# ALIGN-VL: CAN BEING MODEST HELP IN THE ALIGNMENT OF VISION-LANGUAGE MODELS?

Anonymous authors

Paper under double-blind review

# ABSTRACT

011 Multimodal alignment aims to learn a shared latent space between different modal inputs to establish connections across modalities. A prime example is Visual Lan-012 guage Models (VLMs), such as CLIP, which benefit from extensive image-text 013 pre-training and excel in image recognition tasks. These models are emblematic 014 of successful multimodal alignment. Subsequent work has successfully aligned 015 multimodal data on limited datasets using feature mixing enhancement methods. 016 However, these models encounter significant challenges: The presence of am-017 biguous samples (either partially matched or completely unmatched) in datasets 018 with weakly associated, low-quality image-text pairs causes models to become 019 overconfident (in training) and confused (in inference), ultimately reducing performance. Current contrastive learning methods, which rely on single positive 021 pairs, exacerbate this issue by encouraging overconfidence when the model encounters such ambiguous samples. To overcome these challenges, we developed Align-VL, a multimodal alignment enhancement method that operates on the la-023 tent spaces of pre-trained unimodal encoders. This approach adjusts the matching degree of the data and moderates model overconfidence, promoting more appro-025 priate and effective alignments. Align-VL incorporates Random Perturbation and 026 *Embedding Smoothing* strategies to enhance input feature robustness and reduce model overconfidence, improving the model's ability to manage uncertainty and 028 generalize to new data. In our experiments, Align-VL outperformed existing state-029 of-the-art (SoTA) methods in image-text retrieval tasks, demonstrating its superior effectiveness. Align-VL also offers significant reductions in training time and data 031 requirements compared to methods like CLIP, using substantially fewer GPU days 032 and image-text pairs. Code will be publicly available.

034 035

004

010

# 1 INTRODUCITON

Multimodal learning, by integrating different types of data modalities, enhances a model's percep-037 tion and understanding capabilities, facilitating cross-modal information interaction and integration Radford et al. (2021); Vouitsis et al. (2024); Zhai et al. (2022); Liu et al. (2023b;a); Yang et al. (2021); Girdhar et al. (2023). Recent advancements in multimodal machine learning have shown 040 unprecedented potential across various application fields, with some applications even attracting 041 mainstream attention Girdhar et al. (2023); Radford et al. (2021). The cornerstone of multimodal 042 learning is multimodal alignment, which maps information from multiple modalities, such as text 043 and images, into a unified multimodal vector space Radford et al. (2021); Alayrac et al. (2022). 044 Researchers have made numerous efforts in multimodal alignment, with Visual Language Models 045 (VLMs) being particularly representative. VLMs like CLIP Radford et al. (2021), which undergo extensive image-text pre-training, excel in image recognition tasks, showcasing the potential of VLMs 046 in establishing effective cross-modal connections. 047

The success of multimodal alignment largely relies on large-scale training mechanisms like CLIP, which often require extensive GPU resources and rely on billions of multimodal data pairs Zhai
et al. (2022); Radford et al. (2021); Alayrac et al. (2022). However, the high computational costs are impractical for scenarios with limited computing resources or scarce multimodal data. Therefore, designing a cost-effective and efficient multimodal alignment framework is crucial. Inspired by Mixup Zhang et al. (2018), Fusemix introduces an efficient strategy for multimodal alignment Vouitsis et al. (2024) by augmenting the latent spaces of pre-trained unimodal encoders, allowing

for model creation with significantly reduced data and computational requirements. However, ambiguous samples—whether partially matched or completely unmatched—in datasets with weakly
associated (see Figure 4), low-quality image-text pairs can lead to overconfidence and confusion in
models, ultimately degrading performance. Moreover, current contrastive learning methods, which
rely on single positive examples, exacerbate this issue by further encouraging overconfidence in the
presence of ambiguous samples Radford et al. (2021); Vouitsis et al. (2024).

060 To overcome these issues, building on existing foundation Vouitsis et al. (2024), we propose 061 Align-VL, a multimodal alignment enhancement method designed to adjust the matching de-062 gree of the data and moderate model overconfidence, incorporating two key components: 1) 063 Random Perturbation: This introduces normally distributed perturbations at the visual-text fea-064 ture level to simulate uncertainty, enhancing the model's generalization capability and helping it learn more robust feature representations. 2) Embedding Smoothing: This aims to 065 smooth the model's prediction of output distributions, moderating model overconfidence in pos-066 itive samples and increasing the smoothness for predictions on uncertain samples, thereby en-067 hancing generalization. By using Align-VL to align the latent spaces of pre-trained uni-068 modal encoders, we have developed a highly competitive visual-language (V-L) model. In re-069 trieval tasks, this model not only surpasses existing state-of-the-art (SoTA) methods but also significantly reduces the need for computational resources and data, as detailed in Figure 1. 071 Our study makes two significant contributions:

072 073

074

075

076

077

078

079

081

082

084

085

087

090

091

092 093

- Statistical analysis reveals that the quality of existing image-text paired datasets is suboptimal, causing VLMs to become confused and overly confident when faced with ambiguous positive pairs (either partially matched (Figure 4 (a), (b), (c)) or completely unmatched (Figure 4) (d)). This significantly impacts the performance of multimodal alignment.
  - We propose a novel V-L alignment algorithm, Align-VL, which incorporates *Random Perturbation* to simulate input uncertainty and *Embedding Smoothing* to mitigate overconfidence in positive samples. This Align-VL enhances model generalization and robustness, effectively addressing the challenges posed by the suboptimal quality of existing datasets.



Figure 1: Image-to-Text retrieval performance on the Flickr30K test set Young et al. (2014) is plotted against the number of training pairs on a log-scale x-axis, illustrating how training volume impacts effectiveness.



Figure 2: Three Multimodal Alignment Paradigms: CLIP for Contrastive Learning, FuseMix for Enhanced Mixing Embeddings, and Align-VL for Moderating Model Overconfidence and Enhancing Robustness. In Fusemix and Align-VL, the text anchor and visual positive sample are derived from mixed features. In Align-VL, the visual and text positive pairs are the embedding augmented with random perturbation.

# 108 2 RELATED WORK

110

122

144 145

146

147

148

149

150

151

152

153

154

111 Multimodal alignment achieves cross-modal synchronization not through direct correspondences 112 between modalities, but implicitly via internal model mechanisms that discern latent semantic connections within the data. The primary objective of these models is to learn a shared latent space 113 capable of jointly encoding multiple modalities, thereby facilitating effective multimodal alignment 114 Tan & Bansal (2019); Li et al. (2020); Yuan et al. (2021); Wang et al. (2022a); Bao et al. (2022); 115 Wang et al. (2022b); Girdhar et al. (2022); Likhosherstov et al. (2023); Zhang et al. (2023); Wu 116 et al. (2023). Image-language alignment is a pivotal area of study in multimodal alignment, aiming 117 to create universal models capable of interpreting both image and language data. Standard multi-118 modal models usually undergo end-to-end training on image-text pairs. Yet, training these large-119 scale models from scratch demands substantial computational and data resources, which can restrict 120 scalability. Arandjelovic & Zisserman (2017); Lu et al. (2019); Sun et al. (2019); Su et al. (2020); 121 Chen et al. (2020); Li et al. (2021; 2022).

Pioneered by CLIP and ALIGN Radford et al. (2021); Jia et al. (2021), this approach uses a dual-123 encoder architecture, jointly embedding text and images into the same latent space through con-124 trastive target training. 3T aligns text and image encoders with the latent space of a pretrained 125 classifier Kossen et al. (2023). LiT uses a frozen pretrained image classifier as the image encoder 126 and aligns a text encoder to it Zhai et al. (2022). Although these methods have seen success, they 127 mostly train one or two encoders from scratch, relying on expensive cross-GPU gradient computa-128 tions. ImageBind Girdhar et al. (2023) uses images as anchors to learn a shared latent space across 129 six modalities through contrastive learning, jointly training various modality encoders from scratch. Moreover, the large-scale image-text paired datasets they use, ranging from 400 million to 5 billion 130 pairs, mostly sourced from the internet, are generally not public Vouitsis et al. (2024). In contrast to 131 these works, Fusemix boosts computational and data efficiency through feature augmentation tech-132 niques, using frozen pre-trained unimodal encoders and fewer multimodal paired data, requiring 133 fewer resources Vouitsis et al. (2024). However, ambiguous positive samples in weakly associated 134 datasets (see Figure 4) lead to model overconfidence and degraded performance, exacerbated by 135 contrastive learning methods that focus on single positive examples Radford et al. (2021). 136

Figure 2 compares three multimodal alignment paradigms: the CLIP Paradigm, which utilizes contrastive learning to manage data point relationships; the Fusemix Paradigm, which enhances embeddings by mixing features of image-text pairs; and the Align-VL Paradigm, which reduces overconfidence through perturbations and embedding smoothing, enhancing generalization and robustness (CLIP and FuseMix train on 3M data pairs, Align-VL uses about 3.5M pairs to achieve its results). The Align-VL specifically reduces model overconfidence and ensures experiments are computationally and data efficient, requiring only a reasonable amount of GPU resources (see Figure 6).



