

# Image-and-Language Understanding from Pixels Only

Michael Tschannen, Basil Mustafa, Neil Houlsby  
Google Research, Brain Team, Zürich

## Abstract

Multimodal models are becoming increasingly effective, in part due to unified components, such as the Transformer architecture. However, multimodal models still often consist of many task- and modality-specific pieces and training procedures. For example, CLIP (Radford et al., 2021) trains independent text and image towers via a contrastive loss. We explore an additional unification: the use of a pure pixel-based model to perform image, text, and multimodal tasks. Our model is trained with contrastive loss alone, so we call it CLIP-Pixels Only (CLIPPO). CLIPPO uses a single encoder that processes both regular images and text rendered as images. CLIPPO performs image-based tasks such as retrieval and zero-shot image classification almost as well as CLIP, with half the number of parameters and no text-specific tower or embedding. When trained jointly via image-text contrastive learning and next-sentence contrastive learning, CLIPPO can perform well on natural language understanding tasks, without any word-level loss (language modelling or masked language modelling), outperforming pixel-based prior work. Surprisingly, CLIPPO can obtain good accuracy in visual question answering, simply by rendering the question and image together. Finally, we exploit the fact that CLIPPO does not require a tokenizer to show that it can achieve strong performance on multilingual multimodal retrieval without modifications.

## 1. Introduction

In recent years, large-scale multimodal training of Transformer-based models has led to improvements in the state-of-the-art in different domains including vision [2, 9, 67–69], language [5, 10], and audio [4]. In particular, in computer vision and image-language understanding, a single large pre-trained model can outperform task-specific expert models [9, 67, 68]. However, large multimodal models often use modality or dataset-specific encoders and decoders, and accordingly lead to involved protocols. For example, such models frequently involve training different parts of the model in separate phases on their respective datasets, with dataset-specific preprocessing, or transferring

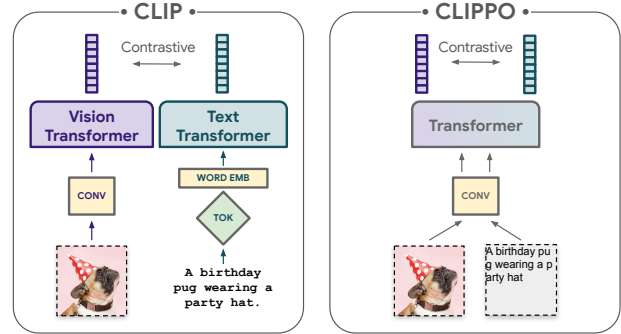


Figure 1. CLIP [50] trains separate image and text encoders, each with a modality-specific preprocessing and embedding, on image/alt-text pairs with a contrastive objective. CLIPPO trains a pure pixel-based model with equivalent capabilities by rendering the alt-text as an image, encoding the resulting image pair using a shared vision encoder (in two separate forward passes), and applying same training objective as CLIP.

different parts in a task-specific manner [68]. Such modality and task-specific components can lead to additional engineering complexity, and poses challenges when introducing new pretraining losses or downstream tasks. Developing a single end-to-end model that can process any modality, or combination of modalities, would be a valuable step for multimodal learning. Here, we focus on images and text.

A number of key unifications have accelerated the progress of multimodal learning. First, the Transformer architecture has been shown to work as a universal backbone, performing well on text [5, 14], vision [15], audio [4, 22, 48], and other domains [6, 32]. Second, many papers have explored mapping different modalities into a single shared embedding space to simplify the input/output interface [20, 21, 42, 63], or develop a single interface to many tasks [29, 35]. Third, alternative representations of modalities allow harnessing in one domain neural architectures or training procedures designed for another domain [26, 45, 48, 54]. For example, [54] and [26, 48] represent text and audio, respectively, by rendering these modalities as images (via a spectrogram in the case of audio).

In this paper, we explore the use of a pure pixel-based model for multimodal learning of text and images. Our

model is a single Vision Transformer [15] that processes visual input, or text, or both together, all rendered as RGB images. The same model parameters are used for all modalities, including low-level feature processing; that is, there are no modality-specific initial convolutions, tokenization algorithms, or input embedding tables. We train our model using only a single task: contrastive learning, as popularized by CLIP [50] and ALIGN [30]. We therefore call our model **CLIP-Pixels Only** (CLIPPO).

We find that CLIPPO performs similarly to CLIP (within 1-2%) on the main tasks CLIP was designed for—image classification and text/image retrieval—despite not having modality-specific towers. Surprisingly, CLIPPO can perform complex language understanding tasks to a decent level without any left-to-right language modelling, masked language modelling, or explicit word-level losses. In particular, on the GLUE benchmark [66] CLIPPO outperforms classic NLP baselines, such as ELMO+BiLSTM+attention, outperforms prior pixel-based masked language models [54], and approaches the score of BERT [14]. Interestingly, CLIPPO can also obtain good performance on VQA when simply rendering the image and text together, despite never having been pre-trained on such data. An immediate advantage of pixel-based models over regular language models is the lack of requirement for pre-determining the vocabulary; consequently, we observe improved performance on multilingual retrieval compared to an equivalent model that uses a classical tokenizer. Finally, we observe that in some circumstances, the previously-observed modality gap [41] is bridged when training CLIPPO.

## 2. Related work

**Multimodal and contrastive pretraining** Most closely related to CLIPPO are CLIP [61] and ALIGN [30] which developed the paradigm of large-scale contrastive training on noisy data from the web. Follow-ups [49, 77] have scaled further and employed state-of-the-art image representation learning to further boost performance.

A number of works have explored model unification via weight-sharing. In the contrastive context, LIMoE [47] explores a one-tower model similar to ours, studying the use of mixture of experts. Outside contrastive training, co-training distinct tasks [1, 42] is a popular strategy, with some approaches [40] involving knowledge distillation and gradient masking. Other works use self-supervised learning algorithms to unify task training [20]. These broadly use discriminative tasks to learn representations for various downstream modalities; generative approaches to multimodal modelling have been scaled to billions of parameters, generating text [2, 9, 67, 74], images [52, 56, 75], videos [25, 65] or audio [4] from (combinations of) various modalities.

Another related domain is document and user interface

(UI) understanding. Corresponding models are trained on diverse multimodal data sets and can usually solve a range of document/UI understanding tasks. Many models rely on text extracted using an off-the-shelf OCR pipeline in combination with document images [3, 27], but models taking only images as an input are getting more popular [33, 38]. While these models can understand visual cues and text from a single input image, they still rely on a tokenized text for training and inference.

**Contrastive training in NLP** There is a sizable body of work on contrastive pretraining on sentence pairs (see [53] for a recent survey), which we explore as an auxiliary objective for CLIPPO. Popular augmentations to generate text pairs involve word deletion, span deletion, reordering, synonym substitution, and next-sentence-prediction [19, 43, 70]. Other methods use different realizations of dropout masks in the model to emulate sentence pairs, or supervised labels to obtain positive and negative pairs [18].

**Visual text and tokenization in NLP** The most closely related method to CLIPPO from the NLP domain is PIXEL [54], which is an MAE [24] trained on rendered text. It obtains strong performance on multilingual syntactic (part-of-speech tagging, dependency parsing) and semantic language understanding (named entity recognition, sentence understanding) tasks, while being more robust to noise in the text than BERT. Other applications for which visual text has been explored include sentiment analysis [62] and machine translation [45, 57].

Visual text side-steps the design and construction of an appropriate tokenizer, which is a large area of research of its own, and can hence simplify text processing in certain—in particular multilingual—scenarios. We refer to [46] for a survey. Popular tokenizer models include WordPiece [14], Byte-Pair Encoding [59], and SentencePiece [36].

Subword-based vocabularies are popular in monolingual setups and usually lead to a good performance trade-off compared to word and character based vocabularies for certain languages including English. In multilingual contexts, appropriately representing the vocabulary of all languages becomes challenging as the number of languages increases [12, 55], which in turn can lead to poor performance in tasks involving underrepresented languages. A variety of mitigation strategies has been developed, for example subword mapping, transliteration, and vocabulary clustering and relocation. We refer to [54, Sec. 5.1] for a more detailed discussion of these strategies.

## 3. Contrastive language-image pretraining with pixels

Contrastive language-image pretraining has emerged as a powerful, scalable paradigm to train versatile vision models on web-scale data sets [50]. Concretely, this approach

relies on image/alt-text pairs which can be automatically collected at large scale from the web. Thereby, the textual descriptions are usually noisy, and can e.g. consist of single keywords, sets of keywords, or potentially lengthy descriptions with many attributes describing the image content. Using this data, two encoders are jointly trained, namely a text encoder embedding the alt-texts and an image encoder embedding the corresponding images into a shared latent space. These two encoders are trained with a contrastive loss, encouraging the embeddings of corresponding images and alt-text to be similar, and at the same time to be dissimilar from all other image and alt-text embeddings.

Once trained, such an encoder pair can be used in many ways: It can be specialized to classifying a fixed set of visual concepts via their textual descriptions (zero-shot classification); the embeddings can be used to retrieve images given a textual description and vice-versa; or the vision encoder can be transferred in supervised fashion to a downstream task by fine-tuning on a labeled data set or by training a head on top of the frozen image encoder representation. In principle, the text encoder can be used as a standalone text embedding, but this application—to our knowledge—has not been explored in-depth, with some authors citing the low quality of the alt-texts leading to weak language modeling performance of the text encoder [61].

Previous works [42, 47] have shown that the image and text encoder can be realized with a single shared transformer model (henceforth referred to as single tower model, or 1T-CLIP), where the images are embedded using a patch embedding, and the tokenized text is embedded using a separate word embedding. Apart from the modality-specific embeddings, all model parameters are shared for the two modalities. While this type of sharing usually leads to a minor performance drop on image/image-language tasks it also halves the number of model parameters.

