# MLA: MultiLingual Acquisition on Multimodal Pre-training

Anonymous Author(s) Affiliation Address email

## Abstract

Vision and diverse languages are important information sources in our living world. 1 A model that understands multi-modalities and multi-languages can be applied to 2 a wider range of real-life scenarios. To build such a multimodal and multilingual 3 model, existing works try to ensemble vision-language data from multiple lan-4 guages in pre-training. However, due to the large number of languages, these works 5 often require huge computing resources and cannot be flexibly extended to new 6 languages. In this work, we propose a MultiLingual Acquisition (MLA) frame-7 work that can easily empower a monolingual Vision-Language Pre-training (VLP) 8 model with multilingual capability. Specifically, we design a lightweight language 9 acquisition encoder based on state-of-the-art monolingual VLP models. We further 10 propose a two-stage training strategy to optimize the language acquisition encoder, 11 namely the Native Language Transfer stage and the Language Exposure stage. 12 With much less multilingual training data and computing resources, our model 13 achieves state-of-the-art performance on multilingual image-text and video-text 14 retrieval benchmarks. 15

# 16 **1 Introduction**

We live in a multimodal and multilingual world. 17 The information we receive in our daily lives 18 may come from different modalities and lan-19 guages. Therefore, building multimodal and 20 multilingual models to effectively understand 21 22 such information has attracted much research attention [12, 36, 20, 3]. Recently, Multilin-23 gual Vision-Language Pre-training (M-VLP) 24 achieves convincing performance in various 25 cross-lingual cross-modal tasks such as multilin-26 gual image-text retrieval [27, 42, 11, 15, 17] and 27 multimodal machine translation [34]. As shown 28 in Figure 1(a), M-VLP models handle multiple 29



(a) Multilingual Vision-Language Pre-training (b) MultiLingual Acquisition

Figure 1: Comparison of data usage between M-VLP and MLA. The size of a circle reflects the amount of training data.

languages and modalities simultaneously during pre-training. Despite their successes, M-VLP mod-30 els suffer from two problems. First, pre-training on vision and multilingual data consumes huge 31 computing resources. For example, the state-of-the-art M-VLP model MURAL [17] is pre-trained on 32 128 Cloud TPUv3 for four days. It could support multimodal tasks on 100+ languages. However, 33 considering there are 6,900+ languages worldwide [42], building such a single model to handle all 34 languages will be highly expensive. Second, M-VLP models cannot be flexibly extended to new 35 languages. Additional training is required for M-VLP models to achieve satisfactory performance 36 on a new language. However, this training process will cause performance degeneration of M-VLP 37

Submitted to 36th Conference on Neural Information Processing Systems (NeurIPS 2022). Do not distribute.

models on the original languages due to the limited model capacity. For example, the limited model

<sup>39</sup> capacity even results in M-VLP models performing worse than their monolingual counterparts on

40 English [27, 42].

To build multimodal and multilingual models with low-cost and high-flexibility, we refer to our 41 human learning habits when acquiring new languages. We humans normally learn our native language 42 during childhood and practice it through interactions with the multimodal living environments. When 43 learning a new language, we humans initially tend to align it with the native language, as we can 44 easily map words in the native language to real-world objects and concepts. After having a certain 45 language foundation, we could further master it by interacting with the environment directly using 46 the new language. This is known as the language exposure [5]. The whole learning process rarely 47 degrades our native language capability. 48

Inspired by this, we propose a new framework, MultiLingual Acquisition (MLA), which constructs 49 multimodal and multilingual models based on monolingual VLPs. The topology of the MLA-based 50 multimodal and multilingual model is illustrated in Figure 1(b). Unlike M-VLPs, which handle data 51 from multiple languages and modalities in a single model, MLA empowers monolingual VLPs with 52 multilingual capability using much less training data through a language acquisition encoder. The 53 language acquisition encoder is realized by inserting our proposed lightweight language acquirers 54 into the pre-trained monolingual encoder of the VLP model. During training, original parameters in 55 the pre-trained monolingual encoder are fixed, only multi-lingual embeddings and language acquirers 56 for each new language are optimized. Following the human learning habits, we propose a two-stage 57 training strategy to train the language acquisition encoder. In the Native Language Transfer (NLT) 58 stage, the model is trained to establish the correspondence between the new languages with the 59 native language. In the Language Exposure (LE) stage, the model is optimized to build cross-modal 60 alignment between new languages and images. We apply our proposed MLA to the monolingual 61 VLP model CLIP [30] and achieve state-of-the-art results on both multilingual image-text and 62 video-text retrieval benchmarks with much less training data and computing resources. Ablation 63 studies demonstrate the effectiveness of our training strategy. Owing to the independence merit of 64 the language acquirers, the MLA-based models can be easily extended to new languages without 65 compromising the performance of their original languages. 66

The main contributions of our work are as follows: 1) We propose a lightweight MultiLingual Acquisition (MLA) framework that can easily empower monolingual VLPs with multilingual capability. 2) We propose a two-stage training strategy to optimize the MLA-based models inspired by the language learning habits of humans. Ablation studies prove the effectiveness of the strategy. 3) We apply MLA to the monolingual VLP model CLIP and achieve the new state-of-the-art results on both multilingual image-text and video-text retrieval benchmarks with much less training data and parameters.

# 73 2 Related Work

Vision-Language Pre-training: There are increasing interest in building Vision-Language Pre-74 training (VLP) models. From the perspective of how to interact between vision and language 75 modalities, existing models can be divided into two categories: single-stream and dual-stream 76 models. The single-stream models perform interaction on image and text directly with a cross-77 78 modal transformer [7, 25, 21]. In contrast, the dual-stream models encode image and text with two independent encoders and optimize via simple objectives like image-text contrastive learning 79 80 [30, 18, 41]. Compared with the single-stream models, the dual-stream models are more efficient to 81 utilize noisy image-text data harvested from the web [16], and thus achieve better performance on downstream tasks. Meanwhile, the dual-stream models are more flexible for extension. Since the 82 dual-stream models process images and text through independent encoders, we can fix the vision 83 encoders and focus on extending the text encoders to support new languages. Therefore, we focus on 84 empowering dual-stream VLPs with multilingual capability in this work. 85

Multilingual Vision-Language Pre-training: To achieve both multilingual and multimodal capability, many works try to learn the relationship between multiple languages and modalities simultaneously through pre-training. M<sup>3</sup>P [27] introduces the multimodal code-switched training method
 to enhance multilingual transferability. UC<sup>2</sup> [42] augments the English image-text data to other
 languages through machine translation and proposes fine-grained pre-training objectives to encourage
 alignment between image regions and multilingual tokens. More recently, MURAL [17] adopts



