Curriculum Learning for Data-Efficient Vision-Language Alignment

Anonymous ACL submission

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 with a curriculum learning algorithm to learn fine-grained vision-language alignments. TONICS (Training with Ontology-Informed Contrastive Sampling) 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 using <1% as much training data.

1 Introduction

002

011

013

014

017

020

021

034

Aligned representations for language and vision which encode texts and images in a common latent space—are necessary to perform effective crossmodal 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. They demonstrate accurate zero-shot retrieval due to strong cross-modal alignment. However, these models were trained on proprietary datasets of 400M and 1B image-text pairs 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.

Since samples within a minibatch act as negative samples for each other in the InfoNCE objective,

Image Encoder Latent image concept: dog bike pizza I1 I2 I3 Latent text concept: dog a black dog running Text across a green field with T1•I1 T1·I2 T1·I3 Encoder a frisbee in its mouth. Later Stages of Training: Harder contrastive task

Initial Stages of Training: Easy contrastive task



Figure 1: TONICS is a contrastive, curriculum learning algorithm for aligning language and vision encoders.

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. To correctly match a *dog*related caption to its image, it is sufficient to identify that the retrieved image must contain a dog, since most randomly sampled negative images will not contain a dog. Random minibatch sampling reduces the contrastive task to just object-matching.

043

044

046

047

050

051

054

060

061

062

063

When minibatches are sampled such that the images contain the same objects, object-level alignments no longer suffice (Figure 1). 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, we leverage rich single-modality pre-trained models—BERT (Devlin et al., 2019) for language, VinVL (Zhang et al., 2021)¹ for vision—and align them to each other using the In-

¹We use VinVL to refer to their pre-trained object detector.

foNCE contrastive objective. We propose TOnICS, a curriculum learning algorithm which initiates training with an easy contrastive task by sampling 066 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 alignmentoutperforming CLIP on zero-shot Flickr30K image retrieval-while using less than 1% as much paired image-text training data.

064

065

071

081

091

100

101

103

105

106

107

108

109

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).

During training, the input to the model is a minibatch of N_B triplets, where each triplet $X_i =$ $\{t^i, v^i, w\}$ represents an image-text pair. Image caption t^i is encoded using BERT, and contains a noun w with word representation h^i . A set of region features v^i are extracted from a frozen pretrained VinVL object detector.² 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 image v^i , we use the InfoNCE loss (Oord et al., 2018) to maximize the lower bound of 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 . We refer to the objective in this setup as the image retrieval loss, \mathcal{L}_{IB} :

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



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

The training loss \mathcal{L}_{IR} is the mean loss $\mathcal{L}_{IR}(i)$ over all minibatch instances $i = \{1...N_B\}$. We define a text retrieval loss, \mathcal{L}_{TR} , where the image v^i is used to retrieve the correct noun representation h^i : 113

$$\mathcal{L}_{TR}(i) = -\log \frac{\exp(s(i,i))}{\sum_{j=1}^{N_B} \exp(s(j,i))}$$
 114

110

111

112

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

140

141

142

143

144

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

TOnICS: Training with Ontology 3 **Informed Contrastive Sampling**

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 gives object-level alignments, while sampling harder minibatches containing the same object in the image gives fine-grained contextual alignments.

We introduce TOnICS, Training with Ontology-Informed Contrastive Sampling (Figure 3), a curriculum learning algorithm that first performs object level alignment via random minibatches, and later learns contextual alignments through minibatches with shared objects.

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 Θ . 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 η_o , each corresponding to an object class o.

²Region 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 η 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).

145 Every object node η_o has a corresponding set of 146 triplet instances $X(\eta_o)$, a subset of the full training 147 dataset whose triplet instances all contain the same 148 object class o in the image, and all containing a 149 noun from the noun set w(o) in the caption.