CLIPPO takes this idea one step further: text inputs are rendered on blank images, and are subsequently dealt with entirely as images, including the initial patch embedding (see Fig. 1 for an illustration). By training this single vision transformer contrastively as prior works, we obtain a single vision transformer model that can understand both images and text through the single interface of vision and provides a single representation which can be used to solve image, image-language, and pure language understanding tasks.

Alongside multimodal versatility, CLIPPO alleviates common hurdles with text processing, namely the development of an appropriate tokenizer and vocabulary. This is particularly interesting in the context of a massively multilingual setup, where the text encoder has to handle dozens of languages.

We find that CLIPPO trained on image/alt-text pairs performs comparably with its 1T-CLIP counterpart on common image and image-language benchmarks, and is com-

petitive with strong baseline language models on the GLUE benchmark [66]. However, due to the low quality of the alt-texts which are often not grammatical sentences, learning language understanding exclusively from alt-texts is fundamentally limited. Therefore, we augment image/alt-text contrastive pretraining with language-based contrastive training. Specifically, we consider positive pairs of consecutive sentences sampled from a text corpus, pairs of translated sentences for different languages, pairs of back-translated sentences, as well as pairs of sentences with word dropout. Such text/text pairs can be seamlessly integrated into the contrastive training by supplementing batches of image/alt-texts with (rendered) text/text pairs.

## 4. Experiments

### 4.1. Training details and models

We rely on a single training setup for all our baselines and visual text models. This setup was tuned to produce good results for standard image/alt-text contrastive training as in [50] (using exactly the same loss function as [50], following the pseudocode in [50, Fig. 3]) and we found that it readily transfers to 1T-CLIP and CLIPPO (including variants with text/text co-training).<sup>1</sup>

Our default architecture is a ViT-B/16 [15] and we perform a subset of experiments with a ViT-L/16 architecture to study the effect of scale (we equip both models a MAP head [37] to pool embeddings). In all cases, the representation dimension used for the contrastive loss is 768. We set the batch size to 10,240 and train the main models for 250k steps, using a minimum 100k training steps for ablations. For models co-trained with a certain percentage of text/text data, we scale the number of iterations such that the number of image/alt-text pairs seen matches the number of iterations of the corresponding model without text/text data (e.g. when 50% of the data is text/text pairs we increase the number of iterations from 250k to 500k). The contrastive loss is computed across the full batch. We use the Adafactor optimizer [60] with a learning rate of  $10^{-3}$  and decoupled weight decay with weight  $10^{-4}$ .

Baseline CLIP-style models are trained using the T5-en SentencePiece tokenizer [51]; we use the abbreviation CLIP\* for the two tower model from [50] trained from scratch using the setup described above, to avoid confusion with the model released by [50]. A sequence length of 196 is used, as this matches the number of visual text “tokens” CLIPPO can process with patch size 16 has at 224px resolution (which we use throughout unless noted otherwise).

**Visual text** For visual text rendering [54, 57] relied on the Google Noto font family<sup>2</sup> which supports the majority

<sup>1</sup>Code will be released as part of the Big Vision code base  
[https://github.com/google-research/big\\_vision](https://github.com/google-research/big_vision).

<sup>2</sup><https://fonts.google.com/noto>

	#param.	training dataset	I1k 10s.	I1k 0s.	C I→T	C T→I	F I→T	F T→I
CLIP*	203M	WebLI	55.8	65.1	48.5	31.3	79.2	59.4
1T-CLIP	118M	WebLI	53.9	62.3	48.0	30.3	77.5	58.2
CLIPPO	93M	WebLI	53.0	61.4	47.3	30.1	76.4	57.3
CLIPPO	93M	WebLI + 25%C4	52.1	57.4	40.7	26.7	68.9	51.8
CLIPPO	93M	WebLI + 50%C4	48.0	53.1	35.2	23.4	64.8	47.2
1T-CLIP L/16	349M	WebLI	60.8	67.8	50.7	32.5	81.0	61.0
CLIPPO L/16	316M	WebLI	60.3	67.4	50.6	33.4	79.2	62.6
CLIPPO L/16	316M	WebLI + 25%C4	60.5	66.0	44.5	29.8	72.9	57.3
CLIPPO L/16	316M	WebLI + 50%C4	56.8	61.7	39.7	27.3	70.1	54.7

Table 1. Vision and vision-language cross-modal results. We report ImageNet-1k 10-shot linear transfer validation accuracy (I1k 10s.), ImageNet-1k zero-shot transfer validation accuracy (I1k 0s.), image-to-text and text-to-image retrieval recall@1 on MS-COCO (C I→T and C T→I) and on Flickr30k (F T→I and F I→T). CLIPPO and 1T-CLIP incur a minor drop in these evaluations compared to CLIP\*, while only using about half of the model parameters. Co-training with text pairs from C4 (models with + xx%C4) degrades performance on some cross-modal tasks (but leads to improved language understanding capabilities, see Table 2).

of Unicode code points. Here, we use the GNU Unifont bitmap font<sup>3</sup>, which has a similar coverage but allows for efficient, lookup-based on-the-fly rendering in our preprocessing pipeline. We emphasize that this rendering strategy does not slow down training compared to tokenizer-based models. In preliminary explorations, we found this to be performance-neutral when compared to the Noto font.

**Image/alt-text data** We use the WebLI data set introduced in [9] which comprises 10 billion images with 12 billion corresponding alt-texts. Importantly, WebLI comprises alt-texts in 109 languages (unlike previous data sets such as LAION-400M [58] which only contain English alt-texts) and it is therefore a great foundation to study multilingual language-image pretraining and its applications. Please refer to [9, Fig. 3] for details on the alt-text language distribution. For English-only models we obtain English versions of non-English alt-texts via GCP Translation API<sup>4</sup>. In addition to alt-text, WebLI also provides OCR annotations, which we do not use in this paper. Finally, WebLI was processed with a de-duplication step removing all images from various splits of the image evaluation sets used in this paper. Please refer to [9, Sec. 3.2] for more details on the WebLI data set and to [9, Appendix B] for a datasheet.

We also present a subset of results based on LAION-400M [58] as an additional comparison point, which can be found in Appendix C.1.

**Text/text data** For co-training with text/text pairs we primarily rely on the publicly available Colossal Clean Crawled Corpus (C4; default/English split) [51]. We randomly sample pairs of consecutive sentences and contrastively train on these pairs, i.e., the model is trained for embedding-based next sentence prediction (NSP) [43]. We also experiment with pairs of parallel sentences in different languages from the WMT19 data set [17] as well as back-

translated English sentences derived from C4 following the strategy described in [11].

## 4.2. Evaluations and metrics

To evaluate the vision and vision/language understanding capabilities of our models we use standard metrics from the literature [47, 50, 77]: “zero-shot” transfer, which uses (embedded) textual description of the classes to be classified/retrieved and compares these with image embeddings. We report the classification accuracy on ImageNet-1k [13] as well as the recall@1 for cross-modal retrieval on MS-COCO [8] and Flickr30k [73]. Furthermore, we test the low-data transfer performance of the models by means of the linear adaptation protocol from [15], reporting the 10-shot accuracy on ImageNet-1k.

We also evaluate CLIPPO and baselines on the popular VQA benchmark VQAv2 [23]. To construct a VQA model using a single pretrained ViT we render the question at the top end of the corresponding image (using the same Unifont renderer as used for CLIPPO training) and follow the standard prediction setup where the answer is predicted as the most likely answer from the training set, i.e. by classification. Specifically, we replace the last layer of our pretrained CLIPPO and baselines with a randomly initialized one with the appropriate number of outputs and fine-tune on VQAv2. This setup tests the ability of the pretrained ViT to combine image and text in intermediate layers as it has produce a single output from a fused image/text input image, unlike in the other cross-modal tasks (and pretraining), where image and text representations are computed with two separate forward passes. Please refer to Appendix A in the supplementary material for examples images with rendered questions and Appendix B.1 for details on the fine-tuning protocol.

Multilingual capabilities are assessed via zero-shot retrieval on CrossModal3600 [64], which is a geographically diverse set comprising 3600 images each human-annotated

<sup>3</sup><http://unifoundry.com/unifont>

<sup>4</sup><https://cloud.google.com/translate>



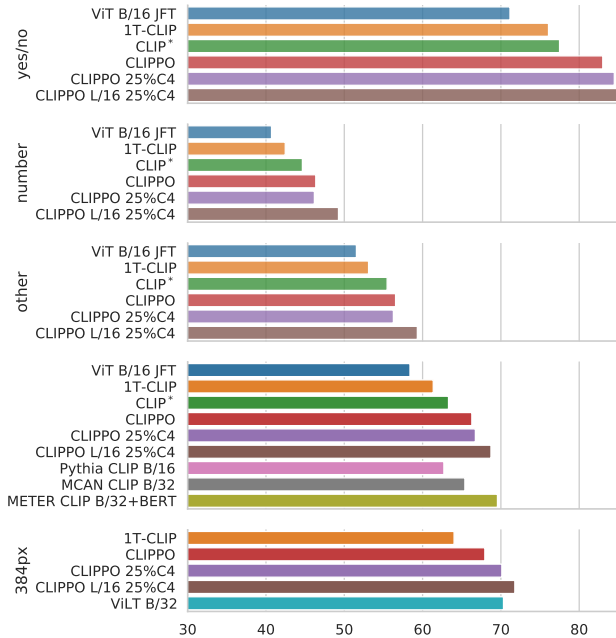


Figure 2. Results on the VQA2 benchmark (test-dev set). In addition to CLIPPO and baselines produced in this work, we also compare to Pythia and MCAN models with ViT encoders from [61], and with comparably sized METER [16] and ViLT [34] models. CLIPPO outperforms CLIP\* and 1T-CLIP clearly on “yes/no” questions and gets similar performance as task-specific models.

with captions in 36 languages.