Figure 2: Model illustration: (a) The overview of MLA framework. (b) The structure of a language acquirer

the dual-stream structure. It is pre-trained with image-text and text-text contrastive objectives on 92 93 multilingual image-text pairs and translation pairs. M-VLP models significantly outperform previous non-pretraining models [12, 36, 20, 3] on multilingual image-text retrieval. Despite their success, 94 these models typically consume huge computing resources and large-scale multilingual training data. 95 Moreover, they fail to take full advantage of the cross-modal knowledge learnt in monolingual VLP, 96 and building cross-modal cross-lingual representations from scratch can be very hard. In contrast, our 97 MLA framework aims to empower VLP models with multilingual capability and it builds multimodal 98 and multilingual models with much less data and computing cost. 99

**Multilingual Extension:** Some works explore making pre-trained monolingual language models 100 multilingual. Reimers et al. [31] extend sentence embeddings from monolingual to multilingual by 101 Multilingual Knowledge Distillation (MKD). Artetxe et al. [2] extend monolingual models by training 102 additional word embeddings. MAD-X [29] extends multilingual pre-training models to support 103 low-resource languages through adapters [14]. By extending state-of-the-art pre-trained language 104 models, these works have achieved impressive results in NLP tasks such as bitext retrieval [31], 105 cross-lingual QA and NER [29, 31]. However, few works focus on making VLP models multilingual. 106 Work in [28] is the first to extend single-stream VLP model OSCAR [25]. It adopts a similar 107 108 strategy with MAD-X [29] that trains language adapters with Masked Language Modeling (MLM) for each language. During inference, it replaces the English language adapters with the target language 109 adapters to achieve zero-shot cross-lingual transfer. However, it generalizes poorly on other languages 110 since the MLM-based training strategy can only implicitly establish the correspondence between 111 other languages and English, let alone vision correspondences. In contrast, MLA directly builds the 112 connection of other languages with English and then with vision in the two-stage training strategy. 113 Therefore, MLA achieves comparable results on other languages as on English in downstream tasks. 114

# 115 3 Method

The MultiLingual Acquisition (MLA) framework is proposed to empower a dual-stream monolingual
VLP model with multilingual capability. We define the *native language* of a VLP as its pre-training
language. In this paper, we choose CLIP-ViT-B [30] as the VLP model. It is pre-trained with 400M
image-text pairs in English [30]. Note that MLA can also be applied to VLP models with different
native languages.

Since the state-of-the-art VLP models can project vision and native language into a shared multimodal space, we design a language acquisition encoder to process non-native languages. We then simulate the learning habits of human beings and propose a two-stage training strategy to optimize the language acquisition encoder. We first introduce the architecture of the MLA framework in Sec.3.1. Then, we describe our training strategy in Sec.3.2.

## 126 3.1 Architecture

Figure 2(a) illustrates the overview of the MLA framework, which consists of three modules: the pre-trained text encoder, the pre-trained vision encoder, and the language acquisition encoder.

**Pre-trained Text Encoder.** Given a sentence S in the native language, the corresponding sentence 129 representation  $s = \Phi(S; \theta_{\Phi})$  is generated through the pre-trained text encoder  $\Phi$ . To preserve the 130 cross-model knowledge of VLP,  $\theta_{\Phi}$  is keep fixed during training. As shown in the top part of Figure 131 2(a), the pre-trained text encoder contains a native embedding block and *l* transformer layers [35]. 132 The native embedding block first tokenizes S with byte pair encoding (BPE) [32]. Then, it converts 133 words into embeddings  $E_S = [e_{0=[SOS]}, e_1, \dots, e_{M=[EOS]}]$ . [SOS] and [EOS] are special tokens 134

denoting the boundary of 
$$S$$
. The word embeddings are then passed through the transformer layers:

$$H^{0} = [e_{0} = [sos], e_{1}, \dots, e_{M} = [eos]] + E_{pos}$$
(1)

$$H^{i} = \texttt{TransformerLayer}(H^{i-1}; \theta_{\Phi}^{i})$$
(2)

where  $H^i = [h_0^i, \dots, h_M^i]$  is the hidden state of the layer *i*.  $\theta_{\Phi}^i$  denotes the parameters of the layer *i*.  $E_{pos}$  is the positional encoding. Note that the causal self-attention mask is used in the transformer 136

137 layers [30]. The last hidden state of the [EOS] token is chosen to generate the sentence representation: 138

$$\boldsymbol{s} = W_a h_M^l \tag{3}$$

where s is the sentence representation of S, and  $W_a$  denotes a linear projection. 139

**Pre-trained Vision Encoder.** We extract the representation  $v = \Psi(V; \theta_{\Psi})$  of an image V with 140 the pre-trained vision encoder  $\Psi$ . Similar with the pre-trained text encoder,  $\theta_{\Psi}$  is also frozen. The 141 pre-trained vision encoder is implemented as a Vision Transformer [9]. As shown in the bottom part 142 of Figure 2(a), it consists of a image embedding block and l transformer layers. Given an image V, 143 the image embedding block first divides V into patches  $V' = [v'_1, \ldots, v'_N]$  following [9]. Then, they are linearly projected into patch embeddings  $E_p = [e_{\text{[CLASS]}}, W_p v'_1, \ldots, W_p v'_N]$ , where  $e_{\text{[CLASS]}}$  is a special embedding for the whole image and  $W_p$  is the linear projection. The patch embeddings are 144 145 146 then fed into transformer layers: 147

$$Z^{0} = [e_{[\text{CLASS}]}, W_{p}v'_{1}, \dots, W_{p}v'_{N}] + E_{pos}$$
(4)

$$Z^{i} = \operatorname{TransformerLayer}(Z_{i-1}; \theta_{V}^{i})$$
(5)

where  $Z^i = [z_0^i, \ldots, z_N^i]$  is the hidden state of the layer *i*. The last hidden state of the [CLASS] 148 embedding  $z_0^l$  is selected to produce the representation of image V: 149

$$\boldsymbol{v} = W_b \boldsymbol{z}_0^l \tag{6}$$

where v is the image representation of V, and  $W_b$  denotes a linear projection. 150