Figure 3: A pipeline of the Align-VL showcases the process of aligning the latent spaces of pre-trained unimodal encoders using a fewer dataset of paired data. The unimodal encoders remain frozen throughout the process, with their latent encodings pre-computed only once for efficiency. In this framework, both Random Perturbation and Embedding Smoothing are applied to each latent space to enhance robustness and reduce model overconfidence. Lightweight V-L adapters are then trained to meticulously align these augmented latents into a cohesive, shared latent space, effectively bridging the semantic gap between different modalities.

# <sup>162</sup> 3 METHODOLOGY

163 164

166

167 168 In this section, we introduce the Align-VL framework, designed to facilitate visual-text modal alignment in the latent space while addressing key considerations such as model overconfidence, and computational and data efficiency. Align-VL entire process is illustrated in Figure 3.

169 3.1 PRELIMINARIES

171 Notation: We define the task of V-L alignment from an alignment perspective. The goal is to learn a 172 shared latent space between visual and textual modal inputs. Formally, given any two data modalities 173 (images  $\mathcal{X}$  and text  $\mathcal{Y}$ ), our objective is to learn two networks,  $f_X : \mathcal{X} \to \mathcal{S}$  and  $f_Y : \mathcal{Y} \to \mathcal{S}$ , that 174 embed each modality into a shared latent space  $\mathcal{S}$ .

175 We take our two encoders as  $f_X = F_X \circ A_X$  and  $f_Y = G_Y \circ A_Y$ . That is, we define 176  $F_X(Frozen): \mathcal{X} \to \mathcal{U}_{\mathcal{X}}$  and  $G_Y(Frozen): \mathcal{Y} \to \mathcal{U}_{\mathcal{Y}}$ , where  $\mathcal{U}_{\mathcal{X}}$  and  $\mathcal{U}_{\mathcal{Y}}$  are intermediate la-177 tent spaces. We then have  $A_X(Learnable): \mathcal{U}_{\mathcal{X}} \to \mathcal{S}$  and  $A_Y(Learnable): \mathcal{U}_{\mathcal{Y}} \to \mathcal{S}$ , which 178 we hereafter refer to as V-L adapters. Our insight here is to take both  $F_X$  and  $G_Y$  as pre-trained 179 unimodal encoders which we keep frozen throughout, and treat our V-L adapters  $A_X$  and  $A_Y$  as 180 learnable heads for multimodal alignment. Therefore, we can define our learning objective using the 176 InfoNCE loss function as follows:

183

185

 $\mathcal{L}_{\text{NCE}} = -\log \frac{\exp(\sin(f_X(\mathcal{X}_i), f_Y(\mathcal{Y}_i))/\tau)}{\sum_{j=1}^N \exp(\sin(f_X(\mathcal{X}_i), f_Y(\mathcal{Y}_j))/\tau)}$ (1)

where  $sim(\cdot, \cdot)$  denotes the similarity function,  $\tau$  is the temperature parameter, and N is the number of all samples. Where  $\mathcal{X}_i$  and  $\mathcal{Y}_i$  are positive pairs.

188 Motivation: Scaling up multimodal models boosts performance but incurs substantial computa-189 tional costs, especially when jointly training networks  $F_X$  and  $G_Y$ , which rapidly increases mem-190 ory and compute demands. Besides, acquiring high-quality paired data is costly, while high-quality 191 unimodal data is more readily available and can provide rich supervision through self-supervised 192 learning. To address these challenges, we aim to design a computationally efficient, V-L alignment 193 adapters  $A_X$  and  $A_Y$ , that reduces reliance on paired data by leveraging unimodal signals, while 194 allowing independent updates to visual and textual components with minimal retraining.

Additionally, the existing image-text matching datasets, mostly sourced from the internet (Figure 4 Upper Figure), generally have low pairing quality. In contrastive learning that emphasize the quality of positive sample pairings, ambiguously matched positive pairs can lead to excessive model confidence, despite the pairs not being correctly matched. Therefore, there is a need for a method to adjust the matching degree of positive pairs and mitigate model overconfidence in these samples.

200 201

202

# 3.2 OVERVIEW

203 We introduce Align-VL, a multimodal alignment method that adjusts data matching by simulating 204 input uncertainty and moderating model overconfidence in ambiguous positive samples, operating 205 on the latent spaces  $\mathcal{U}_{\mathcal{X}}$  and  $\mathcal{U}_{\mathcal{Y}}$  derived from pre-trained unimodal encoders. 1) The method initially 206 employs unimodal encoders to encode V-L modalities into intermediate latent spaces. 2) Subse-207 quently, it utilizes enhanced features based on Fusemix in the Align-VL latent spaces, incorporating Random Perturbation to adjust data matching by simulating input uncertainty. 3) After training 208 through VL-Adapters, Embedding Smoothing is applied, aiming to smooth the model's prediction 209 of output distributions, reduce overconfidence in positive pairs, and enhance the smoothness of pre-210 dictions on uncertain samples. 4) Finally, the smoothed embeddings are used in contrastive learning 211 training, facilitating the learning of two networks,  $A_X$  and  $A_Y$ , as shown in Figure 3. 212

Align-VL utilizes the existing semantics encoded by unimodal encoders, reducing the reliance on
 extensive real paired data and simplifying computational requirements. It effectively mitigates the
 issue of model overconfidence, making the model more "modest" and robust, thereby optimizing
 multimodal alignment, learning efficiency, and generalization capabilities.

# 216 3.3 RANDOM PERTURBATION HELPS ENHANCE FEATURE ROBUSTNESS.

We introduce random perturbation in the Align-VL latent space to adjust data matching by simulating input uncertainty. Given a visual input  $\mathcal{X}$  and a textual input  $\mathcal{Y}$ , their respective embeddings  $\mathbf{z}_v = F_X(\mathcal{X})$  and  $\mathbf{z}_t = G_Y(\mathcal{Y})$  are obtained from the visual encoder  $F_X$  and the text encoder  $G_Y$ . In training phase, we apply Gaussian noise Willett et al. (2000) to these embeddings by sampling noise vectors  $\mathbf{n}_v \sim \mathcal{N}(0, \sigma^2)$  and  $\mathbf{n}_t \sim \mathcal{N}(0, \sigma^2)$ , resulting in perturbed embeddings. The perturbed embeddings are defined as:

224

240

$$\widetilde{\mathbf{z}}_{v} = \mathbf{z}_{v} + \sigma \cdot \epsilon_{v}, \quad \epsilon_{v} \sim \mathcal{N}(0, I) 
\widetilde{\mathbf{z}}_{t} = \mathbf{z}_{t} + \sigma \cdot \epsilon_{t}, \quad \epsilon_{t} \sim \mathcal{N}(0, I)$$
(2)

where  $\sigma$  represents the noise level, and  $\epsilon_v$  and  $\epsilon_t$  are random noise vectors sampled from the standard normal distribution. This noise is added only during the training phase to prevent the model from becoming overconfident in its learned representations, ensuring that it explores a wider range of possible embeddings and becomes more robust to small variations in the input data. By perturbing the embeddings during training, the model is forced to learn representations that are invariant to these perturbations, thus improving generalization (See the Appendix subsection A.3 for more analysis).

To show the regularizing effect of Gaussian noise, we examine the variance it introduces, noting that the expected value of the perturbed embedding remains unchanged from the original.

$$\mathbb{E}[\tilde{\mathbf{z}}_v] = \mathbb{E}[\mathbf{z}_v + \sigma \epsilon_v] = \mathbf{z}_v \tag{3}$$

However, the variance of the perturbed embedding increases due to the Gaussian noise, which adds a regularization effect to the model. The variance of the perturbed embedding is given by:

$$\operatorname{Var}[\tilde{\mathbf{z}}_{v}] = \operatorname{Var}[\mathbf{z}_{v} + \sigma\epsilon_{v}] = \sigma^{2} \cdot I \tag{4}$$

(5)

Thus, the total variance of the embeddings with Gaussian noise becomes:

245 246

249

250

251

252 253 254

255 256

257 258

259

260

261

262

264 265

266 267

Introducing the noise adds variance, helping the model avoid overfitting by promoting smoother decision boundaries and better generalization. This encourages the learning of robust features that perform well on unseen data by broadening the explored feature space during training. For inference, the noise is removed to ensure accurate predictions, returning embeddings to their original state.

