  $P_S(\eta_e)$  does not fall below a fixed threshold  $\beta$ . Thresholding  $P_S(\eta_e)$  ensures the model still sees random minibatches and does not forget the initially learned object-level alignments.

## 4 Experiment Details

We train our BERT-VinVL model on MS-COCO and Conceptual Captions. We compare our model against CLIP on downstream retrieval tasks.

### 4.1 Training Data and Ontology

We train our model on image-text pairs from a combination of MS-COCO (Chen et al., 2015) and Conceptual Captions (Sharma et al., 2018). Our triplet instances only contain nouns which we wish

Model	# Image-Text Pairs	Minibatch Sampling Method	$\mathcal{L}_{TR}$	Zero-Shot Flickr30K				MS-COCO			
				Image Retrieval R@1	Image Retrieval R@5	Text Retrieval R@1	Text Retrieval R@5	Image Retrieval R@1	Image Retrieval R@5	Text Retrieval R@1	Text Retrieval R@5
CLIP-ViT-B/32	400M	Random	-	58.66	83.38	<b>79.2</b>	<b>95</b>	30.45	56.02	50.12	75.02
BERT-VinVL Aligner	2.84M	Random	✗	58.18	84.24	22.2	47.9	42.67	74.43	15.5	37.7
	2.84M	TOnICS	✗	<b>60.32</b>	85.14	24.4	49	47.94	77.38	16.1	35.1
	2.84M	Random	✓	58.9	84.6	76.1	93.3	42.74	74.37	59.84	86.46
	2.84M	TOnICS	✓	59.7	<b>85.24</b>	76.6	94.1	<b>48.26</b>	<b>77.87</b>	<b>65.44</b>	<b>89.36</b>

Table 1: Results of our BERT-VinVL Aligner model on image and text retrieval, compared to a CLIP model. Numbers in bold represent the best results among our model and CLIP.

to explicitly align with the visual modality. Each noun in the training data is initially mapped to the object class with maximum noun-object PMI, calculated over training pairs with object detections, and then adjusted by hand to correct erroneous mappings. Object classes containing fewer than 5000 instances in the training dataset are filtered out. This finally results in a set of 406 nouns, each noun corresponding to one of the 244 object categories  $\Theta$ . For every image-text pair in the original training dataset, we create one triplet for each noun in our set of 406 nouns that the text contains.

Our final training data consists of 5.8M triplet instances corresponding to 2.84M image-text pairs (2.26M from Conceptual Captions, 580K from MS-COCO) from 2.4M unique images. The ontology for TOnICS is constructed by creating an object node for each of the 244 object categories, which are children of the root *entity* node.

## 4.2 Implementation Details

We use pre-trained BERT-base as our text encoder. For our image encoder, we use VinVL, a pre-trained object detector that detects regions of interest (ROIs) in the image and outputs pooled CNN features for all ROIs. We use pre-extracted ROI features and treat the VinVL encoder as frozen, as we cannot backpropagate through the object detector.

All our models are trained for 500K iterations with a batch size of  $N_B = 256$ , yielding 255 negative pairs for every positive pair. Each model was trained on a single V100 GPU for 6 days, compared to CLIP which used 256 V100 GPUs for 12 days.

After every 5K iterations, we evaluate retrieval over a set of held-out instances and perform a curriculum refresh step if the held-out accuracy is at least 90%. When performing a refresh step, we retain  $\alpha = 90\%$  of *entity*'s sampling probability, so long as the probability does not fall below  $\beta = 0.2$ .

## 4.3 Baselines and Evaluation

To compare the effect of using pre-trained unimodal encoders at the start of the alignment process, we compare our model against CLIP (Radford et al., 2021). CLIP also uses separate image and text encoders, aligned using a contrastive loss with image-text data. Unlike our BERT-VinVL Aligner model, CLIP trains the two encoders from scratch, and uses significantly more paired image-text data—400M pairs, compared to our 2.84M pairs. Since we use the base variant of BERT, we compare against the CLIP-ViT-B/32 variant.<sup>3</sup> We do not compare against ALIGN as they have not released their base model checkpoint.