Language Acquisition Encoder. As shown in the middle part of Figure 2(a), the language ac-151 quisition encoder is built upon the pre-trained text encoder. Suppose T is a sentence written in 152 a non-native language L, we get the representation of T through language acquisition encoder 153  $t = \Phi'(T; \theta_{\Phi}, \theta_{emb}, \theta_L)$ , where  $\theta_{\Phi}$  are fixed parameters of the pre-trained text encoder,  $\theta_{emb}$ 154 refers to a shared non-native embedding block and  $\theta_L$  represents specialized language acquirers 155 for language L. Non-native sentence T is first tokenized and processed into word embeddings 156  $E_T = [u_{0=[SOS]}, \ldots, u_{M=[EOS]}]$  through the non-native embedding block. The word embeddings are 157 then encoded through the pre-trained transformer layers and language acquirers: 158

$$X^{0} = [W_{e}u_{0}=[\text{sos}], W_{e}u_{1}, \dots, W_{e}u_{m}=[\text{eos}]] + E_{pos}$$
(7)

$$H^{i} = \operatorname{TransformerLayer}(X^{i-1}; \theta_{\Phi}^{i})$$
(8)

$$X^{i} = \mathsf{LA}(H^{i}; \theta_{L}^{i}) \tag{9}$$

where  $X^i = [x_0^i, \ldots, x_m^i]$  is the hidden state of the layer *i*.  $W_e$  is a linear projection to keep dimension consistency.  $\theta_L^i$  denotes the parameters of the *i*-th language acquirer for language L. 159 160 As shown in Figure 2(b), the language acquirer is implemented as a bottleneck MLP with residual 161 connection [13]: 162

$$LA(X) = W_{upper} ReLU(W_{down}X) + X$$
(10)

Similar with the pre-trained text encoder, the last hidden state of the [EOS] token is projected into 163 the sentence representation t: 164

$$\boldsymbol{t} = W_a \boldsymbol{x}_m^l \tag{11}$$

Eq.11 shares the same linear projection  $W_a$  with Eq.3. The main advantage of the language acquisition 165 encoder is that it can extend the VLP models to support new languages without influencing the existing 166 languages, as it handles different languages with independent language acquirers. 167

#### 168 3.2 Training Strategy

To simulate the language learning habits of humans, we optimize the model in two stages: the Native Language Transfer (NLT) stage and the Language Exposure (LE) stage.

**Native Language Transfer.** When learning a new language, we humans initially tend to align it with the native language. To simulate this learning phase, we align the non-native representations to the native representations during the Native Language Transfer (NLT) stage. Specifically, suppose  $\{(S_1, T_1), ..., (S_n, T_n)\}$  are translation pairs, where  $S_i$  is in the native language, and  $T_i$  is in a non-native language L. The objective in the NLT stage is minimizing the Mean Square Error (MSE) between the native representation  $s_i = \Phi(S_i; \theta_{\Phi})$  and the non-native representation  $t_i =$  $\Phi'(T_i; \theta_{\Phi}, \theta_L, \theta_{emb})$ :

$$\mathcal{L}_{\rm NLT} = \frac{1}{B} \sum_{i=1}^{B} \| \boldsymbol{s}_i - \boldsymbol{t}_i \|^2$$
(12)

where B is the batch size. Note that  $\theta_{\Phi}$  is loaded from the VLP model and is kept frozen.  $\theta_L$  is trained for non-native language L.  $\theta_{emb}$  is shared among non-native languages.

During the NLT stage, the non-native language correspondence with vision can be built by pivoting
 on the native language, since the correspondence between the native language and vision is well
 established through VLP.

Language Exposure. After the NLT stage, the model has built an implicit connection between 183 non-native languages and vision. However, due to the existence of synonyms, two same words in 184 the native language may correspond to different images. Thus, ambiguity may arise when learning 185 non-native languages solely by relying on the native language. Actually, we can regard the language 186 acquisition encoder after the NLT stage as a person with a certain language foundation. He/She has 187 learned the basic usage of a language through native language teaching. To master it, he/she may 188 practice the non-native language by interacting with the multimodal living environments. Inspired by 189 this learning phase, we directly establish the cross-modal alignment between non-native languages 190 and vision during the Language Exposure (LE) stage. Given image-text pairs  $\{(V_1, T_1), ..., (V_n, T_n)\}$ 191 where  $T_i$  is in a non-native language L, the sentence representation  $t_i = \Phi'(T_i; \theta_{\Phi}, \theta_L, \theta_{emb})$  should be closer to the aligned image representation  $v_i = \Psi(V_i; \theta_{\Psi})$ , and away from the misaligned one 192 193  $v_j = \Psi(V_j; \theta_{\Psi}), j \neq i$ . This can be achieved by performing contrastive learning between non-native 194 languages and images. For a non-native sentence  $T_i$ , we treat the corresponding image  $V_i$  as a positive 195 sample, and other images in the same batch  $V_j$ ,  $j \neq i$  as negative samples. Vice versa for images. The objective in the LE stage is minimizing the NCE loss defined as follows: 196 197

$$\mathcal{L}_{\rm LE} = \frac{1}{2} (\mathcal{L}_{v2t} + \mathcal{L}_{t2v}) \tag{13}$$

$$\mathcal{L}_{v2t} = -\frac{1}{B} \sum_{i=1}^{B} \log \frac{\exp(\operatorname{sim}(\boldsymbol{v}_i, \boldsymbol{t}_i)/\tau)}{\sum_{k=1}^{N} \exp(\operatorname{sim}(\boldsymbol{v}_i, \boldsymbol{t}_k)/\tau)}$$
(14)

$$\mathcal{L}_{t2v} = -\frac{1}{B} \sum_{i=1}^{B} \log \frac{\exp(\operatorname{sim}(\boldsymbol{v}_i, \boldsymbol{t}_i)/\tau)}{\sum_{k=1}^{N} \exp(\operatorname{sim}(\boldsymbol{v}_k, \boldsymbol{t}_i)/\tau)}$$
(15)

where *B* is the batch size.  $sim(\boldsymbol{x}, \boldsymbol{y}) = \frac{\boldsymbol{x}^{\top} \boldsymbol{y}}{\|\boldsymbol{x}\| \|\boldsymbol{y}\|}$  is the cosine similarity between two vectors.  $\tau$  is a temperature hyper-parameter to scale the logits. Note that though the image-to-text loss  $\mathcal{L}_{v2t}$  is optimized, the pre-trained vision encoder is kept frozen during training. Similar to NLT, the trainable parameters in LE come from the language acquirers and the non-native embedding block.