150

152

153

154

155

156

157

158

159

160

TONICS Minibatch Sampling At every training step, TONICS proceeds in two stages. First, a node η 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 η_e , we sample the minibatch by sampling N_B instances from the full training data at random. If we sample an object node η_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 probabil-162 ity distribution throughout training. We initialize 163 training by setting $P_S(\eta_e) = 1$ and $P_S(\eta_o) = 0$ 164 for all object nodes. After every fixed number of 165 training steps, we evaluate the model's image retrieval performance on a set of held-out instances. 167 If the held-out retrieval accuracy is greater than a 168 certain threshold, we start introducing harder mini-169 batches in the training by *refreshing* the curriculum. 170 The refresh step is performed by multiplying the 171 entity node's current sampling probability $P_S(\eta_e)$ 172 by a factor α ; $\alpha < 1$. The remaining probability 173 mass $(1 - \alpha) \times P_S(\eta_e)$ is distributed among the 174 175 object nodes. For each object node η_o , we update its sampling probability: 176

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

178Object classes that are more common in the train-179ing data have more sampling probability mass dis-

tributed to their object node η_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 β . Thresholding $P_S(\eta_e)$ ensures the model still sees random minibatches and does not forget the initially learned object-level alignments.

180

181

182

183

184

185

186

187

189

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

207

208

209

210

211

212

213

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 to explicitly align with the visual modality. We select a set of 406 nouns, each noun corresponding to one of the 244 object categories Θ (more details in Appendix A.1). Our final training data consists of 5.8M triplet instances corresponding to 2.84M image-text pairs 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 a text encoder and frozen VinVL, a pre-trained object detector returns pooled CNN features for all regions-ofinterest (ROIs), as an image encoder. We use preextracted ROI features, as we cannot backpropagate through the object detector. Further details can be found in Appendix A.2.

		Minibatch	Zero-Shot Flickr30K				MS-COCO				
	# Image-	Sampling	_	Image l	Retrieval	Text R	etrieval	Image I	Retrieval	Text R	etrieval
Model	Text Pairs	Method	\mathcal{L}_{TR}	R@1	R@5	R@1	R@5	R@1	R@5	R@1	R@5
CLIP-ViT-B/32	400M	Random	-	58.66	83.38	79.2	95	30.45	56.02	50.12	75.02
BERT-VinVL Aligner	2.84M	Random	X	58.18	84.24	22.2	47.9	42.67	74.43	15.5	37.7
	2.84M	TOnICS	X	60.04	84.72	18.8	43.1	47.68	77.14	11.48	27.3
	2.84M 2.84M	Random TOnICS	√ √	58.9 59.68	84.6 84.84	76.1 77.4	93.3 94	42.74 47.15	74.37 76.85	59.84 63.7	86.46 88.5

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.

4.3 Baselines and Evaluation

214

215

216

217

218

219

222

227

228

230

232

234

236

240

241

244

245

246

247

249

We compare our aligned model against CLIP (Radford et al., 2021). CLIP trains the image and text encoders from scratch and uses significantly more paired image-text data—400M pairs, compared to our 2.84M pairs. We use the base variant of BERT, and so compare against CLIP-ViT-B/32.^{3,4}

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 Aligner models and pretrained CLIP on image and text retrieval, using the Recall@1 and Recall@5 metrics. Specifically, we evaluate zero-shot retrieval on the Flickr30K (Plummer et al., 2015) test set, which contains 1,000 images. We also perform retrieval evaluation on the MS-COCO test set, which contains 5,000 images. This evaluation is not zero-shot since our training data contains MS-COCO training images.

5 Results and Discussion

We directly transfer both our trained BERT-VinVL Aligner model and pre-trained CLIP to the downstream task of image and text retrieval (Table 1) using the same task formulation from training time.

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 a 1% improvement in Flickr30K text retrieval. Adding the text retrieval loss slightly hurts image retrieval performance, but still does better than CLIP by 1%.

250

251

252

253

254

255

256

257

259

260

261

262

263

264

265

267

268

269

270

271

272

273

274

275

276

277

279

281

282

283

285

287

Since our model 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 ($\approx 5\%$), 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 4% 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

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

4

³https://huggingface.co/openai

⁴We do not compare against ALIGN because the authors have not released their base model checkpoint.

289

309

310

311

312

313

314

315

316

317

318

319

320

327

329

330 331

333

336 337

7 Ethical Limitations

We use pre-trained language models that can only 290 consume English text, eliding the challenges of 291 multi-lingual language-vision alignment. Further, images in our training data (MS-COCO and Conceptual Captions) are sourced from social media, movie clips and web searches, thus excluding certain image domains, including those taken for ac-296 cessibility needs such as descriptions for people with blindness (Gurari et al., 2018). Such biases 298 in our aligned model, inherited from the datasets and pre-trained models selected, serve the needs of English-speaking, able-bodied folks as a "default." 301 Further, pre-trained language models such as BERT 302 have been known to express gender and race biases. These biases have been shown to compound when 304 multiple modalities are represented (Srinivasan and Bisk, 2021). Our work does not contain analysis of how biases in our aligned BERT model differ from 307 pre-trained BERT.