Finally, we evaluate the language understanding capabilities on the General Language Understanding Evaluation (GLUE) benchmark [66] which comprises natural language inference tasks (MNLI, QNLI, RTE), a sentiment analysis task (SST-2), sentence similarity tasks (QQP, STS-B, MRPC), and a linguistic acceptability task (CoLA). Following common practice, we exclude the WNLI task from the benchmark [14, 70]. We transfer our baselines and CLIPPO models by attaching a 2-hidden layer MLP with 768 units to their representation and following precisely the fine-tuning protocol from BERT [14]. For sentence pair classification tasks we simply render both sentences on the same image, printing [SEP] to mark the start of the second sentence.

### 4.3. Vision and vision-language understanding

**Image classification and retrieval** Table 1 shows the performance of CLIPPO along with the baseline models on the benchmarks described in Sec. 4.2. It can be seen that the CLIPPO and 1T-CLIP incur a drop of a 2-3 percentage points absolute compared to CLIP\*. This is not surprising and can be attributed to the fact that single tower models only have about half the number of parameters of a corresponding two tower model. The difference in per-

formance between the English-only CLIPPO and 1T-CLIP is very small for a B/16 backbone at 100k training steps (see Table 6 in the supplementary material), and vanishes with longer training and/or by increasing the model size, despite the fact that CLIPPO has 25% and 10% fewer parameters than 1T-CLIP for a B/16 and L/16 architecture, respectively (which is due to the absence of the text embedding in CLIPPO).

The multilingual CLIPPO model performs somewhat worse than the corresponding 1T-CLIP, and the gap does not close completely when training longer (see Table 6). However, when evaluated across a broad set of languages on the CrossModal3600 CLIPPO performs on par with or slightly better than 1T-CLIP (see Sec. 4.4 below).

As we add sentence pairs to the training mix the performance on the cross-modal retrieval metrics decreases. This is not surprising as we keep the total batch size constant so that the effective batch size of image/alt-text contrastive training decreases, which is known to impact performance [77]. Interestingly, the 10-shot transfer performance does not move in tandem, but only decreases significantly when half of the training data is sentence pairs. In exchange, co-training with text data leads to significantly improved language understanding performance (Sec. 4.5).

**VQA** In Fig. 2 we report the VQA2 score of our models and baselines. It can be seen that CLIPPO outperforms CLIP\*, 1T-CLIP, as well as a pretrained ViT-B/16 from [15] by a significant margin, achieving a score of 66.3, and co-training with 25% C4 data leads to a slight improvement of the score. The improved score of CLIPPO is mainly due to better performance in “yes/no” questions. Increasing the model size to L/16 adds another 2 points which originate from improvements in the “number” and “other” VQA2 categories. However, note that for an L/16 architecture 1T-CLIP performs competitively with CLIPPO (see Table 7). One possible explanation for this could be that 1T-CLIP develops better OCR capabilities thanks to the higher model capacity (alt-texts can correlate with text in images/scene text, see [9, Fig. 3]). Increasing the resolution to 384px adds 2 to 3 points across models.

We also compare CLIPPO with baselines from the literature. Specifically, [16] proposes framework (called METER) for multimodal tasks, where pretrained transformer-based image and text encoders are combined with a transformer-based fusion module. CLIPPO L/16 achieves performance competitive with their model combining a CLIP B/32 vision backbone with a BERT-Base language backbone, which is roughly comparable in size and computational cost with our L/16 models. Another related work is [61], which combines different CLIP vision backbones with two existing VQA systems, Pythia [31] and MCAN [76]. CLIPPO outperforms different CLIP ViT-based Pythia and MCAN models from [61]. Note, how-

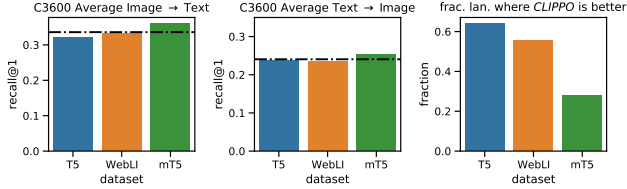


Figure 3. Zero-shot image/text retrieval performance on CrossModal3600 [64]. Although specialized (mc4) tokenizers can be leveraged to improve multilingual performance CLIPPO (dashed black line) broadly matches or exceeds comparable 1T-CLIP models trained with vocabulary size 32,000 (the word embeddings result in a 27% increase in parameter count compared to CLIPPO).

ever, that ResNet-based CLIP backbones lead to better results when combined with these systems. We further note that both [16] and [61] also investigate training their models on a mix of different image-text data sets with multiple objectives such as grounded masked language modeling and text-image matching, before transferring to the VQA task, which leads to significant improvements. ViLT [34] relies on such a strategy to train a single transformer backbone jointly encoding image and text tokens. At 384px resolution, CLIPPO (with 25% C4 data) obtains a VQA score comparable with that of ViLT (and other models from the literature such as ViLBERT [44], VisualBERT [39], and PixelBERT [28]), despite only using a contrastive objective for pretraining.

**The impact of weight sharing** The fact that CLIPPO 1) uses a shared patch embedding for regular images and text images and 2) this embedding has considerably fewer parameters than the text embedding of 1T-CLIP and CLIP\* provokes the question of whether CLIPPO could benefit from separate patch embeddings for text images and regular images. Further, CLIPPO relies on a single head to compute the output representation for images and text, and relaxing this constraint by using separate heads for the two modalities could lead to more expressive representations. To better understand whether a shared patch embedding and/or shared heads are degrading the performance of CLIPPO we train different models with separate embeddings and/or heads. The results (deferred to Appendix C.3) show that neither of these variants lead to improved image classification or retrieval metrics compared to CLIPPO.

#### 4.4. Multilingual vision-language understanding

For typical language models, tokenizer choice can be a challenging process [71]. Commonly used English-language tokenizers generalize poorly to non-latin scripts [77]. This can be alleviated by the use of larger, multilingual vocabularies, at the expense of very large parameter counts. CLIPPO bypasses this issue, removing any

language-related bias stemming from unbalanced or restrictive tokenizers. We consider multilingual image/text retrieval on Crossmodal3600 and compare CLIPPO, trained on WebLI with multilingual alt-texts, against 1T-CLIP with a number of SentencePiece tokenizers; one trained from 300M WebLI multilingual alt-texts, English (T5-en) and multilingual (T5-all) tokenizers from T5 [51], and a multilingual tokenizer (mT5) from mT5 [72], all with a vocabulary size of 32,000. The results are shown in Fig. 3. On average, CLIPPO achieves comparable retrieval performance to these baselines. In the case of mT5, the use of extra data to create the specialized vocabulary can boosts performance above that of CLIPPO; the leveraging of such extra parameters and data in the multilingual context will be an interesting future direction for CLIPPO.

#### Tokenization efficiency

If a tokenizer is well suited to a particular dataset, it will tokenize to shorter sequences—this is especially the case when byte fallback [36] is enabled. SentencePiece

tokenizers have the advantageous ability to tokenize entire—possibly quite long—words to single tokens. CLIPPO cannot learn any such compression, but benefits from equal treatment of all languages and words: it will by definition generalize

equally well to all data, as its tokenization schema has not been trained on a specific dataset. We analyze 20,000 samples for each of the 104 C4 languages. Each CLIPPO token is assumed to be a  $16 \times 16$  patch; though in typical computations all approaches considered here would pad to a fixed length, we compute CLIPPO’s sequence length according to the last patch which contains rendered text. Fig. 4 shows the fraction of C4 languages where CLIPPO processes tokens more efficiently than the vocabularies discussed above. We conservatively define “more efficient” as producing a shorter token sequence for over 75% of examples. Even so, CLIPPO is indeed more efficient across the majority of languages. Per-language breakdowns of multilingual retrieval performance and tokenization efficiency are further discussed in Appendix C.4.

#### 4.5. Language understanding

Table 2 shows the GLUE benchmark results of CLIPPO and baselines. One can observe that CLIPPO trained on We-

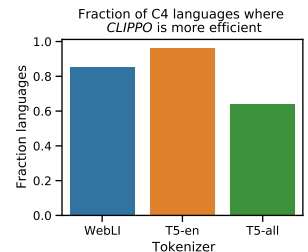


Figure 4. Tokenization efficiency analyzed in terms of the sequence length produced by a given method. CLIPPO produces smaller sequences for the majority of languages compared to 1T-CLIP with alternative tokenizers.