# 202 **4 Experiments**

#### 203 4.1 Dataset Description

We train our model with the Conceptual Captions (CC) dataset [33] and two translation enhanced versions of the CC [42, 4]. We use Multi30K [10], MSCOCO [6, 24, 38] and XTD [1] for multilingual image-text retrieval evaluation, and MSRVTT [37, 15] for multilingual video-text retrieval evaluation. **Conceptual Captions** (CC) [33] contains 3.3 million image-text pairs in English crawled from

the Web.<sup>1</sup> We also randomly select 300K image-text pairs denoted as **CC300K** for training our 208 model to show the low-cost merit of MLA. For multilingual sentences, we leverage two translation 209 augmented CC datasets: (1) CC6L [42] that translates all English captions of the CC into 5 languages 210 (German(de), French(fr), Czech(cs), Chinese(zh));<sup>2</sup> and (2) CC69L [4] that contains 27K captions in 211 each of the 68 languages translated from English.<sup>3</sup> Considering the languages of the downstream 212 datasets, we train the model with CC6L for multilingual image-text retrieval, and with CC69L for 213 214 multilingual video-text retrieval. Multi30K [10] is built upon Flickr30K [39]. The English(en) captions are manually translated into 215 German(de), French(fr) and Czech(cs). It contains 31K images paired with 5 captions per image in 216 English and German, and 1 caption in French and Czech. We use the standard train, dev and test 217 splits defined in [39]. 218

MSCOCO [6] contains 123K images with 5 English captions per image. [38] annotates 5 Japanese captions per image, and [24] extends MSCOCO with Chinese captions for 20K images. We follow the standard train, dev and test splits for English and Japanese as in [19]. For Chinese, we can only perform zero-shot evaluation on the test split defined in [24], as the full splits have overlaps with

223 English and Japanese splits.

XTD [1] provides captions in 11 languages (English(en), German(de), French(fr), Chinese(zh),
 Japanese(ja), Italian(it), Spanish(es), Russian(ru), Polish(pl), Turkish(tr), Korean(ko)) for 1K
 MSCOCO images. Except for Japanese, all non-English captions are translated from the English
 caption directly. We use this dataset for zero-shot image-text retrieval evaluation only.

MSRVTT [37] is a video caption dataset with 10K videos, where each video is annotated with 20
English captions. Huang et al.[15] translates the English captions into 8 languages (German(de),
French(fr), Russian(ru), Spanish(es), Czech(cz), Swahili(sw), Chinese(zh) and Vietnamese(vi)) via
machine translation service. We follow the standard train/dev splits in [37], and evaluate on the 1K

test split as described in [40].

## **4.2 Implementation Details**

We apply MLA on two VLP models: CLIP-ViT-B-32 and CLIP-ViT-B-16 [30], denoted as MLA<sub>CLIP</sub> 234 and MLA<sub>CLIP16</sub> respectively. The hidden dimension of the language acquirers is set to 256, and 235 all language acquirers for each non-native language cost only 3.14 MB parameters. The non-native 236 embedding matrix is initialized with M-BERT [8]. It costs 92.2 MB and shared with all non-native 237 languages. We train two separate models for multilingual image-text retrieval and video-text retrieval. 238 For the image model, we train with CC6L [42]. For the video model, we use multilingual captions 239 from CC69L [4]. For both models, we optimize multiple language acquirers iteratively with a batch 240 size of 128. The NLT stage performs 117,150 steps with a learning rate of 1e-4, and the LE stage 241 performs 11,715 steps with a learning rate of 3e-6. The temperature  $\tau$  is set to 0.01. For both stages, 242 we use the Adam optimizer [22] with a linear warm-up for the first 10% of steps. The whole training 243 process takes about 12 hours to converge on 1 Nvidia V100 GPU. 244

## 245 4.3 Evaluation on Multilingual Image-Text Retrieval

In multilingual image-text retrieval, models are given a sentence in a certain language to find the most 246 semantically relevant image from an image database and vice versa. We compare our model with 247 state-of-the-art multilingual vision-language pre-training methods under three settings: (1) Zero-shot: 248 we directly evaluate the model without fine-tuning on downstream datasets. (2) Fine-tune on English: 249 we first fine-tune the VLP model on downstream English data. We then insert the language acquirers 250 and non-native embedding block into the fine-tuned model and evaluate on other languages directly. 251 (3) Fine-tune on All: after (2), we fine-tune the language acquirers and non-native embedding block 252 on the downstream dataset and freeze other parts of the model. Following previous works [27, 42, 17], 253 we report Average Recall (AR), which is the average score over Recall@1, Recall@5, and Recall@10 254 on two retrieval directions (image $\rightarrow$ text, text $\rightarrow$ image). The results are shown in Table 1. Also, the 255 comparison of computing costs and parameters can be found in Table 2. 256

<sup>&</sup>lt;sup>1</sup>We can only access  $\sim 2.5$  million images due to some broken URLs.

<sup>&</sup>lt;sup>2</sup>Dataset released at https://github.com/zmykevin/UC2, under MIT license.

 $<sup>^{3}</sup>$ Released at https://github.com/FreddeFrallan/Multilingual-CLIP, under MIT license. We remove captions of unaccessible images, leaving  $\sim 20$ K captions for each language.

Table 1: Multilingual image-text retrieval results on Multi30K and MSCOCO. TrTrain: Translatetrain, FT-En: *Fine-tune on English*, FT-All: *Fine-tune on All*. †: Models trained with publicly unavailable datasets. ‡: Models fine-tuned on COCO-CN [24], which has an overlap train split with the test split of English and Japanese. Best results are in bold and second best are underlined.