 $\operatorname{Var}_{\operatorname{total}} = \operatorname{Var}[\mathbf{z}_v] + \sigma^2$ 

$$\mathbf{z}_v = A_X(\mathbf{x}_v), \quad \mathbf{z}_t = A_Y(\mathbf{x}_t) \tag{6}$$

### 3.4 EMBEDDING SMOOTHING HELPS MODERATE MODEL OVERCONFIDENCE

**Embedding Smoothing and Modified Loss Function:** Inspired by label smoothing Gong et al. (2024); Müller et al. (2019), to reduce the model's reliance on hard embeddings and enhance generalization, we design Embedding Smoothing (ES) into our Align-VL framework. In this context, each example within a batch is treated as a unique class, making the number of classes N equivalent to the batch size. For a given target example y, the smoothed target distribution is defined as:

$$\tilde{y}_i = \begin{cases} 1 - \alpha, & \text{if } i = y \\ \frac{\alpha}{N-1}, & \text{if } i \neq y \end{cases}$$
(7)

where  $\alpha \in (0, 1)$  is the smoothing parameter, N is the batch size, and i and y are indices of examples within the batch. ES assigns a small non-zero probability to all other examples in the batch, thereby preventing the model from becoming overconfident and improving its ability to generalize. 270 We incorporate these smoothed targets into a symmetric contrastive loss function. The perturbed 271 embeddings  $\tilde{\mathbf{z}}_v$  and  $\tilde{\mathbf{z}}_t$  from the two modalities (image and text) are projected using their respective 272 V-L Adapters  $A_X$  and  $A_Y$ . The loss function with Embedding Smoothing is defined as: 273

$$\mathcal{L}_{\text{sym}}^{\text{EmbedSmooth}} = \frac{1}{2} \left( \mathcal{L} \left( A_X(\tilde{\mathbf{z}}_v), \tilde{Y}; \tilde{\mathbf{z}}_t \right) + \mathcal{L} \left( A_Y(\tilde{\mathbf{z}}_t), \tilde{Y}; \tilde{\mathbf{z}}_v \right) \right)$$
(8)

The loss  $\mathcal{L}$  is computed using the Kullback-Leibler divergence between the smoothed target distributions and the model's predicted probabilities:

282

283 284 285

286

287

288

289

290

291

278

$$\mathcal{L}\left(A_X(\tilde{\mathbf{z}}_v), \tilde{Y}; \tilde{\mathbf{z}}_t\right) = \mathrm{KL}\left(\tilde{y} \left\| \operatorname{Softmax}\left(\frac{\sin\left(A_X(\tilde{\mathbf{z}}_v), A_Y(\tilde{\mathbf{z}}_t)\right)}{\tau}\right)\right)$$
(9)

$$\mathcal{L}\left(A_{Y}(\tilde{\mathbf{z}}_{t}), \tilde{Y}; \tilde{\mathbf{z}}_{v}\right) = \mathrm{KL}\left(\tilde{y} \left\| \operatorname{Softmax}\left(\frac{\sin\left(A_{Y}(\tilde{\mathbf{z}}_{t}), A_{X}(\tilde{\mathbf{z}}_{v})\right)}{\tau}\right)\right)\right.$$
(9)

where sim( $\cdot, \cdot$ ) denotes a cosine similarity measure.  $\tau$  is a hyperparameter that controls the concentration of the distribution. Softmax( $\cdot$ ) converts similarity scores into a probability distribution over the batch. By utilizing ES, the model is encouraged to produce output distributions that are less peaked and more spread out, which helps in preventing overfitting. The symmetric loss function ensures equal contributions from both modalities during the learning process.

Theoretical Analysis: ES increases the entropy of target distributions by assigning non-zero probabilities to all classes, reducing model overconfidence and reliance on specific training examples, 292 thereby enhancing robustness and generalization. Theoretically, it serves as regularization, prevent-293 ing excessive focus on any single class and promoting better generalization to unseen data. The 294 increase in entropy of the target distribution can be quantified. For the smoothed target distribution 295  $\tilde{y}$ , the entropy is:

296 297 298

299 300

301

302

303

304 305

306 307

$$H(\tilde{y}) = -\left((1-\alpha)\log(1-\alpha) + (N-1)\left(\frac{\alpha}{N-1}\log\frac{\alpha}{N-1}\right)\right)$$
(10)

Higher entropy in the target distribution smooths gradients in training, serving as regularization to prevent overfitting. ES also improves model calibration by using KL divergence to discourage overconfidence, promoting even probability distribution across positive and negative samples (refer to Appendix subsection A.2 for detailed analysis).

#### 4 **EXPERIMENTS**

**Baselines and Metrics:** We conduct comparisons with several multimodal alignment methods: 308 Fusemix Vouitsis et al. (2024), CLIP Radford et al. (2021), LIT Zhai et al. (2022), and 3T Kossen 309 et al. (2024), as shown in Table 1 and Table 2. It's important to note that large-scale image-text 310 datasets used by models like CLIP, LIT, and 3T, ranging from 400 million to 5 billion pairs, are 311 mostly sourced from the internet and not publicly available, with high computational costs making 312 them impractical for limited-resource scenarios. Therefore, Fusemix, detailed in Table 1, offers a 313 more applicable comparison with Align-VL. To assess the performance of multimodal alignment for 314 all considered methods, we use metrics such as R@1, R@5, and R@10. R@1 denotes Recall@1 for 315 either text-to-image or image-to-text, R@5 indicates Recall@5, and R@10 stands for Recall@10, applicable to both text-to-image and image-to-text retrieval scenarios Vouitsis et al. (2024). 316

317 Experimental Setup: To minimize computational demands, all our experiments are conducted on 318 a single 24GB NVIDIA 3090 GPU. We pre-compute latents from pre-trained unimodal encoders, 319 which are then discarded, extracting latents for each modality sequentially to avoid loading more 320 than one encoder at a time. For consistency and fair comparison, we use the same unimodal encoders 321 as Fusemix. V-L adapters are parameterized as lightweight MLPs featuring an inverted bottleneck architecture, inspired by previous studies Lin et al. (2015); Tolstikhin et al. (2021); Bachmann et al. 322 (2023). Each MLP incorporates residual blocks and a default final projection layer with a dimension 323 of 512, embedding each modality into a shared latent space. For the image encoder, we consider



Figure 4: The upper figure mentions "ambiguous samples" in datasets, including partially matched samples (a, b, c) and completely unmatched samples (d). During contrastive learning training, it can lead to positive pairs not being truly positive. The lower figure use CLIP to compute similarity for four datasets (For details, 346 see Table 6), revealing generally moderate alignment between images and texts (mostly ranging from 20-40%). 347 Notably, the COCO dataset shows higher alignment (indicative of superior data quality), which correlates with relatively higher performance metrics. 348

DINOv2 Oquab et al. (2023), and for the text side, we select text encoder with demonstrably seman-350 tic latent spaces, specifically BGE Xiao et al. (2023). We highlight that since our V-L adapters are 351 operating on low-dimensional latents, the computational cost to train them is minimal, and despite 352 training on a single GPU, we can use large batch sizes (up to Batch = 10K on 3090 GPU), which 353 has been shown to benefit contrastive learning Wu et al. (2018); Tian et al. (2020); He et al. (2020). 354 The smoothing parameter  $\alpha$  is set to 0.1, and the Gaussian noise level  $\sigma$  is set to 0.01. 355

**Training Datasets.** To evaluate the effectiveness of the Align-VL method for the task of modality 356 alignment, we conducted extensive comparative experiments against SoTA methods across various 357 datasets. following previous works Chen et al. (2020); Li et al. (2021; 2022; 2023), These datasets 358 include COCO (human-annotated) Lin et al. (2014b), Visual Genome (VG) (human-annotated) Kr-359 ishna et al. (2017a), SBU (web datasets) Ordonez et al. (2011b), and Conceptual Captions 3M 360 (CC3M, web datasets) Sharma et al. (2018b). Table 6 provides detailed information about these 361 four datasets. It is noteworthy that the original CC3M dataset, consisting of images stored as inter-362 net URLs, currently has only 1.5 million data pairs available. Using a single NVIDIA 3090 GPU for training, Align-VL achieved SoTA performance across datasets of comparable sizes. Detailed 364 results and analyses are presented in the subsequent sections.

365 366

367

344

345

349

**RESULTS & DISCUSSIONS** 4.1

368 Benchmark Comparison: As demonstrated in our experiments in Table 1, Align-VL consistently 369 outperforms Fusemix across various dataset sizes and configurations in image-text retrieval tasks. For example, on the COCO dataset ( $\approx 560K$  pairs), Align-VL achieves a significant improvement 370 in text-to-image retrieval, with a 3% higher R@1 score compared to Fusemix. On SBU dataset 371 ( $\approx 840K$  pairs), Align-VL surpasses Fusemix by over 4.48% in R@1 for text-to-image tasks. Fur-372 thermore, on a combined training configuration of four datasets (VG+COCO+SBU+CC3M, totaling 373  $\approx 3.5M$  pairs), Align-VL achieves a significant improvement of over 1.44% in R@1 for text-to-374 image retrieval. These results highlight Align-VL's robust generalization capabilities and superior 375 performance over the SoTA method. 376

The Align-VL demonstrates a significant advantage in achieving high performance with substan-377 tially lower training costs. Compared to training datasets of similar size, Align-VL outperforms



Figure 5: (a) In the original distribution, clusters of image-text pairs are tightly packed, showing high model confidence. After implementing the "Be Modest" Align-VL approach, the embeddings become more dis-396 persed, indicating reduced confidence in individual labels and a more robust, softer probability distribution 397 that enhances model generalization. (b) The ablation study of two techniques (Random Perturbation (RP) and Embedding Smoothing (ES)).