	training dataset	MNLI-M/MM	QQP	QNLI	SST-2	COLA	STS-B	MRPC	RTE	avg
BERT-Base	Wiki + BC	84.0 / 84.0	87.6	91.0	92.6	60.3	88.8	90.2	69.5	83.1
PIXEL	Wiki + BC	78.1 / 78.1	84.5	87.8	89.6	38.4	81.1	88.2	60.5	76.3
BiLSTM		66.7 / 66.7	82.0	77.0	87.5	17.6	72.0	85.1	58.5	68.1
BiLSTM+Attn, ELMo		72.4 / 72.4	83.6	75.2	91.5	44.1	56.1	82.1	52.7	70.0
CLIP* img enc.	WebLI	66.4 / 66.4	78.6	69.4	78.6	0.0	5.2	81.2	52.7	55.5
CLIP* text enc.	WebLI	71.8 / 71.8	82.7	73.0	86.2	6.6	65.0	81.4	53.8	65.9
1T-CLIP text enc.	WebLI	72.6 / 72.6	83.8	80.7	84.9	0.0	79.6	83.3	57.0	68.3
CLIPPO	WebLI	73.0 / 73.0	84.3	81.2	86.8	1.8	80.5	84.1	53.4	68.6
CLIPPO	WebLI + 25%C4	77.7 / 77.7	85.3	83.1	90.9	28.2	83.4	84.5	59.2	74.4
CLIPPO	WebLI + 50%C4	79.2 / 79.2	86.4	84.2	92.9	38.9	83.4	84.8	59.9	76.6
CLIPPO	C4	79.9 / 79.9	86.7	85.2	93.3	50.9	84.7	86.3	58.5	78.4
CLIPPO L/16	WebLI + 25%C4	76.6 / 76.6	87.1	79.9	93.2	48.2	84.1	84.6	56.0	76.1
CLIPPO L/16	WebLI + 50%C4	82.3 / 82.3	87.9	86.7	94.2	55.3	85.8	85.9	59.2	80.0

Table 2. Results for the GLUE benchmark (dev set). The metric is accuracy except for the performance on QQP and MRPC, which is measured using the  $F_1$  score, CoLA which uses Matthew’s correlation, and STS-B which evaluated based on Spearman’s correlation coefficient. “avg” corresponds to the average across all metrics. The results for BERT-Base and PIXEL are from [54, Table 3], and BiLSTM and BiLSTM+Attn, ELMo from [66, Table 6]. All encoders considered here have a transformer architecture comparable to BERT-Base (up to the text embedding layer), except for CLIPPO L/16 which uses a ViT L/16, and the two BiLSTM model variants. Wiki and BC stand for (English) Wikipedia and Bookcorpus [78] data, respectively.

bLI performs competitively with the BiLSTM+Attn+ELMo baseline which relies on deep word embeddings trained on a large language corpus. Also, it can be seen that CLIPPO along with 1T-CLIP outperform the language encoder trained using standard contrastive language vision pretraining (CLIP\*). This indicates that multimodal training in a single encoder benefits language understanding. Furthermore, CLIPPO achieves a much higher GLUE score than the CLIP\* image encoder, which in turn leads to significantly better results than fine-tuning a ViT-B/16 from scratch on GLUE (see Appendix C.2 for additional results). Unsurprisingly, the models pretrained on WebLI cannot do better than random guessing on the CoLA evaluation which requires to assess the grammatical correctness of sentences (recall that alt-texts are rarely grammatical sentences). Also the accuracy of CLIP\* and 1T-CLIP vision encoders we observe for SST-2 is in agreement with what was reported in [50, Table 10] for CLIP with a ViT-B/16 image encoder.

Adding sentence pairs from the C4 corpus gradually improves the GLUE score, and when half of the examples are sentence pairs our model becomes competitive with PIXEL, while still retaining decent image and vision-language understanding capabilities (cf. Table 1). Note, however, that there is a trade-off between language-only tasks and tasks that involve image understanding. Finally, training CLIPPO only on sentence pairs leads to a model which outperforms PIXEL by a significant margin, albeit our model has seen more sentence pairs than PIXEL, so PIXEL might improve as well when training longer.

To corroborate that contrastive NSP is a sensible objective to improve language understanding in the context of CLIPPO, we train CLIPPO without any image/alt-text data

on pairs of parallel translated sentences (this is straightforward in our framework since visual text is language-agnostic), as well as English back-translated data, and evaluate the resulting text representations on GLUE. Table 3 shows that NSP on C4 clearly achieves the highest GLUE score.

	WMT19	WMT19 BT	C4 NSP
GLUE score	61.2	66.6	77.6

Table 3. Ablation of text pair-based contrastive co-training tasks: Training on parallel translated sentences (WMT19), training on parallel back-translated sentences (WMT19 BT), and NSP for sentences sampled from C4 (C4 NSP). C4 NSP leads to the highest GLUE score by a large margin.

## 4.6. Modality gap analysis

Liang et al. [41] discovered that text and image embeddings of CLIP-style models form two distinct clusters rather than both filling the embedding space densely and occupying the same spatial region. They attribute this phenomenon to a combination of initialization conditions and properties of the loss function/training dynamics. Since we consider single tower models here, and also co-train some of these models with text-only pairs it is interesting to see how this affects the modality gap. We compute the gap and visualize it following the recipe from [41] in Fig. 5. CLIPPO attains a slightly lower modality gap than CLIP\*, but clearly features a clustering structure for image and text embeddings. However, when training contrastively with sentence pairs in addition to image/alt-text pairs, the clustering structure

disappears, the image and text embeddings overlap, and the modality gap decreases significantly. A possible explanation for this behavior could be that the additional learning pressure induced by the contrastive loss on sentence pairs encourages text embeddings to spread out more and hence the structure of all embeddings changes.

We refer to the supplementary material for an extended analysis of the modality gap as well as a visualization of the patch embedding of our models and baselines.

## 5. Discussion and limitations

We proposed and evaluated CLIPPO which produces a single ViT that can understand images and language jointly using images as a sole input modality. CLIPPO matches the performance of the 1T-CLIP baseline across many of the considered tasks, and only incurs a minor drop compared to the CLIP\* baseline, in particular taking into account the fact that it has less than half the parameters of a comparable CLIP\*. Nevertheless, several limitations remain, as discussed next.

First, to achieve language understanding performance competitive with PIXEL and BERT on GLUE, contrastive co-training with text pairs is necessary. While adding 25% C4 data to the batch seems to strike a good balance across all tasks considered, it does induce a non-negligible drop in zero-shot image classification and image/text retrieval. This drop becomes more severe as the fraction of C4 examples increases. We observed an associated change in modality gap, and further investigation of the representation in the context of co-training might help to develop models that achieve better overall performance in the co-training setup.

CLIPPO currently relies on cleanly rendered text as an input and its capabilities to handle text from documents or web pages without further adaption is limited (besides the basic OCR capabilities that CLIP-style models learn from image/alt-text pairs). We emphasize that sophisticated OCR and document understanding is not a goal of this paper, and training CLIPPO with augmented noisy rendered text that mimics the distribution of documents and websites is likely to lead to worse performance across the considered tasks, since image/alt-text pairs are less correlated and provide a weaker learning signal. However, developing CLIPPO further to handle less clean visual text will open many additional applications.

CLIPPO, like CLIP, BERT, and PIXEL and many other models, uses an encoder-only design and hence lacks the ability to generate text outputs. A common approach to equip encoder-only models with generation capabilities (e.g., for image captioning or VQA) is to simply combine them with a (potentially pretrained) language model [7, 69]. This approach naturally also applies to CLIPPO and PIXEL, but defeats the advantages of visual text in certain (e.g. multilingual) scenarios. While visual text outputs have previ-

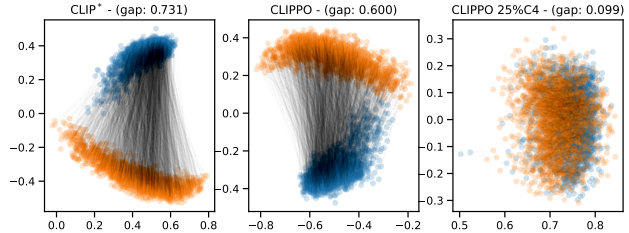


Figure 5. Visualization of the modality gap for CLIP\* and CLIPPO optionally trained with 25% C4 data. The visualization follows the analysis from [41] and shows embedded images (blue dots) and corresponding alt-text (orange dots) from the WebLI validation set, projected to the first two principal components of the validation data matrix. CLIPPO has a slightly smaller modality gap than CLIP\*; co-training with C4 data strongly reduces the gap.

ously been explored by [45] in the context of machine translation, it seems unclear what a scalable tokenizer-free way to generate text is.

Finally, we showed that CLIPPO obtains strong multilingual image/text retrieval performance without requiring the development of an appropriate tokenizer. For fine-grained adjustment and balancing of the retrieval performance further steps will be necessary, including data balancing and potentially co-training with multi-lingual text data. Furthermore, similar to PIXEL, CLIPPO relies on certain ad-hoc design choices w.r.t. the visual representation, for example the left-to-right rendering of Arabic scripts. This approach leads to decent performance on average, but it is unclear what kind of unwanted effects it could be inducing and how these could be mitigated.

## 6. Conclusion

We introduced CLIPPO, a joint model for processing image and text through the lens of vision. This reduces design choices and parameter count, can improve language understanding, and increases generality across multiple languages. We also explored methods of enhancing language understanding, where traditional image/text contrastive models trained on web data fall short (e.g. poor grammatical understanding on CoLA). We demonstrate this is possible by co-training with text pairs, with CLIPPO models outperforming strong NLP baselines while maintaining solid image understanding capabilities.

Although we present a unified contrastive training algorithm, CLIPPO suffers somewhat when co-training on multiple tasks, and future work aiming to further harmonize the co-training setup to ameliorate the trade-off would enhance the models significantly. Deeper understanding of the design choices made in rendering text as images, and the impact on performance, is another interesting avenue to explore.



**Acknowledgments** We would like to thank Lucas Beyer, Josip Djolonga, Alexander Kolesnikov, Mario Lucic, Andreas Steiner, Xiao Wang, and Xiaohua Zhai for inspiring discussions and helpful feedback on this project. We also thank Jeffrey Sorensen for help with the text rendering pre-processing.