	Method Training Data			Mult	i30K		MSCOCO 1K		MSCOCO 5K	
		Training Data	en	de	fr	cs	en	ja	en	ja
	Unicoder-VL	CC3M (English only)	72.0	-	-	-	63.7	-	-	-
	ALIGN	AT-en (English only)	84.3	-	-	-	80.0	-	<u>60.6</u>	-
ot	$M^{3}P$	CC3M+Wiki	57.9	36.8	27.1	20.4	63.1	33.3	-	-
-sh	$UC^2$	TrTrain(CC3M)	66.6	62.5	60.4	55.1	70.9	62.3	-	-
cro	$MKD_{CLIP}$	TrTrain(CC300k)	82.1	77.1	75.2	<u>72.3</u>	78.5	73.6	-	-
Ň	MURAL	TrTrain(CC12M)+EOBT	80.9	76.0	75.7	68.2	78.1	72.5	58.0	49.7
	$MURAL^{\dagger}$	AT+MBT		76.2	75.0	64.6	79.2	73.4	59.5	<u>54.4</u>
	MLA <sub>CLIP</sub>	TrTrain(CC300K)	<u>84.4</u>	<u>78.7</u>	<u>77.7</u>	70.8	79.4	<u>74.9</u>	60.5	54.1
	$MLA_{CLIP16}$	TrTrain(CC300K)	86.4	80.8	80.9	72.9	80.9	76.7	62.6	57.0
	M <sup>3</sup> P	CC3M+Wiki	87.4	82.1	67.3	65.0	88.6	56.0	-	-
Е'n	$UC^2$	TrTrain(CC3M)	87.2	83.8	77.6	74.2	88.1	71.7	-	-
F.	MLA <sub>CLIP</sub>	TrTrain(CC300K)	<u>92.0</u>	82.6	<u>85.1</u>	<u>76.2</u>	<u>89.3</u>	80.4	<u>75.7</u>	<u>62.1</u>
_	$MLA_{\rm CLIP16}$	TrTrain(CC300K)	94.5	86.4	87.3	79.5	91.3	82.6	79.4	65.5
	$M^3 P^{\ddagger}$	CC3M+Wiki	87.7	82.7	73.9	72.2	88.7 <sup>‡</sup>	87.9 <sup>‡</sup>	-	-
_	$UC^{2\ddagger}$	TrTrain(CC3M)	88.2	84.5	83.9	81.2	$88.1^{\ddagger}$	87.5 <sup>‡</sup>	-	-
AI	MURAL	TrTrain(CC12M)+EOBT	91.0	87.3	86.4	82.4	89.4	87.4	73.7	71.9
F.	$MURAL^{\dagger}$	AT+MBT	92.2	88.6	<u>87.6</u>	84.2	88.6	88.4	75.4	74.9
1	MLA <sub>CLIP</sub>	TrTrain(CC300K)	92.0	86.8	85.4	82.3	<u>89.3</u>	88.1	<u>75.7</u>	73.2
	MLA <sub>CLIP16</sub>	TrTrain(CC300K)	94.5	89.7	89.2	85.9	91.3	90.4	79.4	76.5

<sup>257</sup> Under the *Zero-shot* setting, we observe that <sup>258</sup> MLA<sub>CLIP</sub> performs significantly better than

state-of-the-art M-VLP models on English. This
 is because MLA<sub>CLIP</sub> could completely maintain
 the strong English performance of CLIP. In con trast, M-VLP models typically perform worse

than their monolingual counterparts on English (M<sup>3</sup>P 57.9 vs. Unicoder-VL[23] 72.0, MURAL Table 2: Comparison of trainable parameters and computing costs between MLA and M-VLPs.

1 0		
Method	Trainable Params	Computing Costs
M <sup>3</sup> P	566 M	4×V100×7d
$UC^2$	478 M	$8 \times V100 \times 4d$
MURAL	300 M	128×TPUv3×4d
Ours (MLA <sub>CLIP</sub> ) $ $	108 M	$1 \times V100 \times 0.5d$

80.9 vs. ALIGN[18] 84.3). MLA<sub>CLIP</sub> also outperforms M-VLP models on other languages. For 265 example, MLA<sub>CLIP</sub> achieves 78.7 average recall score on German, outperforming MURAL by 266 2.7%. Note that the pre-training dataset of MURAL contains 12 million image-text pairs for each 267 language, while MLA<sub>CLIP</sub> only uses 300K training image-text pairs. It demonstrates that MLA is a 268 269 high-data-efficient method to empower monolingual VLP models with multilingual capability. Under the *Fine-tune on English* setting, MLA shows strong cross-lingual transfer capability. Under the 270 271 *Fine-tune on All* setting, MLA<sub>CLIP</sub> performs slightly worse than MURAL which was pre-trained on publicly unavailable dataset AT+MBT [17]. We consider the reason is that MURAL has more 272 trainable parameters than MLA<sub>CLIP</sub> (300M vs 108M, as shown in Table 2) for fine-tuning, which 273 makes it easier to fit the downstream datasets with a certain scale such as Multi30K and MSCOCO. 274 MLA<sub>CLIP16</sub> achieves state-of-the-art results on all languages under three settings. It indicates that if 275 stronger VLP models such as ALIGN-L2 [18] or Florence [41] are provided, better performance on 276 multilingual image-text retrieval could be reached through MLA. 277

#### 278 4.4 Evaluation on Multilingual Video-Text Retrieval

In multilingual video-text retrieval, the model searches for the most semantically relevant videos 279 given a text query in a certain language. Following [26], we first uniformly sample 12 frames from 280 each video, and use the pre-trained vision encoder to extract representations for each frame. We then 281 perform mean pooling over frame representations to get the video representation. We also evaluate 282 the models under three settings as in Sec.4.3. We report the text $\rightarrow$ video Recall@1 score in Table 283 284 3. Under Zero-shot setting, MLA<sub>CLIP</sub>, which is trained on CC69L without using any video data, achieves comparable or even better results than the fine-tuning results of the state-of-the-art M-VLP 285 model XLM-R-MMP [15] on several languages (de: 20.1 vs. 21.1; fr: 22.0 vs. 21.8; es: 20.2 vs. 286 21.9). Under the *Fine-tune on English* and *Fine-tune on All* settings, MLA<sub>CLIP</sub> also outperforms 287 XLM-R-MMP significantly. We consider the convincing performance comes from two reasons: 1) 288

	Method	en	de	fr	cs	zh	ru	vi	SW	es 1	mean
ZS	$\begin{array}{l} Ours(MLA_{\rm CLIP} \text{ w/o } LE) \\ Ours(MLA_{\rm CLIP}) \end{array}$	30.8 30.8	18.3 <b>20.1</b>	18.9 <b>22.0</b>	14.5 <b>15.7</b>	<b>18.6</b> 18.3	12.6 <b>14.4</b>	7.2 <b>8.2</b>	10.2 10.7	19.3 20.2	16.7 <b>17.8</b>
FT-En	XLM-R-MMP [15]	23.8	19.4	20.7	19.3	18.2	19.1	8.2	8.4	20.4	17.5
	Ours(MLA <sub>CLIP</sub> )	<b>42.5</b>	<b>26.1</b>	<b>26.7</b>	<b>20.5</b>	25.3	<b>18.9</b>	<b>12.9</b>	<b>12.6</b>	27.2	23.6
FT-AII	XLM-R-MMP [15]	23.1	21.1	21.8	20.7	20.0	20.5	10.9	14.4	21.9	19.4
	Ours(MLA <sub>CLIP</sub> )	<b>42.5</b>	<b>33.1</b>	<b>34.5</b>	<b>30.5</b>	<b>31.6</b>	<b>28.9</b>	<b>16.9</b>	<b>24.3</b>	<b>33.5</b>	<b>30.6</b>

Table 3: Multilingual video-text retrieval results on MSRVTT. ZS: Zero-shot, FT-En: Fine-tune on English, FT-All: Fine-tune on All.