Table 1: The performance of Align-VL and Fusemix Vouitsis et al. (2024) on Flickr30K test set. By evaluating on multiple datasets with varying sizes and complexities, we show that our Align-VL model exhibits strong generalization capabilities and achieves SoTA performance on image retrieval tasks. Bold signifies the best.

		Flickr30K (1K test set)						
Size	Training Dataset	Method	$\text{text} \rightarrow \text{image}$			$image \rightarrow text$		
			R@1	R@5	R@10	R@1	R@5	R@10
$\approx 560K$	000	Fusemix (Vouitsis et al., 2024)	57.80	83.38	89.54	71.60	91.10	95.00
	000	Align-VL (ours)	60.80	84.82	90.82	72.60	93.30	95.70
$\approx 820K$	VG	Fusemix (Vouitsis et al., 2024)	51.66	79.20	86.66	66.90	90.20	95.40
	VU	Align-VL (ours)	52.90	80.64	87.86	70.20	90.40	95.90
$\approx 840 K$	SBU	Fusemix (Vouitsis et al., 2024)	43.32	72.48	81.36	61.60	86.30	92.10
	300	Align-VL (ours)	47.80	75.94	84.18	62.10	87.00	92.30
1.38M	VGLCOCO	Fusemix (Vouitsis et al., 2024)	61.56	86.00	91.40	76.40	93.60	97.10
	V0+C0C0	Align-VL (ours)	61.54	86.16	91.72	77.20	94.40	97.80
2M	VG+COCO+SBU	Fusemix (Vouitsis et al., 2024)	61.40	84.82	90.46	77.20	94.40	97.20
	VO+COCO+SBU	Align-VL (ours)	62.10	85.52	90.76	77.80	94.60	97.70
3.5M	VG+COCO+SPU+CC2M	Fusemix (Vouitsis et al., 2024)	64.28	87.60	91.86	81.20	96.40	98.20
	VO+COCO+SBU+CCSM	Align-VL (ours)	65.72	87.82	92.50	81.60	96.00	98.30

414 415 416

417

421

399

400

401

both CLIP and Fusemix, achieving an R@1 improvement of 11.42% in text-to-image retrieval and 418 14.2% in image-to-text retrieval over CLIP (Table 2). Additionally, it surpasses Fusemix with im-419 provements of 5.82% and 7.2% in both tasks, respectively, as detailed in Table 2. Furthermore, de-420 spite using significantly smaller training data (3.5M pairs), Align-VL performs competitively against much larger datasets used by CLIP (400M) and LIT (4B), trailing by only 0.78% in text-to-image 422 and 2.30% in image-to-text retrieval against LIT. This demonstrates the effectiveness and efficiency 423 of Align-VL. We anticipate that with further increases in training data, Align-VL has the potential to surpass these SoTA methods, further underscoring its usability. 424

425 Efficiency in Dataset Quality: To evaluate the dataset quality, we utilized the CLIP-ViT/B-32 426 model Radford et al. (2021) to compute the similarity between images and texts across four datasets, 427 as shown in Table 6. This analysis reveals weak associations and low data quality between images 428 and texts in exisiting datasets. For the CC3M, 93.5% of the image-text pairs have a similarity of less than 35, indicating a low level of alignment between the images and texts in the dataset. Similar 429 trends are observed in the COCO, SBU, and VG datasets, where 91.9%, 94.4%, and 99.1% of image-430 text pairs, respectively, have a similarity of less than 35 (See Figure 4 lower figure). This reflects the 431 prevalence of ambiguous or weakly paired image-text samples across these datasets. It's noteworthy

that COCO exhibits a noticeably higher degree of pairing among these datasets, resulting in higher
performance metrics for alignment models trained on it, underscoring the importance of data quality
(Table 1). In contrastive learning, datasets with poorer quality lead to many ambiguous positive
pairs, causing models to become overly confident during training and confused during inference.
Therefore, Align-VL becomes particularly crucial in addressing these challenges.

Table 2: The performance of SoTA methods and Align-VL on different training datasets is assessed on the Flickr30K dataset's 1K test set, evaluating text-to-image and image-to-text retrieval accuracy using R@1 scores.

Size	Method	Flickr30K (1K test set)			
5120	Wiethou	$\text{text} \rightarrow \text{image}$	image $\rightarrow$ text		
400M	CLIP (Radford et al., 2021)	68.70	88.00		
4B	LIT (Zhai et al., 2022)	66.50	83.90		
5B	3T (Kossen et al., 2024)	72.10	87.30		
3M	CLIP (Radford et al., 2021)	54.30	67.40		
3M	Fusemix $(D,B)$ (Vouitsis et al., 2024)	59.90	74.40		
5M	Fusemix $(U,E)$ (Vouitsis et al., 2024)	64.30	80.20		
3.5M	Align-VL (Ours)	65.72	81.60		

# 4.2 ABLATION STUDY

Effect of Random Perturbation (RP): To validate the effectiveness of RP, we conducted ablation experiments for RP. As shown in Table 3, adding RP consistently improves performance across all datasets. On the COCO dataset, adding Random Perturbation (RP) in two settings-with and without Embedding Smoothing (ES)-results in an R@1 score increase from 57.98% to 60.8%, a 2.82% improvement, and from 57.8% to 60.56%, a 2.76% improvement. Similarly, this positive trend in gains from RP is observable across other datasets as well, as shown in Figure 5 (b). These results demonstrate that RP significantly enhances model performance regardless of whether ES techniques are used.

Effect of Embedding Smoothing (ES): As shown in Table 3, ES improves the R@1 scores on the SBU dataset by 3.8% without Random Perturbation (RP) and by 1.6% with RP. Similar gains from ES are observed across other datasets, regardless of the RP setting, as shown in Figure 5. ES enhances performance on complex, large-scale datasets by improving generalization across diverse image-text pairs. As a regularization strategy in contrastive learning, it smooths target distributions to mitigate overfitting, thereby boosting model robustness and accuracy on validation datasets.

Table 3: Ablation experiments of different techniques, including Random Perturbation (RP) and Embedding Smoothing (ES), are conducted across various datasets, measuring text-to-image retrieval performance (R@1).

DD	EC	Training Dataset							
КГ	ЕЭ	COCO	SUB	VG+COCO+SBU	VG+COCO+SBU+CC3M				
		57.80	43.32	61.40	64.28				
	~	57.98	47.12	61.66	64.94				
~		60.56	46.20	61.66	65.04				
~	~	60.80	47.80	62.10	65.72				

Effect of Smoothing Parameter  $\alpha$  in ES: As shown in Table 4, default parameter ES consistently outperforms Dynamic  $\alpha$  of ES (with  $\alpha$  decreasing over training), demonstrating the default  $\alpha$ 's ef-fectiveness. On the COCO dataset, default ES achieves a marginal but notable 0.34% improvement in R@1 scores for text-to-image retrieval, while the improvement in image-to-text retrieval is more substantial, with a 0.8% increase. These results are consistent across other datasets, where default  $\alpha$  ES consistently delivers performance enhancements. While dynamic ES offers certain benefits, default ES is more effective in optimizing retrieval tasks, particularly by adapting to diverse datasets and handling complex data more efficiently. Both dynamic and default smoothing parameters sig-nificantly outperform models without ES, underscoring ES's effectiveness.

Table 4: Ablation Study of Different ES Types in Various Datasets (R@1). This study compares the performance of different types of Embedding Smoothing (ES) across several datasets. We use  $\checkmark$  to denote without ES,  $\bigcirc$  for dynamic  $\alpha$  in ES, and  $\checkmark$  for default  $\alpha$  in ES.

	Different Type of ES						
	$\text{text} \rightarrow \text{image}$			image $\rightarrow$ text			
Dataset	X	0	1	X	0	1	
COCO	57.80	60.22	60.56	71.60	72.70	73.50	
VG	51.66	52.50	53.80	66.90	68.10	69.20	
SUB	43.32	45.46	47.12	61.60	63.60	62.80	
VG+COCO	61.56	61.98	62.24	76.40	77.60	77.20	
VG+COCO+SUB	61.40	61.32	61.66	77.20	77.90	78.40	
VG+COCO+SUB+CC3M	64.28	65.14	64.94	81.20	81.50	82.50	

496 497 498

522

**Impact of Dataset Quality:** Smaller, human-annotated datasets like COCO ( $\approx 560K$ ) can signif-499 icantly outperform larger, web-sourced ones like the SBU ( $\approx 840K$ ) in retrieval tasks, as reflected 500 in their R@1 scores (57.80% vs 43.32% in text-image R@1, Table 5). Similar observations are 501 made when comparing combined datasets of VG+COCO+SBU with CC3M. This highlights that 502 careful curation of datasets often leads to better model performance and generalization than merely increasing dataset size. However, curating large-scale datasets is impractical, making the capabil-504 ity of Align-VL to adjust the degree of image-text match within the dataset (potentially enhancing 505 data quality) particularly valuable. As shown in Figure 5 (a), after processing with the Align-VL 506 algorithm, embeddings shift from "hard" to "modest" and become smoother, indicating reduced 507 confidence in individual samples and a more dispersed, robust probability distribution that enhances 508 model generalization. Align-VL potentially enhances dataset quality and the subsequent improve-509 ment in performance metrics across all datasets underscores its effectiveness. 510

Table 5: Analysis of the Impact of Dataset Quality and Size with Align-VL. This analysis compares the effects of varying dataset quality and size on the Align-VL model's performance. The smaller datasets of higher quality outperform larger but lower-quality datasets, underscoring the importance of data quality.