## References

- [1] Hassan Akbari, Liangzhe Yuan, Rui Qian, Wei-Hong Chuang, Shih-Fu Chang, Yin Cui, and Boqing Gong. VATT: transformers for multimodal self-supervised learning from raw video, audio and text. In *NeurIPS*, 2021. [2](#)
- [2] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andrew Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. Flamingo: A visual language model for few-shot learning. *CoRR*, abs/2204.14198, 2022. [1](#), [2](#)
- [3] Srikanth Appalaraju, Bhavan Jasani, Bhargava Urala Kota, Yusheng Xie, and R. Manmatha. DocFormer: End-to-end transformer for document understanding. In *ICCV*, pages 973–983, 2021. [2](#)
- [4] Zalán Borsos, Raphaël Marinier, Damien Vincent, Eugene Kharitonov, Olivier Pietquin, Matthew Sharifi, Olivier Teboul, David Grangier, Marco Tagliasacchi, and Neil Zeghidour. Audioldm: a language modeling approach to audio generation. *CoRR*, abs/2209.03143, 2022. [1](#), [2](#)
- [5] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, et al. Language models are few-shot learners. In *NeurIPS*, 2020. [1](#)
- [6] Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Michael Laskin, Pieter Abbeel, Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence modeling. In *NeurIPS*, pages 15084–15097, 2021. [1](#)
- [7] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey E. Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, 2020. [8](#)
- [8] 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. [4](#)
- [9] Xi Chen, Xiao Wang, Soravit Changpinyo, A. J. Piergiovanni, Piotr Padlewski, Daniel Salz, Sebastian Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer, Alexander Kolesnikov, Joan Puigcerver, Nan Ding, Keran Rong, Hassan Akbari, Gaurav Mishra, Linting Xue, Ashish Thapliyal, James Bradbury, Weicheng Kuo, Mojtaba Seyedhosseini, Chao Jia, Burcu Karagol Ayan, Carlos Riquelme, Andreas Steiner, Anelia Angelova, Xiaohua Zhai, Neil Houlsby, and Radu Soricut. PaLI: A jointly-scaled multilingual language-image model. *CoRR*, abs/2209.06794, 2022. [1](#), [2](#), [4](#), [5](#)
- [10] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. PaLM: Scaling language modeling with pathways. *CoRR*, abs/2204.02311, 2022. [1](#)
- [11] Geoffrey Cideron, Sertan Girgin, Anton Raichuk, Olivier Pietquin, Olivier Bachem, and Léonard Hussenot. vec2text with round-trip translations. *CoRR*, abs/2209.06792, 2022. [4](#)
- [12] Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. Unsupervised cross-lingual representation learning at scale. In *ACL*, pages 8440–8451, 2020. [2](#)
- [13] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009. [4](#)
- [14] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT*, 2019. [1](#), [2](#), [5](#)
- [15] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16×16 words: Transformers for image recognition at scale. In *ICLR*, 2021. [1](#), [2](#), [3](#), [4](#), [5](#), [14](#), [16](#), [17](#), [22](#)
- [16] Zi-Yi Dou, Yichong Xu, Zhe Gan, Jianfeng Wang, Shuohang Wang, Lijuan Wang, Chenguang Zhu, Pengchuan Zhang, Lu Yuan, Nanyun Peng, Zicheng Liu, and Michael Zeng. An empirical study of training end-to-end vision-and-language transformers. In *CVPR*, pages 18145–18155, 2022. [5](#), [6](#), [14](#), [17](#)
- [17] Wikimedia Foundation. ACL 2019 Fourth Conference on Machine Translation (WMT19), Shared Task: Machine Translation of News. <http://www.statmt.org/wmt19/translation-task.html>. [4](#)

- [18] Tianyu Gao, Xingcheng Yao, and Danqi Chen. Simcse: Simple contrastive learning of sentence embeddings. In *EMNLP*, pages 6894–6910, 2021. [2](#)
- [19] John M. Giorgi, Osvald Nitski, Bo Wang, and Gary D. Bader. DeCLUTR: Deep contrastive learning for unsupervised textual representations. In *ACL/IJCNLP*, pages 879–895, 2021. [2](#)
- [20] Rohit Girdhar, Alaaeldin El-Nouby, Mannat Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan Misra. OmniMAE: Single model masked pretraining on images and videos. *CoRR*, abs/2206.08356, 2022. [1](#), [2](#)
- [21] Rohit Girdhar, Mannat Singh, Nikhila Ravi, Laurens van der Maaten, Armand Joulin, and Ishan Misra. Omnivore: A single model for many visual modalities. In *CVPR*, pages 16081–16091, 2022. [1](#)
- [22] Yuan Gong, Yu-An Chung, and James R. Glass. AST: audio spectrogram transformer. In *Interspeech*, pages 571–575, 2021. [1](#)
- [23] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the V in VQA matter: Elevating the role of image understanding in visual question answering. In *CVPR*, 2017. [4](#), [13](#)
- [24] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross B. Girshick. Masked autoencoders are scalable vision learners. In *CVPR*, pages 15979–15988, 2022. [2](#)
- [25] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey A. Gritsenko, Diederik P. Kingma, Ben Poole, Mohammad Norouzi, David J. Fleet, and Tim Salimans. Imagen video: High definition video generation with diffusion models. *CoRR*, abs/2210.02303, 2022. [2](#)
- [26] Po-Yao Huang, Hu Xu, Juncheng Li, Alexei Baevski, Michael Auli, Wojciech Galuba, Florian Metze, and Christoph Feichtenhofer. Masked autoencoders that listen. *CoRR*, abs/2207.06405, 2022. [1](#)
- [27] Yupan Huang, Tengchao Lv, Lei Cui, Yutong Lu, and Furu Wei. Layoutlmv3: Pre-training for document AI with unified text and image masking. In *ACMMM*, pages 4083–4091, 2022. [2](#)
- [28] Zhicheng Huang, Zhaoyang Zeng, Bei Liu, Dongmei Fu, and Jianlong Fu. Pixel-BERT: Aligning image pixels with text by deep multi-modal transformers. *CoRR*, abs/2004.00849, 2020. [6](#)
- [29] Andrew Jaegle, Sebastian Borgeaud, Jean-Baptiste Alayrac, Carl Doersch, Catalin Ionescu, David Ding, Skanda Kopula, Daniel Zoran, Andrew Brock, Evan Shelhamer, Olivier J. Hénaff, Matthew M. Botvinick, Andrew Zisserman, Oriol Vinyals, and João Carreira. Perceiver IO: A general architecture for structured inputs & outputs. In *ICLR*, 2022. [1](#)
- [30] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc V. Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *ICML*, 2021. [2](#)
- [31] Yu Jiang, Vivek Natarajan, Xinlei Chen, Marcus Rohrbach, Dhruv Batra, and Devi Parikh. Pythia v0.1: The winning entry to the VQA challenge 2018. *CoRR*, abs/1807.09956, 2018. [5](#)
- [32] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Židek, Anna Potapenko, Alex Bridgland, Clemens Meyer, Simon A A Kohl, Andrew J Ballard, Andrew Cowie, Bernardino Romera-Paredes, Stanislav Nikolov, Rishub Jain, Jonas Adler, Trevor Back, Stig Petersen, David Reiman, Ellen Clancy, Michal Zielinski, Martin Steinegger, Michalina Pacholska, Tamas Berghammer, Sebastian Bodenstein, David Silver, Oriol Vinyals, Andrew W Senior, Koray Kavukcuoglu, Pushmeet Kohli, and Demis Hassabis. Highly accurate protein structure prediction with AlphaFold. *Nature*, 596(7873):583–589, 2021. [1](#)
- [33] Geewook Kim, Teakgyu Hong, Moonbin Yim, JeongYeon Nam, Jinyoung Park, Jinyeong Yim, Wonseok Hwang, Sangdoo Yun, Dongyoon Han, and Seunghyun Park. OCR-Free document understanding transformer. In *ECCV*, pages 498–517, 2022. [2](#)
- [34] Wonjae Kim, Bokyung Son, and Ildoo Kim. ViLT: Vision-and-language transformer without convolution or region supervision. In *ICML*, pages 5583–5594. PMLR, 2021. [5](#), [6](#), [17](#)
- [35] Alexander Kolesnikov, André Susano Pinto, Lucas Beyer, Xiaohua Zhai, Jeremiah Harmsen, and Neil Houlsby. UViM: A unified modeling approach for vision with learned guiding codes. In *NeurIPS*, 2022. [1](#)
- [36] Taku Kudo and John Richardson. SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In *EMNLP*, 2018. [2](#), [6](#)
- [37] Juho Lee, Yoonho Lee, Jungtaek Kim, Adam R. Kosiorek, Seungjin Choi, and Yee Whye Teh. Set transformer: A framework for attention-based permutation-invariant neural networks. In *ICML*, pages 3744–3753, 2019. [3](#), [14](#)
- [38] Kenton Lee, Mandar Joshi, Iulia Turc, Hexiang Hu, Fangyu Liu, Julian Eisenschlos, Urvashi Khandelwal, Peter Shaw, Ming-Wei Chang, and Kristina Toutanova. Pix2struct: Screenshot parsing as pretraining for visual language understanding. *CoRR*, abs/2210.03347, 2022. [2](#)
- [39] Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. VisualBERT: A simple and performant baseline for vision and language. *CoRR*, abs/1908.03557, 2019. [6](#)
- [40] Qing Li, Boqing Gong, Yin Cui, Dan Kondratyuk, Xianzhi Du, Ming-Hsuan Yang, and Matthew Brown. Towards a unified foundation model: Jointly pre-training transformers on unpaired images and text. *CoRR*, abs/2112.07074, 2021. [2](#)
- [41] Weixin Liang, Yuhui Zhang, Yongchan Kwon, Serena Yeung, and James Zou. Mind the gap: Understanding the modality gap in multi-modal contrastive representation learning. In *NeurIPS*, 2022. [2](#), [7](#), [8](#), [20](#), [21](#)
- [42] Valerii Likhoshesterov, Anurag Arnab, Krzysztof Choromanski, Mario Lucic, Yi Tay, Adrian Weller, and Mostafa Dehghani. Polyvit: Co-training vision transformers on images, videos and audio. *CoRR*, abs/2111.12993, 2021. [1](#), [2](#), [3](#)
- [43] Lajanugen Logeswaran and Honglak Lee. An efficient framework for learning sentence representations. In *ICLR*, 2018. [2](#), [4](#)
- [44] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. ViL-BERT: Pretraining task-agnostic visiolinguistic representa-