CLIP is a strong VLP model that can generalize well on video data. 2) The proposed MLA framework
 can well transfer the open-domain knowledge learned by CLIP to other languages. These results
 suggest that MLA could maintain the open-domain capability of the VLP model which generalizes

292 well on different downstream data.

#### 293 4.5 Ablation Studies

#### 294 A. Training Strategy

We conduct an ablation study in Table 4 to val-295 idate the effectiveness of the proposed MLA 296 training strategy. For those settings with NLT 297 and LE at the same stage, we add the loss of 298 the two objectives together during training. By 299 comparing row 1 to row 2&3, we observe that 300 LE at stage one leads to poor performance. This 301 indicates that aligning with the native language 302

Table 4:	Ab	lation	study	on	training	strategy.

Row	Stage	one	Stage	two	N	/ulti30	ĸ	MSCO	DCO 1K
Row	NLT	LE	NLT	LE	de	fr	cs	ja	zh
1	√				76.3	74.2	67.2	72.1	75.7
2		$\checkmark$			68.2	67.7	58.6	65.9	71.7
3	<ul> <li>✓</li> </ul>	$\checkmark$			71.1	69.7	59.8	67.6	73.9
4	√			$\checkmark$	78.7	77.7	70.8	74.9	78.5
5	√		$\checkmark$	$\checkmark$	78.4	77.3	69.9	74.2	78.1

is more important for the VLP model to acquire new languages at an early stage. It is consistent with the learning habits of humans. By comparing row 1 and row 4, we see that LE at stage two could bring improvements on the new languages. Additionally, comparing row 4 and row 5 suggests that optimizing the model with NLT and LE together at stage two does not bring improvements.

### 307 B. Language Acquirers and Embedding Initialization

In order to validate the effectiveness of the pro-308 posed Language Acquirers, we remove the lan-309 guage acquirers and the M-BERT embedding 310 initialization from the model respectively and 311 evaluate on zero-shot multilingual image-text 312 313 retrieval. As shown in Table 5, the performance on all languages drops significantly without lan-314 guage acquirers. Meanwhile, initializing the 315 embedding with M-BERT [8] only brings incre-316

Table 5: Ablation study on la	nguage acquirers and
embedding initialization. LA	: Language Acquirers,
EI: M-BERT Embedding Initiali	zation

	-					
Mathada	N	Aulti301	MSCOCO 1K			
Wieulous	de	fr	cs	ja	zh	
MLA <sub>CLIP</sub>	78.7	77.7	70.8	74.9	78.5	
MLA <sub>CLIP</sub> w/o LA	76.1	74.9	65.7	70.3	76.5	
$MLA_{\rm CLIP}$ w/o EI	77.9	76.2	69.4	74.6	78.1	

mental improvements. It indicates that the language acquirers contribute most to the performance,

and MLA does not depend much on the initialization of non-native embedding.

## 319 C. Low-resource Languages

Image-text pairs may be rare for low-resource languages. To explore the performance of MLA under this situation, we further simulate a **lowresource scenario** using XTD dataset. We finetune MLA<sub>CLIP</sub> and UC<sup>2</sup> (pre-trained on CC6L) with small amount of data from XTD in an un-

Table 6:	Low resource	performance	on image-
Korean re	etrieval.		

	Methods	Data	Training samples 100 / 200 / 600
1	$UC^2$	Img-Txt	47.0/60.1/78.3
2	MLA <sub>CLIP</sub>	Txt-Txt	51.7 / 62.8 / 78.7
3	MLA <sub>CLIP</sub>	Both	56.7 / 66.9 / 80.1

seen language. We randomly sample 600 pairs for finetuning, and the remained 400 samples are evenly divided for validation and testing. Korean is chosen to perform simulation as its script and language family are not covered by CC6L. Experimental results in Table 6 show that MLA can achieve competitive results with **very small amount of text-text pairs only** (row 2), and adding image-text pairs brings further improvement (row 3). It demonstrates that MLA is still an attractive method for low-resource languages even without any image-text pairs.

#### 332 D. Amount of Training Data

<sup>333</sup> We conduct experiments to control the numbers of image-text pairs used for each language.

We train the models with CC6L and evaluate on MSCOCO 1K and Multi30K under the zero-shot

setting. The corresponding mean AR over non-336 English languages (de, fr, cs, ja, zh) are drawn 337 in Figure 3. We observe that MLA performs 338 significantly better than MKD [31] in all cases. 339 Note that when the amount of training data is 340 small, the advantage of MLA is more obvious, 341 which could outperform MKD even without the 342 LE training stage. Additionally, when training 343 with only 30K image-text pairs per language, 344 MLA outperforms  $UC^2$ , which is pre-trained 345 with 3M pairs per language. MLA is thus a 346 data-efficient method to build multilingual and 347 multimodal models. 348



Figure 3: Mean AR vs. number of image-text pairs per language.

#### 349 E. Language Extensibility

Multilingual models often encounter the need to support new languages that do not occur in the 350 training stage. We conduct language extension experiments to compare MLA<sub>CLIP</sub> with M-VLP 351 model UC<sup>2</sup> [42] on the XTD dataset [1]. XTD supports 11 languages, and 5 of them (en, de, fr, cs, 352 zh, ja) are seen in the pre-training stage of  $UC^2$ , while other 6 languages (it, es, ru, pl, tr, ko) are 353 unseen. To make a fair comparison, we first train  $MLA_{CLIP}$  with the same data as  $UC^2$  and then train 354 both of them on unseen languages with CC69L. The zero-shot image-text retrieval results on XTD 355 are shown in Table 7. We observe a significant performance degeneration on the seen languages for 356 UC<sup>2</sup> when training solely with unseen languages (row 1 vs. row 2). Even keep training with the seen 357 languages, the performance is still significantly reduced due to the limited model capacity (row 1 vs. 358 row 3). In contrast, as MLA decoupled multiple languages through acquirers, the performance of the 359 seen languages is rarely affected (row 4 vs. row 5). This suggests that MLA framework can build 360 multimodal multilingual models that are suitable for supporting increasing numbers of languages.