Size	Dataset	te	$xt \rightarrow ima$	age	$image \rightarrow text$						
5120	Dataset	R@1	R@5	R@10	@10 R@1 R@5 F						
$\approx 3M$	CC3M	59.90	86.40	91.60	74.70	94.00	<u>97.40</u>				
$\approx 2M$	VG+COCO+SBU	<u>61.40</u>	84.82	90.46	77.20	<u>94.40</u>	97.20				
$\approx 820K$	VG	51.66	79.20	86.66	66.90	90.20	<u>95.40</u>				
$\approx 840 K$	SUB	43.32	72.48	81.36	61.60	86.30	92.10				
$\approx 560 K$	COCO	<u>57.80</u>	83.38	<u>89.54</u>	71.60	<u>91.10</u>	95.00				

<sup>5</sup> CONCLUSION

523 In this work, we propose a multimodal alignment framework Align-VL that promotes modesty in 524 model predictions while being computationally and data-efficient. Align-VL is a straightforward yet 525 effective method that reduces model overconfidence and enhances the intrinsic connections of paired 526 data in the latent space. It effectively leverages guidance from pretrained visual and text unimodal 527 encoders. Notably, Align-VL excels on datasets with lower data quality (image-text match level) and also enhances performance on datasets with higher match levels. Validated across multiple 528 datasets, Align-VL has consistently demonstrated its robust capability to align V-L models better. 529 For future developments, Align-VL could incorporate data quality assessments to dynamically adjust 530 model confidence, applying stricter constraints on lower quality data and more lenient ones on higher 531 quality data to effectively manage model overconfidence. 532

534 LIMITATIONS

Although Align-VL achieved performance improvements across all datasets and their combinations,
it is notable that the gains are more significant for datasets of poorer quality, while the enhancements
are more modest for relatively better datasets. We are unable to test Align-VL on larger datasets like
the 400M pairs used in CLIP, as it is not publicly available. Consequently, it is difficult to ascertain
Align-VL's performance benefits for extremely large-scale datasets. Besides, it remains unclear how
much Align-VL benefits datasets of very high quality, where image-text pairs are perfectly matched.

#### 540 REFERENCES 541

551

563

564

565

566 567

569

571

574

575

576

577

582

583

584

585

- 542 Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, 543 Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob L Menick, 544 Sebastian Borgeaud, Andy Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikoł aj Bińkowski, 545 Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karén Simonyan. Flamingo: a Visual 546 Language Model for Few-Shot Learning. In Advances in Neural Information Processing Systems, 547 volume 35, pp. 23716–23736, 2022. 548
- 549 Relja Arandjelovic and Andrew Zisserman. Look, listen and learn. In Proceedings of the IEEE 550 International Conference on Computer Vision, 2017.
- Gregor Bachmann, Sotiris Anagnostidis, and Thomas Hofmann. Scaling MLPs: A Tale of Inductive 552 Bias. arXiv:2306.13575, 2023. 553
- 554 Hangbo Bao, Wenhui Wang, Li Dong, Qiang Liu, Owais Khan Mohammed, Kriti Aggarwal, 555 Subhojit Som, Songhao Piao, and Furu Wei. VLMo: Unified Vision-Language Pre-Training 556 with Mixture-of-Modality-Experts. In Advances in Neural Information Processing Systems, volume 35, pp. 32897–32912, 2022. 558
- 559 Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and 560 Jingjing Liu. UNITER: UNiversal Image-TExt Representation Learning. In Andrea Vedaldi, 561 Horst Bischof, Thomas Brox, and Jan-Michael Frahm (eds.), Computer Vision – ECCV 2020, pp. 104-120, 2020. ISBN 978-3-030-58577-8. 562
  - Rohit Girdhar, Mannat Singh, Nikhila Ravi, Laurens van der Maaten, Armand Joulin, and Ishan Misra. Omnivore: A single model for many visual modalities. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 16102–16112, June 2022.
- Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand 568 Joulin, and Ishan Misra. ImageBind: One Embedding Space To Bind Them All. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 15180–15190, 570 June 2023.
- Xiuwen Gong, Nitin Bisht, and Guandong Xu. Does label smoothing help deep partial label learn-572 ing? In Forty-first International Conference on Machine Learning, 2024. 573
  - Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In Proceedings of the IEEE/CVF Conference on *Computer Vision and Pattern Recognition*, pp. 9729–9738, 2020.
- 578 Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan 579 Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning 580 with noisy text supervision. In Proceedings of the 38th International Conference on Machine Learning, volume 139, pp. 4904–4916. PMLR, 18–24 Jul 2021. 581
  - Jannik Kossen, Mark Collier, Basil Mustafa, Xiao Wang, Xiaohua Zhai, Lucas Beyer, Andreas Steiner, Jesse Berent, Rodolphe Jenatton, and Efi Kokiopoulou. Three towers: Flexible contrastive learning with pretrained image models. arXiv:2305.16999, 2023.
- 586 Jannik Kossen, Mark Collier, Basil Mustafa, Xiao Wang, Xiaohua Zhai, Lucas Beyer, Andreas Steiner, Jesse Berent, Rodolphe Jenatton, and Effrosyni Kokiopoulou. Three towers: Flexible 588 contrastive learning with pretrained image models. Advances in Neural Information Processing 589 Systems, 36, 2024.
- 590

Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, Michael S. Bernstein, and Fei-Fei Li. 592 Visual genome: Connecting language and vision using crowdsourced dense image annotations. International Journal of Computer Vision, 123:32–73, 2017a.

598

607

609

616

624

631

637

- 594 Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie 595 Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting lan-596 guage and vision using crowdsourced dense image annotations. International journal of computer 597 vision, 123:32-73, 2017b.
- Junnan Li, Ramprasaath Selvaraju, Akhilesh Gotmare, Shafiq Joty, Caiming Xiong, and Steven Chu Hong Hoi. Align before fuse: Vision and language representation learning with momentum 600 distillation. In Advances in Neural Information Processing Systems, volume 34, pp. 9694–9705, 601 2021. 602
- 603 Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. BLIP: Bootstrapping language-image pre-604 training for unified vision-language understanding and generation. In Proceedings of the 39th 605 International Conference on Machine Learning, volume 162, pp. 12888–12900. PMLR, 17–23 606 Jul 2022.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. BLIP-2: Bootstrapping language-image 608 pre-training with frozen image encoders and large language models. In Proceedings of the 40th International Conference on Machine Learning. PMLR, 2023. 610
- 611 Xiujun Li, Xi Yin, Chunyuan Li, Pengchuan Zhang, Xiaowei Hu, Lei Zhang, Lijuan Wang, Houdong 612 Hu, Li Dong, Furu Wei, Yejin Choi, and Jianfeng Gao. Oscar: Object-Semantics Aligned Pre-613 training for Vision-Language Tasks. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-614 Michael Frahm (eds.), Computer Vision – ECCV 2020, pp. 121–137, 2020. ISBN 978-3-030-58577-8. 615
- Valerii Likhosherstov, Anurag Arnab, Krzysztof Marcin Choromanski, Mario Lucic, Yi Tay, and 617 Mostafa Dehghani. PolyViT: Co-training Vision Transformers on Images, Videos and Audio. 618 Transactions on Machine Learning Research, 2023. ISSN 2835-8856. 619
- 620 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr 621 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In Computer 622 Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, 623 Proceedings, Part V 13, pp. 740-755. Springer, 2014a.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr 625 Dollár, and C. Lawrence Zitnick. Microsoft COCO: Common Objects in Context. In Computer 626 Vision – ECCV 2014, pp. 740–755. Springer International Publishing, 2014b. ISBN 978-3-319-627 10602-1. 628
- 629 Zhouhan Lin, Roland Memisevic, and Kishore Konda. How far can we go without convolution: 630 Improving fully-connected networks. arXiv:1511.02580, 2015.
- Jiaxiang Liu, Jin Hao, Hangzheng Lin, Wei Pan, Jianfei Yang, Yang Feng, Gaoang Wang, Jin Li, 632 Zuolin Jin, Zhihe Zhao, et al. Deep learning-enabled 3d multimodal fusion of cone-beam ct and 633 intraoral mesh scans for clinically applicable tooth-bone reconstruction. Patterns, 4(9), 2023a. 634
- 635 Jiaxiang Liu, Tianxiang Hu, Yan Zhang, Yang Feng, Jin Hao, Junhui Lv, and Zuozhu Liu. Parameter-636 efficient transfer learning for medical visual question answering. *IEEE Transactions on Emerging* Topics in Computational Intelligence, 2023b. 638
- Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. arXiv 639 preprint arXiv:1608.03983, 2016. 640
- 641 Ilya Loshchilov and Frank Hutter. Fixing weight decay regularization in adam. 2018. 642
- 643 Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. ViLBERT: Pretraining Task-Agnostic Visi-644 olinguistic Representations for Vision-and-Language Tasks. In Advances in Neural Information 645 Processing Systems, volume 32, 2019.
- Rafael Müller, Simon Kornblith, and Geoffrey E Hinton. When does label smoothing help? Ad-647 vances in neural information processing systems, 32, 2019.