- tions for vision-and-language tasks. In *NeurIPS*, pages 13–23, 2019. [6](#)
- [45] Elman Mansimov, Mitchell Stern, Mia Xu Chen, Orhan Firat, Jakob Uszkoreit, and Puneet Jain. Towards end-to-end in-image neural machine translation. *CoRR*, abs/2010.10648, 2020. [1](#), [2](#), [8](#)
- [46] Sabrina J. Mielke, Zaid Alyafeai, Elizabeth Salesky, Colin Raffel, Manan Dey, Matthias Gallé, Arun Raja, Chenglei Si, Wilson Y. Lee, Benoît Sagot, and Samson Tan. Between words and characters: A brief history of open-vocabulary modeling and tokenization in NLP. *CoRR*, abs/2112.10508, 2021. [2](#)
- [47] Basil Mustafa, Carlos Riquelme, Joan Puigcerver, Rodolphe Jenatton, and Neil Houlsby. Multimodal contrastive learning with LIMO: the language-image mixture of experts. In *NeurIPS*, 2022. [2](#), [3](#), [4](#)
- [48] Daisuke Niizumi, Daiki Takeuchi, Yasunori Ohishi, Noboru Harada, and Kunio Kashino. Masked spectrogram modeling using masked autoencoders for learning general-purpose audio representation. *CoRR*, abs/2204.12260, 2022. [1](#)
- [49] Hieu Pham, Zihang Dai, Golnaz Ghiasi, Kenji Kawaguchi, Hanxiao Liu, Adams Wei Yu, Jiahui Yu, Yi-Ting Chen, Minh-Thang Luong, Yonghui Wu, Mingxing Tan, and Quoc V. Le. Combined scaling for open-vocabulary image classification. *CoRR*, abs/2111.10050, 2021. [2](#)
- [50] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *ICML*, 2021. [1](#), [2](#), [3](#), [4](#), [7](#), [14](#)
- [51] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21:140:1–140:67, 2020. [3](#), [4](#), [6](#), [19](#)
- [52] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *ICML*, pages 8821–8831, 2021. [2](#)
- [53] Nils Rethmeier and Isabelle Augenstein. A primer on contrastive pretraining in language processing: Methods, lessons learned and perspectives. *ACM Computing Surveys*, 2021. [2](#)
- [54] Phillip Rust, Jonas F. Lotz, Emanuele Bugliarello, Elizabeth Salesky, Miryam de Lhoneux, and Desmond Elliott. Language modelling with pixels. *CoRR*, abs/2207.06991, 2022. [1](#), [2](#), [3](#), [7](#), [18](#)
- [55] Phillip Rust, Jonas Pfeiffer, Ivan Vulic, Sebastian Ruder, and Iryna Gurevych. How good is your tokenizer? on the monolingual performance of multilingual language models. In *ACL/IJCNLP*, pages 3118–3135, 2021. [2](#)
- [56] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S. Sara Mahdavi, Rapha Gontijo Lopes, Tim Salimans, Jonathan Ho, David J. Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language understanding. In *NeurIPS*, 2022. [2](#)
- [57] Elizabeth Salesky, David Etter, and Matt Post. Robust open-vocabulary translation from visual text representations. In *EMNLP*, pages 7235–7252, 2021. [2](#), [3](#)
- [58] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. LAION-400M: Open dataset of CLIP-filtered 400 million image-text pairs. *CoRR*, abs/2111.02114, 2021. [4](#), [15](#)
- [59] Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words with subword units. In *ACL*, 2016. [2](#)
- [60] Noam Shazeer and Mitchell Stern. Adafactor: Adaptive learning rates with sublinear memory cost. In *ICML*, 2018. [3](#), [14](#)
- [61] Sheng Shen, Liunian Harold Li, Hao Tan, Mohit Bansal, Anna Rohrbach, Kai-Wei Chang, Zhewei Yao, and Kurt Keutzer. How much can CLIP benefit vision-and-language tasks? In *ICLR*, 2022. [2](#), [3](#), [5](#), [6](#), [17](#)
- [62] Baohua Sun, Lin Yang, Catherine Chi, Wenhan Zhang, and Michael Lin. Squared english word: A method of generating glyph to use super characters for sentiment analysis. In *Workshop on Affective Content Analysis*, volume 2328, pages 140–151, 2019. [2](#)
- [63] Zineng Tang, Jaemin Cho, Yixin Nie, and Mohit Bansal. TVLT: Textless vision-language transformer. In *NeurIPS*, 2022. [1](#)
- [64] Ashish Thapliyal, Jordi Pont-Tuset, Xi Chen, and Radu Soricut. Crossmodal-3600: A Massively Multilingual Multimodal Evaluation Dataset. In *EMNLP*, 2022. [4](#), [6](#), [19](#)
- [65] Ruben Villegas, Mohammad Babaeizadeh, Pieter-Jan Kin-dermans, Hernan Moraldo, Han Zhang, Mohammad Taghi Saffar, Santiago Castro, Julius Kunze, and Dumitru Erhan. Phenaki: Variable length video generation from open domain textual description. *CoRR*, abs/2210.02399, 2022. [2](#)
- [66] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In *ICLR*, 2019. [2](#), [3](#), [5](#), [7](#), [18](#)
- [67] Jianfeng Wang, Zhengyuan Yang, Xiaowei Hu, Linjie Li, Kevin Lin, Zhe Gan, Zicheng Liu, Ce Liu, and Lijuan Wang. GIT: A generative image-to-text transformer for vision and language. *CoRR*, abs/2205.14100, 2022. [1](#), [2](#)
- [68] Wenhui Wang, Hangbo Bao, Li Dong, Johan Bjorck, Zhiliang Peng, Qiang Liu, Kriti Aggarwal, Owais Khan Mohammed, Saksham Singhal, Subhojit Som, and Furu Wei. Image as a foreign language: BEiT pretraining for all vision and vision-language tasks. *CoRR*, abs/2208.10442, 2022. [1](#)
- [69] Zirui Wang, Jiahui Yu, Adams Wei Yu, Zihang Dai, Yulia Tsvetkov, and Yuan Cao. SimVLM: Simple visual language model pretraining with weak supervision. In *ICLR*, 2022. [1](#), [8](#)
- [70] Zhuofeng Wu, Sinong Wang, Jiatao Gu, Madian Khabsa, Fei Sun, and Hao Ma. CLEAR: contrastive learning for sentence representation. *CoRR*, abs/2012.15466, 2020. [2](#), [5](#)
- [71] Linting Xue, Aditya Barua, Noah Constant, Rami Al-Rfou, Sharan Narang, Mihir Kale, Adam Roberts, and Colin Raffel. Byt5: Towards a token-free future with pre-trained byte-

- to-byte models. *Trans. Assoc. Comput. Linguistics*, 10:291–306, 2022. 6
- [72] Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. mT5: A massively multilingual pre-trained text-to-text transformer. In *NAACL-HLT*, 2021. 6, 19
- [73] Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *Trans. Assoc. Comput. Linguistics*, 2:67–78, 2014. 4
- [74] Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu. Coca: Contrastive captioners are image-text foundation models. *Transactions on Machine Learning Research*, 2022. 2
- [75] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Guntjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, Ben Hutchinson, Wei Han, Zarana Parekh, Xin Li, Han Zhang, Jason Baldridge, and Yonghui Wu. Scaling autoregressive models for content-rich text-to-image generation. *CoRR*, abs/2206.10789, 2022. 2
- [76] Zhou Yu, Jun Yu, Yuhao Cui, Dacheng Tao, and Qi Tian. Deep modular co-attention networks for visual question answering. In *CVPR*, pages 6281–6290, 2019. 5
- [77] 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 *CVPR*, pages 18102–18112, 2022. 2, 4, 5, 6
- [78] Yukun Zhu, Ryan Kiros, Richard S. Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. In *ICCV*, pages 19–27, 2015. 7, 18



## A. Example input images

Fig. 6 shows two examples of consecutive sentences from the C4 corpus, rendered using our Unifont renderer. The alt-texts for contrastive pretraining are rendered in the same way.

Fig. 7 shows example images from the VQAv2 training set [23] with rendered text in the format we use to adapt CLIPPO (and our baselines) to VQA. The question is rendered with line height of 16px (which is identical to the line height used during pretraining) and the image is resized as to fill the remaining space (with a total image size of  $224 \times 224$  or  $384 \times 384$ ).