Dow	Method		Seen languages				Unseen languages					
ROW			de	fr	zh	ja	it	es	ru	pl	tr	ko
1	UC <sup>2</sup> w/o unseen language training	71.8	67.5	68.4	61.9	51.5	-	-	-	-	-	
2	UC <sup>2</sup> w/ unseen language training	63.6	57.8	57.6	57.6	48.4	56.4	56.2	51.3	56.4	51.62	51.3
3	UC <sup>2</sup> w/ all language training	65.2	59.3	59.7	60.1	50.5	57.7	56.5	50.9	55.3	53.2	50.2
4	MLA <sub>CLIP</sub> w/o unseen language training	75.9	72.6	72.9	73.7	67.2	-	-	-	-	-	-
5	$MLA_{\rm CLIP}$ w/ unseen language training	76.0	72.6	72.9	73.8	67.2	64.7	62.8	58.1	63.0	56.5	57.3

Table 7: Language extention experiments on XTD dataset.

361

## **362 5 Conclusion and Limitations**

363 In this paper, we propose the MultiLingual Acquisition (MLA) framework that can empower multilingual capability on monolingual Vision-Language Pre-training models with low-cost and high-364 flexibility. MLA injects language acquirers and a non-native embedding block into VLPs to support 365 new languages. Inspired by the language learning habits of humans, we propose a two-stage training 366 strategy to optimize the language acquirers and non-native embedding block. MLA applied on CLIP 367 368 achieves state-of-the-art performances on multilingual image-text and video-text retrieval benchmarks 369 with much less computing costs and training data. Extensive ablation studies demonstrate that MLA is a flexible, effective, and efficient method to empower multilingual capability on multimodal models. 370

Though MLA has shown high performance in our experiments, it has one limitation that it learns multilingual representations at a coarse grained level. Therefore, our future works include exploring fine-grained alignment between different languages. Furthermore, the majority of our training data is automatically constructed through machine translation, so the ethical prejudice from the machine translation service may potentially affect the behavior of multilingual models produced by MLA. One way to mitigate such concern is to use human annotated or reviewed data for training.

# 377 **References**

- [1] Pranav Aggarwal and Ajinkya Kale. Towards zero-shot cross-lingual image retrieval. *CoRR*, abs/2012.05107, 2020.
- [2] Mikel Artetxe, Sebastian Ruder, and Dani Yogatama. On the cross-lingual transferability of
   monolingual representations. In *58th Annual Meeting of the Association for Computational Linguistics*, pages 4623–4637. Association for Computational Linguistics, 2020.
- [3] Andrea Burns, Donghyun Kim, Derry Wijaya, Kate Saenko, and Bryan A Plummer. Learning
   to scale multilingual representations for vision-language tasks. In *European Conference on Computer Vision*, pages 197–213. Springer, 2020.
- [4] Fredrik Carlsson. Multilingual clip. https://github.com/FreddeFrallan/
   Multilingual-CLIP, 2021.
- [5] Dominic Castello. First language acquisition and classroom language learning: Similarities and
   differences. *ELAL College of Arts & Law*, pages 1–18, 2015.
- [6] Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár,
   and C. Lawrence Zitnick. Microsoft COCO captions: Data collection and evaluation server.
   *CoRR*, abs/1504.00325, 2015.
- [7] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng,
   and Jingjing Liu. Uniter: Universal image-text representation learning. In *European conference* on computer vision. Springer, 2020.
- [8] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of
   deep bidirectional transformers for language understanding. In *NAACL-HLT (1)*, 2019.
- [9] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,
   Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al.
   An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2020.
- [10] Desmond Elliott, Stella Frank, Khalil Sima'an, and Lucia Specia. Multi30k: Multilingual
   english-german image descriptions. In *VL@ ACL*, 2016.
- [11] Hongliang Fei, Tan Yu, and Ping Li. Cross-lingual cross-modal pretraining for multimodal re trieval. In 2021 Conference of the North American Chapter of the Association for Computational
   Linguistics: Human Language Technologies, pages 3644–3650, 2021.
- [12] Spandana Gella, Rico Sennrich, Frank Keller, and Mirella Lapata. Image pivoting for learning
   multilingual multimodal representations. In *EMNLP 2017: Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 2017.
- [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
   recognition. In *IEEE conference on computer vision and pattern recognition*, 2016.
- [14] Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe,
   Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning
   for NLP. In *36th International Conference on Machine Learning*, 2019.
- [15] Po-Yao Huang, Mandela Patrick, Junjie Hu, Graham Neubig, Florian Metze, and Alexander G
   Hauptmann. Multilingual multimodal pre-training for zero-shot cross-lingual transfer of vision language models. In 2021 Conference of the North American Chapter of the Association for
   Computational Linguistics: Human Language Technologies, 2021.
- <sup>419</sup> [16] Yuqi Huo, Manli Zhang, Guangzhen Liu, Haoyu Lu, Yizhao Gao, et al. Wenlan: Bridging <sup>420</sup> vision and language by large-scale multi-modal pre-training, 2021.
- [17] Aashi Jain, Mandy Guo, Krishna Srinivasan, Ting Chen, Sneha Kudugunta, Chao Jia, Yinfei
   Yang, and Jason Baldridge. MURAL: Multimodal, multitask representations across languages.
   In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 3449–3463, 2021.