- 648 Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, 649 Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Mahmoud Assran, Nico-650 las Ballas, Russel Galuba, Wojciech Howes, Po-Yao Huang, Li Shang-Wen, Ishan Misra, Michael 651 Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Hervé Jegou, Julien Mairal, Patrick Labatut, 652 Armand Joulin, and Piotr Bojanowski. DINOv2: Learning robust visual features without supervision. arXiv:2304.07193, 2023. 653 654 Vicente Ordonez, Girish Kulkarni, and Tamara Berg. Im2text: Describing images using 1 million 655 captioned photographs. Advances in neural information processing systems, 24, 2011a. 656 Vicente Ordonez, Girish Kulkarni, and Tamara Berg. Im2Text: Describing Images Using 1 Million 657 Captioned Photographs. In Advances in Neural Information Processing Systems, volume 24, 658 2011b. 659 660 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agar-661 wal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya 662 Sutskever. Learning transferable visual models from natural language supervision. In Proceedings of the 38th International Conference on Machine Learning, volume 139, pp. 8748–8763. 663 PMLR, 18-24 Jul 2021. 664 665 Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, 666 hypernymed, image alt-text dataset for automatic image captioning. In Iryna Gurevych and 667 Yusuke Miyao (eds.), Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 2556-2565, Melbourne, Australia, July 668 2018a. Association for Computational Linguistics. doi: 10.18653/v1/P18-1238. URL https: 669 //aclanthology.org/P18-1238. 670 671 Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, 672 hypernymed, image alt-text dataset for automatic image captioning. In Proceedings of the 56th 673 Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 674 2556-2565, July 2018b. doi: 10.18653/v1/P18-1238. 675 Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. VL-BERT: Pre-676 training of Generic Visual-Linguistic Representations. In International Conference on Learning 677 Representations, 2020. 678 Chen Sun, Austin Myers, Carl Vondrick, Kevin Murphy, and Cordelia Schmid. VideoBERT: A 679 Joint Model for Video and Language Representation Learning. In Proceedings of the IEEE/CVF 680 International Conference on Computer Vision, 2019. 681 682 Hao Tan and Mohit Bansal. LXMERT: Learning cross-modality encoder representations from 683 transformers. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan (eds.), Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th 684 International Joint Conference on Natural Language Processing, pp. 5100–5111, 2019. doi: 685 10.18653/v1/D19-1514. 686 687 Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding. In Computer 688 Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, 689 Part XI 16, pp. 776–794. Springer, 2020. 690 Ilya O Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas Un-691 terthiner, Jessica Yung, Andreas Steiner, Daniel Keysers, Jakob Uszkoreit, Mario Lucic, and 692 Alexey Dosovitskiy. MLP-Mixer: An all-MLP architecture for vision. In Advances in Neural 693 Information Processing Systems, volume 34, pp. 24261–24272, 2021. 694 Noël Vouitsis, Zhaoyan Liu, Satya Krishna Gorti, Valentin Villecroze, Jesse C Cresswell, Guangwei 695 Yu, Gabriel Loaiza-Ganem, and Maksims Volkovs. Data-efficient multimodal fusion on a single 696 gpu. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 697 pp. 27239-27251, 2024. 698 Yi Wang, Kunchang Li, Yizhuo Li, Yinan He, Bingkun Huang, Zhiyu Zhao, Hongjie Zhang, Jilan 699
- Yi Wang, Kunchang Li, Yizhuo Li, Yinan He, Bingkun Huang, Zhiyu Zhao, Hongjie Zhang, Jilan Xu, Yi Liu, Zun Wang, Sen Xing, Guo Chen, Junting Pan, Jiashuo Yu, Yali Wang, Limin Wang, and Yu Qiao. InternVideo: General video foundation models via generative and discriminative learning. *arXiv:2212.03191*, 2022a.

702 703 704	Zirui Wang, Jiahui Yu, Adams Wei Yu, Zihang Dai, Yulia Tsvetkov, and Yuan Cao. SimVLM: Simple visual language model pretraining with weak supervision. In <i>International Conference on Learning Representations</i> , 2022b.
706 707 708	Peter Willett, Peter F Swaszek, and Rick S Blum. The good, bad and ugly: Distributed detection of a known signal in dependent gaussian noise. <i>IEEE Transactions on signal processing</i> , 48(12): 3266–3279, 2000.
709 710	Shengqiong Wu, Hao Fei, Leigang Qu, Wei Ji, and Tat-Seng Chua. NExT-GPT: Any-to-any multi- modal LLM. <i>arXiv:2309.05519</i> , 2023.
711 712 713 714	Zhirong Wu, Yuanjun Xiong, Stella X Yu, and Dahua Lin. Unsupervised feature learning via non- parametric instance discrimination. In <i>Proceedings of the IEEE Conference on Computer Vision</i> <i>and Pattern Recognition</i> , pp. 3733–3742, 2018.
715 716	Shitao Xiao, Zheng Liu, Peitian Zhang, and Niklas Muennighof. C-pack: Packaged resources to advance general chinese embedding. <i>arXiv:2309.07597</i> , 2023.
717 718 719 720	Yong Yang, Jiaxiang Liu, Shuying Huang, Weiguo Wan, Wenying Wen, and Juwei Guan. Infrared and visible image fusion via texture conditional generative adversarial network. <i>IEEE Transactions on Circuits and Systems for Video Technology</i> , 31(12):4771–4783, 2021.
721 722 723	Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. <i>Transactions of the Association for Computational Linguistics</i> , 2:67–78, 2014. doi: 10.1162/tacl_a_00166.
724 725 726 727 728	Lu Yuan, Dongdong Chen, Yi-Ling Chen, Noel Codella, Xiyang Dai, Jianfeng Gao, Houdong Hu, Xuedong Huang, Boxin Li, Chunyuan Li, Ce Liu, Mengchen Liu, Zicheng Liu, Yumao Lu, Yu Shi, Lijuan Wang, Jianfeng Wang, Bin Xiao, Zhen Xiao, Jianwei Yang, Michael Zeng, Luowei Zhou, and Pengchuan Zhang. Florence: A new foundation model for computer vision. <i>arXiv:2111.11432</i> , 2021.
729 730 731 732 733	Xiaohua Zhai, Xiao Wang, Basil Mustafa, Andreas Steiner, Daniel Keysers, Alexander Kolesnikov, and Lucas Beyer. LiT: Zero-Shot Transfer With Locked-Image Text Tuning. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 18123–18133, June 2022.
734 735	Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond Empir- ical Risk Minimization. In <i>International Conference on Learning Representations</i> , 2018.
737 738 739 740 741 742 743 744 745 746	Yiyuan Zhang, Kaixiong Gong, Kaipeng Zhang, Hongsheng Li, Yu Qiao, Wanli Ouyang, and Xi- angyu Yue. Meta-transformer: A unified framework for multimodal learning. <i>arXiv</i> :2307.10802, 2023.
747 748 740	
749 750 751	
752 753	
754 755	

#### 756 APPENDIX А

In this section, we present additional implementation details, experiment results, theoretical analysis, pseudo code and supplements. The content structure is outlined as follows:

Section A.1 - Assessing the Match Quality of Image-text Datasets with CLIP
– Section A.1.1 - Similarity Calculation
<ul> <li>Section A.1.2 - Similarity in Four Datasets</li> </ul>
<ul> <li>Section A.2 - Theoretical Analysis for Embedding Smoothing</li> </ul>
• Section A.3 - Theoretical Analysis for Random Perturbation
Section A.4 - Implementation Details
• Section A.5 - Pseudocode of Align-VL
A.1 Assessing the Match Quality of Image-text Datasets with CLIP
A.1.1 SIMILARITY CALCULATION
In the context of contrastive learning models such as CLIP, the similarity produced image-text pair is closely related to the cosine similarity of their respective embedding

such as CLIP, the similarity produced for a given 775 imilarity of their respective embeddings, modulated 776 by a temperature scaling factor  $\tau$ . Specifically, let Similarity(i, j) denote the similarity score for 777 image i and text j. This score can be mathematically expressed as: 778

758

759

760

761 762

763

764

765 766

772 773 774

780

781 782

783

784

where  $\tau$  is the temperature coefficient, which is set to 0.01. Consequently, when multiplied by the temperature coefficient  $\frac{1}{\tau}$ , the similarity score will be constrained within the new range.