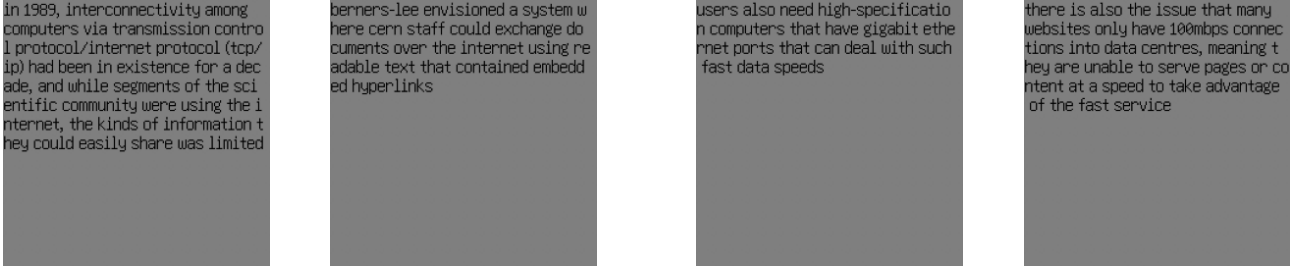


Figure 6. Two examples for rendered consecutive sentences from C4 (image size  $224 \times 224$ ). The rendering is identical for alt-texts.

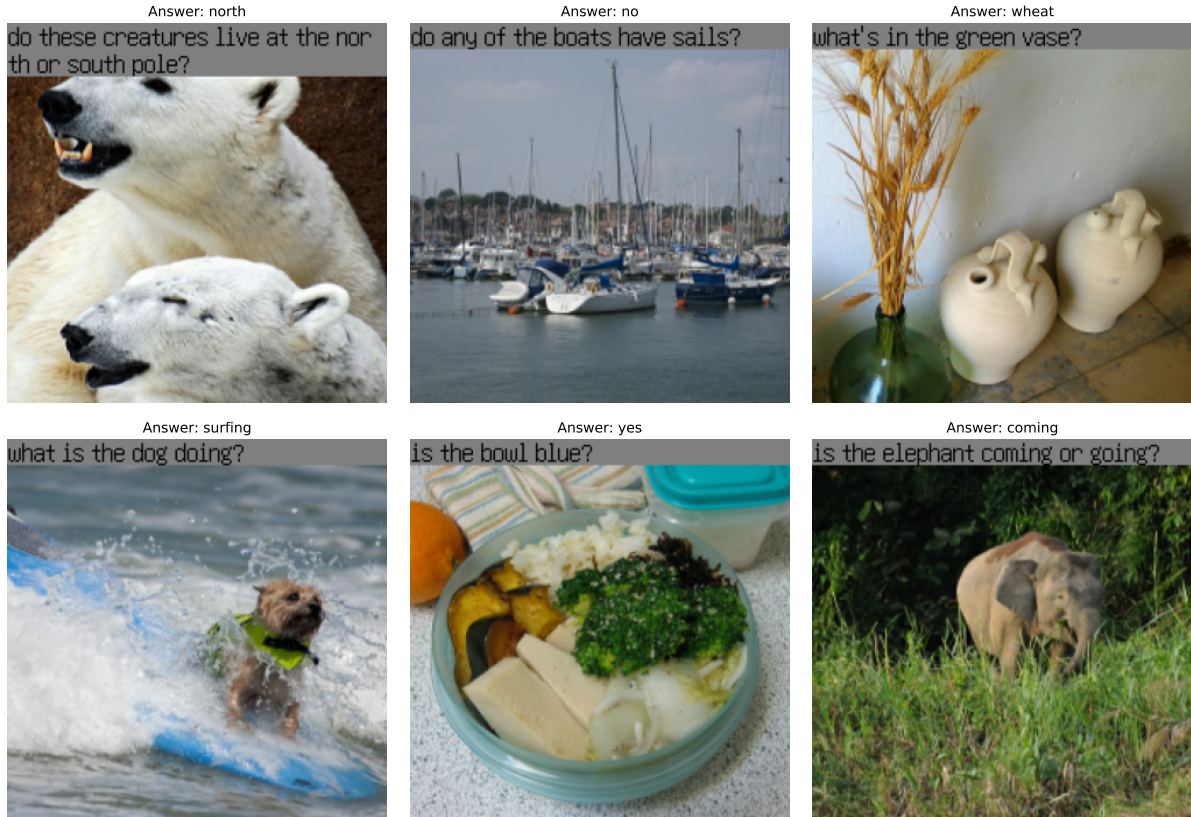


Figure 7. Example training images with rendered questions (black letters on gray background) from the VQAv2 dataset (image size  $224 \times 224$ ). After fine-tuning CLIPPO on VQAv2 it can process images and question jointly in this form. Note that the answers (on white background) are not part of the image.

## B. Training details

We rely on a single training setup for all our baselines and visual text models. This setup was tuned to produce good results for standard image/alt-text contrastive training as in [50] (using exactly the same loss function as [50], following the pseudocode in [50, Fig. 3]) and we found that it readily transfers to 1T-CLIP and CLIPPO (including variants with text/text co-training).

Our default architecture is a ViT-B/16 [15] and we perform a subset of experiments with a ViT-L/16 architecture to study the effect of scale (we equip both models a MAP head [37] to pool embeddings). In all cases, the representation dimension used for the contrastive loss is 768. We set the batch size to 10,240 and train the main models for 250k steps, using a minimum 100k training steps for ablations. For models co-trained with a certain percentage of text/text data, we scale the number of iterations such that the number of image/alt-text pairs seen matches the number of iterations of the corresponding model without text/text data (e.g. when 50% of the data is text/text pairs we increase the number of iterations from 250k to 500k). The contrastive loss is computed across the full batch.

We use the Adafactor optimizer [60] with a learning rate of  $10^{-3}$ , parameters  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$ , and decoupled weight decay with weight  $10^{-4}$ . Gradients are clipped to a norm of 1. We initialize the learned temperature parameter in the contrastive loss with a value of 10. We employ a reciprocal square root schedule with 10k steps linear warmup and 10k steps linear cooldown. This schedule has the advantage that it allows resuming training before cooldown to train a subset of models for more steps (unlike e.g. a cosine schedule which is scaled to a predefined target number of steps). Apart from the learning rate, the training setup is static for all models except for the CLIPPO L/16 models co-trained with 25% and 50% C4 data. To save compute, we do not co-train these with C4 from scratch, but we take the checkpoints pretrained for 150k steps without C4 and continue training these with mixed batches for 350k more steps (i.e. we deviate from the rule described above to adapt the number of training steps with mixed batches).

For all CLIPPO and 1T-CLIP experiments with ViT B/16-scale architecture (i.e. the majority of experiments) we train on 64 Cloud TPUv2 chips. For larger models (CLIP\* B/16 and CLIPPO/1T-CLIP L/16) we use 64 Cloud TPUv3 or Cloud TPUv4 chips to accommodate the increased memory requirements.

### B.1. Fine-tuning details for VQA tasks

Our fine-tuning protocol is inspired by the one described in [16, Sec. 4.1.1]. After replacing the last linear layer of the model with a randomly initialized one with an appropriate number of outputs, we fine-tune for 8,000 steps on a combination of the VQAv2 training set and 90% of the validation set, using the remaining 10% for learning rate selection (recall that we report results on the test-dev set). We rely on SGD with momentum 0.9 and a cosine schedule with 800 linear warmup steps, selecting the learning rate for each model from  $\{0.03, 0.1, 0.2\}$ . The learning rate for the parameters of the freshly initialized head is multiplied by a factor of 10. Gradients are clipped to a norm of 1.

As it is common in the VQA literature to perform evaluation at high resolution, we also evaluate our models on  $384 \times 384$  images (rendering the question at the top of the image following the same strategy as for  $224 \times 224$  images, see Appendix A). To adapt the models to this resolution before fine-tuning, we train a subset of models for 30k iterations at a resolution of 384px, starting from the corresponding 224px checkpoints stored right before cooldown.

## C. Further results and ablations

### C.1. Results on LAION-400M

In Tables 4 and 5 show results on vision and vision-language benchmarks as well as the GLUE benchmark, for the most important CLIPPO and 1T-CLIP models trained on the publicly available LAION-400M dataset [58] (see Appendix C.2 for these results in the context of all other results in the paper). We also show the corresponding models trained on WebLI.

For all the benchmarks/metrics, models trained on LAION-400M exhibit the same ranking as the models trained on WebLI. The ImageNet-1k zero shot and 10-shot results are a few percentage points lower for the models trained on LAION-400M compared to the models trained on WebLI, but the retrieval results on MS-COCO and Flickr30k are consistently a few points better. The GLUE average scores seem largely independent of whether WebLI or LAION-400M is used as a pretraining data set, except for 1T-CLIP, where WebLI-based pretraining leads to a better GLUE score.