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)*.
- Danna Gurari, Qing Li, Abigale J Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo, and Jeffrey P Bigham. 2018. Vizwiz grand challenge: Answering visual questions from blind people. In *Computer Vision and Pattern Recognition (CVPR)*.
- 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 imageto-sentence models. In *International Conference on Computer Vision (ICCV)*. 339

341

342

344

345

347

348

349

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

- 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)*.
- Tejas Srinivasan and Yonatan Bisk. 2021. Worst of both worlds: Biases compound in pre-trained visionand-language models. In *Workshop on Gender Bias in Natural Language Processing (GeBNLP)*.
- 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)*.
- Tian Yun, Chen Sun, and Ellie Pavlick. 2021. Does vision-and-language pretraining improve lexical grounding? In *Empirical Methods in Natural Language Processing (EMNLP)*.
- 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)*.

A Implementation Details

374

375

379

380

384

390

391

396

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

A.1 Which Nouns do we Align?

The triplet instances in our training data only contain nouns which we wish to explicitly align with the visual modality. We decide this set of nouns in the following procedure. 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 Θ . 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.

A.2 Training Hyperparameters

All our models are trained for 500K iterations with a batch size of $N_B = 256$, yielding 255 negative pairs for every positive pair. We select the model checkpoint which has maximum Recall@1 on the Flickr30K validation set, evaluated after every 5K iterations.

After every 5K iterations, we also evaluate retrieval over a set of 100 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 $\beta = 0.2$.

Each model was trained on a single V100 GPU for 6 days, compared to CLIP which used 256 V100 GPUs for 12 days.

B Analysis of Aligned Language Representations

We hypothesize that by aligning pre-trained BERT to visual representations from a pre-trained VinVL model, our aligned BERT's representations of visually-groundable objects will contain more visual context information. Similar to (Yun et al., 2021), we investigate whether noun representations extracted from our Aligned-BERT contain information about their visual attributes that are also described in the caption. Specifically, we look at representations of the word *shirt* in Flickr30K captions where the color of the shirt is also mentioned. We extract 275 such captions where the shirt is described as being one of ten colors, and extract the



Figure 4: TSNE projections of contextual representations of the word *shirt* occuring in different color contexts. Each dot corresponds to a contextual representations of the word *shirt*, where the color of the dot corresponds to the color of the shirt described in the caption (grey dots represent representations of white shirts). We compare the TSNE visualizations of pretrained BERT and the Aligned-BERT from our BERT-VinVL Aligner model.

Model	Homogeneity	Completeness	V-Score
BERT Aligned-BERT CLIP	$\begin{array}{c} 9.79 \pm 1.48 \\ 42.64 \pm 5.51 \\ 98.39 \pm 0.00 \end{array}$	$\begin{array}{c} 9.13 \pm 1.39 \\ 40.59 \pm 5.24 \\ 98.28 \pm 0.00 \end{array}$	$\begin{array}{c} 9.45 \pm 1.43 \\ 41.58 \pm 5.37 \\ 98.33 \pm 0.00 \end{array}$

Table 2: K-Means Clustering metrics (K=10) for *shirt* representations, across five different initializations. We present mean and standard deviation of all metrics, across the different templates.

word *shirt*'s contextual representations from both pre-trained BERT and our BERT-VinVL Aligner's text encoder, which we refer to as Aligned-BERT.

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

Figure 4 compares the TSNE visualizations of representations extracted from BERT and Aligned-BERT. We see clear clusters formed by representations of the same colored shirt in Aligned-BERT's visualization, whereas no such clusters exist in the BERT representations.

We also provide a quantitative analysis of the clustering in the representations, by performing K-Means clustering with K = 10. We evaluate the Homogeneity and Completeness of these clusters, which are equivalent to Set-Precision and Set-Recall respectively, as well as V-Score which is their harmonic mean. In Table 2, we see that Aligned-BERT's clusters are much more homogenous and complete than pre-trained BERT, but pre-trained CLIP's clusters are much better than both.