To evaluate the utility of our TOnICS algorithm, we also train our BERT-VinVL Aligner using a **Random** minibatch sampling baseline, where the minibatch instances are always randomly sampled throughout the training process.

We directly evaluate our trained Aligner model's (as well as pre-trained CLIP) on image and text retrieval. Specifically, we perform zero-shot retrieval on the test set of Flickr30K (Plummer et al., 2015), which contains 1,000 images. We also perform retrieval evaluation on the MS-COCO test set, which contains 5,000 images. This latter evaluation is not zero-shot since our training data contains images from the MS-COCO train set. We compare the Recall@1 and Recall@5 of all models.

## 5 Results and Discussion

In Table 1, we directly transfer both our trained BERT-VinVL Aligner model and pre-trained CLIP to the downstream task of image and text retrieval. Since our models are trained using retrieval objectives, we perform the retrieval evaluation using the same setup as training.

<sup>3</sup>Checkpoint provided at <https://huggingface.co/openai>

The Flickr30K evaluation is zero-shot for both CLIP and our BERT-VinVL Aligner model since neither model’s training data contains images from the Flickr30K train set. We see that even with the Random minibatch sampling and only the image retrieval loss,  $\mathcal{L}_{IR}$ , our BERT-VinVL Aligner achieves approximately the same image retrieval performance as CLIP. When the Aligner is trained with our TOnICS curriculum learning algorithm, we get a 1.5% improvement on R@1 over CLIP.

However, this model fails to do well at the text retrieval task. Adding the text retrieval loss  $\mathcal{L}_{TR}$  leads to substantial improvements in downstream text retrieval, with the Random baseline performing only 3% worse than CLIP. We further see that training with TOnICS leads to only slight improvements in Flickr30K text retrieval. Adding the text retrieval loss slightly hurts image retrieval performance, but still does better than CLIP by 1%.

Since our model, unlike CLIP, includes MS-COCO training images in the training data, it significantly outperforms CLIP on the MS-COCO retrieval evaluation. Hence, we compare our TOnICS algorithm to the Random baseline on the MS-COCO evaluation. We see that TOnICS leads to significant improvements in image retrieval (> 5%), both when the text contrastive loss is and isn’t used. We once again see that the text retrieval performance is very poor without the text retrieval objective during training, but improves significantly with it. TOnICS results in a 5% improvement over the Random baseline in text retrieval as well.

Minibatch sampling with TOnICS results in large gains in in-distribution retrieval evaluation (MS-COCO) as well as small improvements in zero-shot retrieval (Flickr30K). Training BERT-VinVL with TOnICS yields better zero-shot image retrieval performance than CLIP, even with substantially less training data.

## 6 Conclusions and Future Work

In this work, we align individually pre-trained language and vision encoders—BERT and VinVL, respectively—using a novel curriculum learning algorithm called TOnICS. Our aligned model is able to achieve better downstream zero-shot image retrieval performance than CLIP, in spite of being trained with less than 1% as many image-text training pairs. We further show that our TOnICS algorithm leads to gains in both in-domain and zero-shot retrieval tasks.

## References

- 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*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *North American Chapter of the Association for Computational Linguistics (NAACL)*.
- Tanmay Gupta, Arash Vahdat, Gal Chechik, Xiaodong Yang, Jan Kautz, and Derek Hoiem. 2020. Contrastive learning for weakly supervised phrase grounding. In *European Conference on Computer Vision (ECCV)*.
- Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. 2021. Scaling up visual and vision-language representation learning with noisy text supervision. In *International Conference on Machine Learning (ICML)*.
- Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2018. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*.
- 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 *International Conference on Computer Vision (ICCV)*.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. *arXiv preprint arXiv:2103.00020*.
- Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. 2018. Conceptual Captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Association for Computational Linguistics (ACL)*.
- 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 *Neural Information Processing Systems (NeurIPS)*.
- Pengchuan Zhang, Xiujun Li, Xiaowei Hu, Jianwei Yang, Lei Zhang, Lijuan Wang, Yejin Choi, and Jianfeng Gao. 2021. VinVL: Revisiting visual representations in vision-language models. In *Computer Vision and Pattern Recognition (CVPR)*.