 $\label{eq:Similarity} \text{Similarity}(i,j) = \frac{\text{image\_embedding}(i) \cdot \text{text\_embedding}(j)}{\|\text{image\_embedding}(i)\|\|\text{text\_embedding}(j)\|} \times \frac{1}{\tau}$ 

(11)

#### 785 A.1.2 SIMILARITY IN FOUR DATASETS 786

787 To assess the match quality of image-text datasets in the main text, we conducted similarity calculation experiments on four widely used image-text datasets: COCO, CC3M, SBU, and VG. These 788 datasets cover a diverse range of image content, from everyday objects to complex scenes, and vary 789 in terms of size and annotation quality. Table 6 presents a summary of these datasets, detailing 790 the number of image-text pairs in each and the ranges of similarity scores calculated for different 791 image-text pairings. 792

793 In a similarity analysis, the COCO dataset exhibited the highest mean similarity per image at 30.48, indicating a strong and consistent association between images and text. The CC3M and SBU 794 datasets demonstrated similar mean similarities of 28.80, reflecting comparable levels of alignment 795 quality. This highlights how dataset structure and complexity significantly influence model perfor-796 mance in multimodal tasks, as shown in Figure 4 (b). Table 6 shows the distribution of similarity 797 for Image-Text pairs across four datasets. It can be observed that the values are concentrated around 798 the 30 range, indicating that when CLIP is used as a scoring model, the resulting similarity tends to 799 cluster within a relatively modest range. 800

801

Table 6: Similarity Per Image-Text Across Ranges for Different Datasets. The table shows the count of sim-802 ilarity scores within specific similarity ranges for each dataset (CC3M, COCO, SBU, and VG). Each row 803 corresponds to a dataset, and each column represents a range of similarity values.

805		[5,10)	[10, 15)	[15, 20)	[20, 25)	[25, 30)	[30, 35)	[35, 40)	[40, 45)	[45, 50)
806	CC3M	11	400	14877	189355	616912	408777	70836	3389	32
807	COCO	1	29	1187	27566	216780	275238	44607	1331	8
808	SBU	15	647	15268	127793	365869	284130	45279	1702	24
809	VG	39	3149	68441	267798	342860	132259	7151	77	0



Figure 6: In terms of training efficiency, CLIP requires 3000 GPU days for training on 400 million data pairs, while Align-VL needs only approximately 5 GPU days for 3.5 million data pairs.

## A.2 THEORETICAL ANALYSIS FOR EMBEDDING SMOOTHING

Embedding Smoothing effectively increases the entropy of the target distributions by assigning nonzero probabilities to all classes (examples in the batch). This reduction in confidence prevents the model from becoming overly reliant on specific training examples, thereby enhancing its robustness. The inclusion of the smoothing parameter  $\alpha$  allows for control over the degree of smoothing applied, enabling a balance between model confidence and generalization ability. To provide a theoretical understanding of how Embedding Smoothing improves generalization, we analyze its impact on the loss function and the model's predictions.

In standard contrastive learning without smoothing, the loss for a positive pair is:

839

828

829 830

831 832

843

844 845

846

847

854 855 856

857 858 This loss encourages the model to maximize the similarity between positive pairs and minimize it between negative pairs. However, it can lead to overconfident predictions, as the model focuses heavily on the positive pair. With Embedding Smoothing, the loss incorporates the smoothed target distribution  $\tilde{y}$ , and the KL divergence becomes:

 $\mathcal{L}_{\text{pos}} = -\log \frac{\exp\left(\sin\left(A_X(z_x), A_Y(z_y)\right)/\tau\right)}{\sum_{j=1}^{N} \exp\left(\sin\left(A_X(z_x), A_Y(z_j)\right)/\tau\right)}$ 

$$\mathcal{L}\left(A_X(z_x), \tilde{Y}; Z_Y\right) = -\sum_{i=1}^N \tilde{y}_i \log p_i \tag{13}$$

(12)

where  $p_i$  is the predicted probability for the *i*-th example in the batch:

$$p_{i} = \frac{\exp\left(\sin\left(A_{X}(z_{x}), A_{Y}(z_{i})\right)/\tau\right)}{\sum_{i=1}^{N}\exp\left(\sin\left(A_{X}(z_{x}), A_{Y}(z_{j})\right)/\tau\right)}$$
(14)

By assigning non-zero probabilities  $\tilde{y}_i$  to all classes, the loss function penalizes the model not only for the positive pair but also for negative pairs, albeit to a lesser extent. This encourages the model to produce a probability distribution that is more uniform and less confident.

Analyzing from the perspective of information entropy: The entropy  $H(\tilde{y})$  of the smoothed target distribution is higher than that of a one-hot distribution. The entropy of  $\tilde{y}$  is:

864

866 867

872

873 874 875

876 877

889 890 891

900

901

$$H(\tilde{y}) = -\left((1-\alpha)\log(1-\alpha) + (N-1)\left(\frac{\alpha}{N-1}\log\frac{\alpha}{N-1}\right)\right)$$
(15)

Higher entropy in the target distribution leads to smoother gradients during training, which can prevent the model from fitting noise in the training data. This smoothing effect acts as a form of regularization, reducing overfitting.

**Reduction in Overconfident Predictions:** Embedding Smoothing reduces the Kullback-Leibler divergence between the predicted distribution p and the uniform distribution u, where  $u_i = \frac{1}{N}$ :

$$\mathrm{KL}(u\|p) = \sum_{i=1}^{N} u_i \log \frac{u_i}{p_i} \tag{16}$$

By making p closer to  $\tilde{y}$ , which has higher entropy, the model's predictions become less confident. This can be beneficial because overconfident predictions on training data often lead to poor generalization on unseen data.

Connection to Label Smoothing Theory: Embedding Smoothing in our context is analogous to
 label smoothing in classification tasks. Previous works have shown that label smoothing has the
 following effects: 1). Margin Maximization: It implicitly increases the decision margin between
 classes, which can improve generalization. 2). Penalization of Confident Wrong Predictions: By
 smoothing the targets, the loss function penalizes overconfident incorrect predictions more heavily.

887 Besides, consider the gradient of the loss with respect to the logits z:

$$\frac{\partial \mathcal{L}}{\partial z_i} = p_i - \tilde{y}_i \tag{17}$$

892 When using Embedding Smoothing,  $\tilde{y}_i$  is never exactly 0 or 1. This means that the gradients are non-893 zero for all classes, encouraging the model to adjust its predictions across all examples in the batch. This leads to more generalized feature representations. By incorporating Embedding Smoothing 894 into the loss function, we introduce a regularization effect that enhances the model's generalization 895 capabilities. The smoothing parameter  $\alpha$  provides a mechanism to control this effect, allowing for a 896 trade-off between fitting the training data and maintaining robustness to unseen data. This theoretical 897 understanding aligns with our experimental observations, where models trained with Embedding 898 Smoothing demonstrate improved performance on validation datasets. 899

A.3 THEORETICAL ANALYSIS FOR RANDOM PERTURBATION

902 Why Choose Gaussian noise?: This paper introduces Gaussian noise as a perturbation in Align-903 VL for several reasons: 1). Well-Defined Mathematical Properties: Gaussian distribution exhibits 904 continuous and smooth probability density functions across the real number line, facilitating theo-905 retical analysis and calculations. 2). Zero-Mean Symmetry: By choosing a Gaussian distribution 906 with a mean of zero, the added noise is symmetrically balanced around zero, introducing no sys-907 tematic bias and only increasing the variance, thereby preserving the expected value of embeddings. 908 3). Adjustable Perturbation Intensity: The standard deviation of the Gaussian distribution can be 909 precisely controlled, allowing for careful calibration of noise intensity. This flexibility is crucial for introducing an appropriate level of uncertainty to enhance model robustness. 4). Alignment with 910 Natural Phenomena: According to the Central Limit Theorem, the sum of many independent random 911 variables tends toward a Gaussian distribution. Thus, Gaussian noise effectively simulates random 912 disturbances or measurement errors prevalent in natural and engineering contexts. 5). Facilitation 913 of Optimization and Training: In deep learning, incorporating Gaussian noise helps to smooth the 914 loss function landscape, avoiding local minima and promoting more effective training processes. 915

In summary, choosing Gaussian noise as the source of perturbation provides theoretical soundness
 and practical convenience, aiding in the development of more robust feature representations, preventing overfitting, and enhancing the generalization capabilities of models.

Loss Function with Perturbed Embeddings: The perturbed embeddings are used in the training loss, specifically in contrastive learning with the InfoNCE loss. The objective is to maximize the similarity between the noisy visual and textual embeddings:

$$\mathcal{L}_{\text{NCE}} = -\log \frac{\exp(\sin(\tilde{\mathbf{z}}_v, \tilde{\mathbf{z}}_t)/\tau)}{\sum_{j=1}^N \exp(\sin(\tilde{\mathbf{z}}_v, \tilde{\mathbf{z}}_t^j)/\tau)}$$
(18)

where sim $(\cdot, \cdot)$  represents a similarity function,  $\tau$  is a temperature parameter, and N is the number of negative samples. This formulation ensures that the model learns representations that are invariant to noise, thus improving generalization.