- [18] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan
   Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning
   with noisy text supervision. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139, pages 4904–4916. PMLR, 2021.
- [19] Andrej Karpathy and Li Fei-Fei. Deep visual-semantic alignments for generating image
   descriptions. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June
   2015.
- [20] Donghyun Kim, Kuniaki Saito, Kate Saenko, Stan Sclaroff, and Bryan Plummer. Mule:
   Multimodal universal language embedding. In *AAAI Conference on Artificial Intelligence*,
   volume 34, pages 11254–11261, 2020.
- [21] Wonjae Kim, Bokyung Son, and Ildoo Kim. Vilt: Vision-and-language transformer without
   convolution or region supervision. In *38th International Conference on Machine Learning*,
   volume 139 of *Proceedings of Machine Learning Research*. PMLR, 2021.
- [22] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR* (*Poster*), 2015.
- [23] Gen Li, Nan Duan, Yuejian Fang, Ming Gong, and Daxin Jiang. Unicoder-vl: A universal
   encoder for vision and language by cross-modal pre-training. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 11336–11344, 2020.
- [24] Xirong Li, Chaoxi Xu, Xiaoxu Wang, Weiyu Lan, Zhengxiong Jia, Gang Yang, and Jieping
   Xu. Coco-cn for cross-lingual image tagging, captioning, and retrieval. *IEEE Transactions on Multimedia*, 21(9):2347–2360, 2019.
- [25] Xiujun Li, Xi Yin, Chunyuan Li, Pengchuan Zhang, Xiaowei Hu, Lei Zhang, Lijuan Wang,
   Houdong Hu, Li Dong, Furu Wei, et al. Oscar: Object-semantics aligned pre-training for
   vision-language tasks. In *European Conference on Computer Vision*. Springer, 2020.
- [26] Huaishao Luo, Lei Ji, Ming Zhong, Yang Chen, Wen Lei, Nan Duan, and Tianrui Li. Clip4clip:
   An empirical study of clip for end to end video clip retrieval. *arXiv preprint arXiv:2104.08860*, 2021.
- [27] Minheng Ni, Haoyang Huang, Lin Su, Edward Cui, Taroon Bharti, Lijuan Wang, Dongdong
   Zhang, and Nan Duan. M3p: Learning universal representations via multitask multilingual
   multimodal pre-training. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
   pages 3977–3986, 2021.
- [28] Jonas Pfeiffer, Gregor Geigle, Aishwarya Kamath, Jan-Martin O Steitz, Stefan Roth, Ivan Vulić,
   and Iryna Gurevych. xgqa: Cross-lingual visual question answering. *arXiv e-prints*, 2021.
- [29] Jonas Pfeiffer, Ivan Vulić, Iryna Gurevych, and Sebastian Ruder. Mad-x: An adapter-based
   framework for multi-task cross-lingual transfer. In 2020 Conference on Empirical Methods in
   Natural Language Processing (EMNLP), pages 7654–7673, 2020.
- [30] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agar wal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya
   Sutskever. Learning transferable visual models from natural language supervision. In *38th International Conference on Machine Learning*, volume 139, pages 8748–8763, 2021.
- [31] Nils Reimers and Iryna Gurevych. Making monolingual sentence embeddings multilingual
   using knowledge distillation. In 2020 Conference on Empirical Methods in Natural Language
   *Processing (EMNLP)*, pages 4512–4525, 2020.
- [32] Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words
   with subword units. In *54th Annual Meeting of the Association for Computational Linguistics* (Volume 1: Long Papers), pages 1715–1725, 2016.
- [33] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A
  cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages
  2556–2565, 2018.

- [34] Yuqing Song, Shizhe Chen, Qin Jin, Wei Luo, Jun Xie, and Fei Huang. Product-oriented
   machine translation with cross-modal cross-lingual pre-training. In *29th ACM International Conference on Multimedia*, 2021.
- [35] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
   Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.
- [36] Jonatas Wehrmann, Douglas M Souza, Mauricio A Lopes, and Rodrigo C Barros. Language agnostic visual-semantic embeddings. In *IEEE/CVF International Conference on Computer Vision*, pages 5804–5813, 2019.
- [37] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. Msr-vtt: A large video description dataset for
   bridging video and language. In *IEEE Conference on Computer Vision and Pattern Recognition* (CVPR), June 2016.
- [38] Yuya Yoshikawa, Yutaro Shigeto, and Akikazu Takeuchi. Stair captions: Constructing a
   large-scale japanese image caption dataset. In *55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, 2017.
- [39] Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. From image descriptions
   to visual denotations: New similarity metrics for semantic inference over event descriptions.
   *Transactions of the Association for Computational Linguistics*, 2:67–78, 2014.
- <sup>493</sup> [40] Youngjae Yu, Jongseok Kim, and Gunhee Kim. A joint sequence fusion model for video <sup>494</sup> question answering and retrieval. In *European Conference on Computer Vision (ECCV)*, 2018.
- [41] Lu Yuan, Dongdong Chen, Yi-Ling Chen, Noel Codella, Xiyang Dai, Jianfeng Gao, Houdong
   Hu, Xuedong Huang, Boxin Li, Chunyuan Li, et al. Florence: A new foundation model for
   computer vision. *arXiv preprint arXiv:2111.11432*, 2021.
- [42] Mingyang Zhou, Luowei Zhou, Shuohang Wang, Yu Cheng, Linjie Li, Zhou Yu, and Jingjing
   Liu. Uc2: Universal cross-lingual cross-modal vision-and-language pre-training. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4155–4165, June 2021.

# 501 Checklist

502	1. For all authors
503 504	(a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
505	(b) Did you describe the limitations of your work? [Yes] See Sec.5
506	(c) Did you discuss any potential negative societal impacts of your work? [Yes] See Sec.5
507 508	<ul><li>(d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]</li></ul>
509	2. If you are including theoretical results
510	(a) Did you state the full set of assumptions of all theoretical results? [N/A]
511	(b) Did you include complete proofs of all theoretical results? [N/A]
512	3. If you ran experiments
513	(a) Did you include the code, data, and instructions needed to reproduce the main experi-
514	mental results (either in the supplemental material or as a URL)? [Yes] Please check
515	the supplementary material for code and model.
516	(b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
517	were chosen)? [Yes] See Sec.4.2
518	(c) Did you report error bars (e.g., with respect to the random seed after running experi-
519	ments multiple times)? [No]
520	(d) Did you include the total amount of compute and the type of resources used (e.g., type
521	of GPUs, internal cluster, or cloud provider)? [Yes] Sec.4.2

522	4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets
523	(a) If your work uses existing assets, did you cite the creators? [Yes]
524	(b) Did you mention the license of the assets? [Yes] See the footnote of page 6.
525	(c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
526	Please check the supplementary material for code and model.
527	(d) Did you discuss whether and how consent was obtained from people whose data you're
528	using/curating? [Yes]
529	(e) Did you discuss whether the data you are using/curating contains personally identifiable
530	information or offensive content? [N/A]
531	5. If you used crowdsourcing or conducted research with human subjects
532	(a) Did you include the full text of instructions given to participants and screenshots, if
533	applicable? [N/A]
534	(b) Did you describe any potential participant risks, with links to Institutional Review
535	Board (IRB) approvals, if applicable? [N/A]
536	(c) Did you include the estimated hourly wage paid to participants and the total amount
537	spent on participant compensation? [N/A]