Noise and Regularization Effect: To understand the impact of Gaussian noise on regularization, we first look at the expected value of the perturbed embeddings. Since the noise is zero-mean, the expected value of the perturbed embeddings is identical to the original embeddings:

$$\mathbb{E}[\tilde{\mathbf{z}}_v] = \mathbb{E}[\mathbf{z}_v + \sigma \epsilon_v] = \mathbf{z}_v \tag{19}$$

However, the variance of the perturbed embeddings increases due to the added noise. The varianceof the perturbed embeddings can be calculated as:

$$\operatorname{Var}[\tilde{\mathbf{z}}_{v}] = \operatorname{Var}[\mathbf{z}_{v} + \sigma \epsilon_{v}] = \operatorname{Var}[\mathbf{z}_{v}] + \sigma^{2} \operatorname{Var}[\epsilon_{v}]$$
(20)

Given that  $Var[\epsilon_v] = I$ , where I is the identity matrix, the total variance of the perturbed embeddings becomes:

$$\operatorname{Var}_{\text{total}} = \operatorname{Var}[\mathbf{z}_{v}] + \sigma^{2}I \tag{21}$$

The additional term  $\sigma^2 I$  acts as a regularizer, which spreads out the embeddings and prevents the model from becoming overconfident in its predictions.

**Minimizing the Generalization Error:** The added noise effectively smooths the decision boundary of the model, which reduces overfitting. By introducing noise, we minimize the generalization error. Assuming the model's prediction function is f(z) and the true function is  $f^*(z)$ , the goal is to minimize the expected generalization error:

$$\mathbb{E}_{\mathbf{z}}[(f(\mathbf{z}) - f^*(\mathbf{z}))^2] \tag{22}$$

With noise perturbation, the variance in the embeddings increases, which forces the model to learn smoother decision boundaries. The regularization effect introduced by the noise helps bind the generalization error:

$$\mathbb{E}_{\mathbf{z}}[(f(\mathbf{z}) - f^*(\mathbf{z}))^2] \le \operatorname{Var}_{\operatorname{total}} = \operatorname{Var}[\mathbf{z}_v] + \sigma^2 \tag{23}$$

Thus, the noise helps control the generalization error by ensuring that the model does not overfit to specific features of the training data, which is especially important in cases where the training data contains noise or is limited in size.

964 Noise-Induced Gradient Regularization: We can also analyze the effect of noise on the gradient 965 of the loss function. Given a loss function  $\mathcal{L}(\mathbf{z}_v, \mathbf{z}_t)$ , the gradient with respect to the perturbed 966 embeddings can be expressed as:

$$\nabla_{\tilde{\mathbf{z}}_{v}} \mathcal{L}(\tilde{\mathbf{z}}_{v}, \tilde{\mathbf{z}}_{t}) = \nabla_{\mathbf{z}_{v}} \mathcal{L}(\mathbf{z}_{v}, \mathbf{z}_{t}) + \sigma \cdot \nabla_{\epsilon_{v}} \mathcal{L}(\tilde{\mathbf{z}}_{v}, \tilde{\mathbf{z}}_{t})$$
(24)

970 The second term,  $\sigma \cdot \nabla_{\epsilon_v} \mathcal{L}(\tilde{\mathbf{z}}_v, \tilde{\mathbf{z}}_t)$ , acts as a regularizer that prevents the gradient from becoming 971 too large. The gradient is smoothed by the presence of noise, which further prevents overfitting and encourages the model to learn more generalizable patterns. Inference Phase: During the inference phase, we remove the Gaussian noise to ensure accurate predictions on unseen data. The embeddings revert to their original clean form: 

$$\mathbf{z}_v = F(\mathbf{x}_v), \quad \mathbf{z}_t = G(\mathbf{x}_t) \tag{25}$$

Without the added noise, the model makes precise predictions based on the robust features it learned during training. Therefore, by adding Gaussian noise to the embeddings during training, we in-troduce a form of regularization that improves the generalization ability of the model. The noise prevents overfitting by increasing the variance of the embeddings, ensuring that the model learns smoother decision boundaries. This leads to better performance on unseen data and helps minimize the generalization error. The noise-induced gradient regularization further contributes to preventing the model from overfitting to the training data, making it more robust in real-world applications. 

#### A.4 IMPLEMENTATION DETAILS

For all experiments, we use the AdamW Loshchilov & Hutter (2018) optimizer during training. We perform learning rate warmup by linearly increasing the learning rate from  $10^{-6}$  to  $10^{-3}$ . We then decay the learning rate using a cosine schedule Loshchilov & Hutter (2016). We use a depth of 4 for both V-L adapters which we train for 500 epochs with a batch size of Batch = 10K. We set the learning rate as  $1r = 10^{-3}$  and use weight decay of 0.1 during optimization. The image encoder is DINOv2 ViT-G/14, and for the text side, the text encoder is the BGE large version. To evaluate the effectiveness of the Align-VL method for the task of modality alignment, we conducted extensive comparative experiments against SoTA methods across various datasets. These datasets include COCO Lin et al. (2014a), VG Krishna et al. (2017b), SBU Ordonez et al. (2011a), and CC3M Sharma et al. (2018a). Table 6 provides detailed information about these four datasets. Additionally, we compared different schemes such as Fusemix, CLIP, and LIT. Utilizing a single NVIDIA 3090 GPU for training, Align-VL demonstrated SoTA performance across datasets of varying sizes. 

- A.5 PSEUDOCODE OF ALIGN-VL

1026

1027 Algorithm 1: PyTorch-style pseudocode for Align-VL 1028 **Input:**  $A_X, A_Y$ : learnable V-L adapters; 1029 Batch: batch size; 1030  $D_x, D_y$ : latent dimensions of unimodal encoders; 1031  $D_s$ : latent dimension of shared space; 1032  $\beta$ : Mixup Beta distribution hyperparameter; 1033 t: learnable temperature parameter; 1034  $\alpha$ : smoothing parameter; 1035  $\sigma$ : Gaussian noise level. 1036 1 for  $z_x, z_y$  in loader do 1037 // FuseMix Split  $z_x$  and  $z_y$  into two parts: 2 1039  $z_{x1}, z_{x2} \in \mathbb{R}^{B \times D_x}, z_{y1}, z_{y2} \in \mathbb{R}^{B \times D_y};$ 1040 Sample mixing coefficient  $\lambda \sim \text{Beta}(\beta, \beta)$ ; 4 Mix embeddings: 1041 5  $z_x = \lambda z_{x1} + (1 - \lambda) z_{x2};$ 1042 6  $z_y = \lambda z_{y1} + (1 - \lambda) z_{y2};$ 1043 7 // Add Gaussian noise perturbation (training only) 1044 if training then 8 1045  $z_x = z_x + \sigma \times \mathcal{N}(0, I);$ 1046  $z_y = z_y + \sigma \times \mathcal{N}(0, I);$ 10 1047 end 11 1048 // Project into joint space and normalize 1049  $s_x = \text{normalize}(A_X(z_x)) \in \mathbb{R}^{B \times D_s};$ 12 1050  $s_y = \text{normalize}(A_Y(z_y)) \in \mathbb{R}^{B \times D_s};$ 13 1051 // Compute pairwise cosine similarities with temperature 1052  $\text{logits}_{xy} = (s_x s_y^{\top}) \times \exp(t) \in \mathbb{R}^{B \times B};$ 14 1053  $\text{logits}_{yx} = (s_y s_x^{\top}) \times \exp(t) \in \mathbb{R}^{B \times B};$ 15 1054 // Embedding Smoothing 1055 if  $\alpha > 0$  then 16 1056 Create smoothed targets: 17 1057  $\tilde{Y} = (1 - \alpha) \times I_B + \frac{\alpha}{B - 1} \times (\mathbf{1} - I_B) \in \mathbb{R}^{B \times B};$ 1058 18 1059 Compute losses: 19 1060  $loss_{xy} = KLDivLoss (log softmax(logits_{xy}), \tilde{Y});$ 20 1061  $loss_{yx} = KLDivLoss (log softmax(logits_{yx}), \tilde{Y});$ 1062 21 1063 end 22 1064 23 else 1065 // Standard symmetric alignment loss 1066 Labels: labels = [0, 1, ..., B - 1];24 1067 25 Compute losses: 1068  $loss_{xy} = CrossEntropyLoss (logits_{xy}, labels);$ 26 1069 27  $loss_{yx} = CrossEntropyLoss (logits_{ux}, labels);$ 1070 28 end 1071 29 Compute average loss: 1072  $\log = \frac{\log s_{xy} + \log s_{yx}}{2};$ 30 1073 21074 // Optimize optimizer.zero\_grad(); 1075 31 loss.backward(); 1076 32 optimizer.step(); 1077 33 34 end 1078 1079