	#param.	training dataset	I1k 10s.	I1k 0s.	C I→T	C T→I	F I→T	F T→I
1T-CLIP	118M	WebLI	50.9	60.1	46.2	28.2	76.1	55.2
CLIPPO	93M	WebLI	49.7	58.0	44.9	29.0	73.1	55.4
CLIPPO	93M	WebLI + 25%C4	49.4	55.4	40.2	25.3	69.0	50.5
CLIPPO	93M	WebLI + 50%C4	45.6	51.1	34.3	21.7	61.7	43.2
1T-CLIP	118M	LAION	46.0	54.3	49.0	31.5	77.5	59.7
CLIPPO	93M	LAION	45.3	53.6	46.7	30.3	76.9	58.9
CLIPPO	93M	LAION + 25%C4	44.9	50.6	41.8	27.2	71.1	53.7
CLIPPO	93M	LAION + 50%C4	41.4	46.0	38.2	24.3	66.3	49.0

Table 4. Vision and vision-language cross-modal results obtained when training on LAION-400M [58], along with the corresponding models trained on WebLI. We report ImageNet-1k 10-shot linear transfer validation accuracy (I1k 10s.), ImageNet-1k zero-shot transfer validation accuracy (I1k 0s.), image-to-text and text-to-image retrieval recall@1 on MS-COCO (C I→T and C T→I) and on Flickr30k (F T→I and F I→T). All models have a ViT B/16 architecture (with separate text embedding for 1T-CLIP) and are trained for 100k iterations (with adapted number of steps for models co-trained with C4, see Sec. 4.1).

	training dataset	MNLI-M/MM	QQP	QNLI	SST-2	COLA	STS-B	MRPC	RTE	avg
1T-CLIP text enc.	WebLI	71.6 / 71.6	83.5	80.5	85.0	0.0	74.1	82.8	54.2	67.0
CLIPPO	WebLI	72.2 / 72.2	84.0	81.2	86.7	0.0	81.0	84.0	57.8	68.8
CLIPPO	WebLI + 25%C4	77.0 / 77.0	85.4	82.8	90.9	20.1	83.1	83.6	54.5	72.7
CLIPPO	WebLI + 50%C4	78.8 / 78.8	86.0	84.8	92.0	34.4	83.1	84.2	58.8	75.6
1T-CLIP text enc.	LAION	72.2 / 72.2	84.1	79.8	86.9	0.0	38.0	81.4	54.2	63.3
CLIPPO	LAION	73.2 / 73.2	84.2	80.9	86.5	0.0	75.3	82.2	53.8	67.7
CLIPPO	LAION + 25%C4	77.0 / 77.0	85.5	83.3	91.1	22.0	83.3	84.6	57.0	73.4
CLIPPO	LAION + 50%C4	78.8 / 78.8	86.1	84.3	92.2	38.3	83.7	83.9	55.2	75.7

Table 5. Results for the GLUE benchmark (dev set) when training on LAION-400M [58], along with the corresponding models trained on WebLI. The metric is accuracy except for the performance on QQP and MRPC, which is measured using the  $F_1$  score, CoLA which uses Matthew’s correlation, and STS-B which evaluated based on Spearman’s correlation coefficient. “avg” corresponds to the average across all metrics. All models have a ViT B/16 architecture (with separate text embedding for 1T-CLIP) trained for 100k iterations (with adapted number of steps for models co-trained with C4, see Sec. 4.1).

### C.2. All image, vision-language, and language understanding results

**Image classification and retrieval** Table 6 shows the full set of image classification and image/text retrieval results, including models trained for 100k and 250k steps.

In addition to the results presented in the main paper, we also show results for pretraining with multilingual alt-texts. In this context, CLIP\*, 1T-CLIP, and CLIPPO all obtain a somewhat worse scores on these English-based metrics, but perform much better when evaluated on multilingual image/text retrieval.

We also show results CLIPPO models that were initialized with a ViT trained for image classification. We observe that this improves ImageNet-1k-based classification metrics, but cannot prevent the image and image/text metrics from degrading when co-training with C4 data.

	lan.	#param.	training dataset	steps	I1k 10s.	I1k 0s.	C I→T	C T→I	F I→T	F T→I
CLIP*	EN	203M	WebLI	100k	52.9	62.8	47.2	29.7	76.8	57.2
1T-CLIP	EN	118M	WebLI	100k	50.9	60.1	46.2	28.2	76.1	55.2
CLIPPO	EN	93M	WebLI	100k	49.7	58.0	44.9	29.0	73.1	55.4
CLIPPO untied	EN	186M	WebLI	100k	52.4	61.8	47.2	29.5	76.6	55.0
CLIPPO	EN	93M	WebLI + 25%C4	133k	49.4	55.4	40.2	25.3	69.0	50.5
CLIPPO	EN	93M	WebLI + 50%C4	200k	45.6	51.1	34.3	21.7	61.7	43.2
1T-CLIP L/16	EN	349M	WebLI	100k	58.0	65.6	49.5	31.6	80.2	57.8
CLIPPO L/16	EN	316M	WebLI	100k	56.6	64.9	50.2	33.0	77.0	61.5
CLIP*	ML	203M	WebLI	100k	50.8	59.0	43.6	27.4	71.1	53.2
1T-CLIP	ML	118M	WebLI	100k	49.2	55.2	41.6	25.4	70.9	51.0
CLIPPO	ML	93M	WebLI	100k	47.3	52.0	38.9	24.4	67.7	48.3
CLIPPO JFT init	EN	93M	WebLI	100k	57.1	59.9	43.9	29.2	71.1	55.0
CLIPPO JFT init	EN	93M	WebLI + 25%C4	133k	54.5	56.3	37.0	24.3	64.4	47.3
CLIPPO JFT init	EN	93M	WebLI + 50%C4	200k	50.9	51.8	34.3	22.1	60.5	45.1
1T-CLIP	EN	118M	LAION	100k	46.0	54.3	49.0	31.5	77.5	59.7
CLIPPO	EN	93M	LAION	100k	45.3	53.6	46.7	30.3	76.9	58.9
CLIPPO	EN	93M	LAION + 25%C4	133k	44.9	50.6	41.8	27.2	71.1	53.7
CLIPPO	EN	93M	LAION + 50%C4	200k	41.4	46.0	38.2	24.3	66.3	49.0
CLIP*	EN	203M	WebLI	250k	55.8	65.1	48.5	31.3	79.2	59.4
1T-CLIP	EN	118M	WebLI	250k	53.9	62.3	48.0	30.3	77.5	58.2
CLIPPO	EN	93M	WebLI	250k	53.0	61.4	47.3	30.1	76.4	57.3
CLIPPO	EN	93M	WebLI + 25%C4	333k	52.1	57.4	40.7	26.7	68.9	51.8
CLIPPO	EN	93M	WebLI + 50%C4	500k	48.0	53.1	35.2	23.4	64.8	47.2
1T-CLIP L/16	EN	349M	WebLI	250k	60.8	67.8	50.7	32.5	81.0	61.0
CLIPPO L/16	EN	316M	WebLI	250k	60.3	67.4	50.6	33.4	79.2	62.6
CLIPPO L/16	EN	316M	WebLI + 25%C4	500k	60.5	66.0	44.5	29.8	72.9	57.3
CLIPPO L/16	EN	316M	WebLI + 50%C4	500k	56.8	61.7	39.7	27.3	70.1	54.7
1T-CLIP 384px	EN	118M	WebLI	270k	57.8	66.2	51.5	32.7	81.7	63.0
CLIPPO 384px	EN	93M	WebLI	270k	57.2	64.7	51.0	32.9	79.9	61.9
CLIPPO 384px	EN	93M	WebLI + 25%C4	350k	56.0	61.0	44.3	27.9	73.4	55.0
1T-CLIP L/16 384px	EN	349M	WebLI	270k	64.5	70.9	52.6	34.8	81.6	63.8
CLIPPO L/16 384px	EN	317M	WebLI	270k	63.9	70.5	54.4	35.3	83.6	64.9
CLIPPO L/16 384px	EN	317M	WebLI + 25%C4	520k	64.2	69.0	47.5	31.9	76.2	59.7
CLIP*	ML	203M	WebLI	250k	53.7	62.1	46.9	29.4	76.9	57.8
1T-CLIP	ML	118M	WebLI	250k	52.6	58.4	44.9	27.7	72.2	53.7
CLIPPO	ML	93M	WebLI	250k	51.1	56.1	42.5	26.6	69.9	52.9

Table 6. We report ImageNet-1k 10-shot linear transfer validation accuracy (I1k 10s.), ImageNet-1k zero-shot transfer validation accuracy (I1k 0s.), image-to-text and text-to-image retrieval recall@1 on MS-COCO (C I→T and C T→I) and on Flickr30k (F T→I and F I→T). “CLIPPO untied” is a two tower model where two separate ViT B/16 models (i.e. with separate parameters) are used to encode images and rendered alt-texts. “CLIPPO JFT init” are CLIPPO models that were initialized with the parameters of ViT B/16 from [15] trained on JFT-300M. Models with the suffix “384px” are models trained for 30k iterations at a resolution of 384px, starting from the corresponding 224px checkpoints stored right before cooldown.



**VQA** Table 7 shows results for all our models and baselines on VQAv2 (test-dev set). In addition to what is discussed in the main paper, we observe that co-training with 50% C4 data does not lead to improvements over co-training with 25% C4 data. Further, the gap between 1T-CLIP and CLIPPO becomes narrow as the model size grows. Increasing the resolution from 224px to 384px leads to a substantial improvement across models.

	res.	yes/no	number	other	overall
ViT B/16 JFT	224	71.16	40.71	51.55	58.39
1T-CLIP	224	76.08	42.46	53.1	61.36
CLIP*	224	77.49	44.65	55.47	63.31
CLIPPO 50%C4	224	83.81	45.45	55.62	66.08
CLIPPO	224	83.01	46.36	56.55	66.29
CLIPPO 25%C4	224	84.48	46.18	56.27	66.74
CLIPPO L/16 50%C4	224	84.33	48.2	58.68	68.05
CLIPPO L/16	224	83.74	49.33	58.9	68.05
1T-CLIP L/16	224	84.03	49.41	59.53	68.48
CLIPPO L/16 25%C4	224	84.91	49.26	59.33	68.73
1T-CLIP	384	77.92	45.21	56.45	64.02
CLIPPO	384	84.22	47.94	58.62	67.95
CLIPPO 25%C4	384	86.91	49.34	60.52	70.12
CLIPPO L/16	384	86.26	51.91	61.89	70.79
1T-CLIP L/16	384	86.3	52.01	62.32	71.03
CLIPPO L/16 25%C4	384	86.85	53.57	63.05	71.78
METER CLIP B/32+BERT	224				69.56
ViLT B/32	384				70.33
Pythia CLIP B/16	600				62.72
MCAN CLIP B/32	600				65.40

Table 7. Results on the VQAv2 benchmark (test-dev set). Our 224px and 384px models and baselines are pretrained for 250k and 270k steps (or an appropriately adapted number of steps when co-trained with C4), respectively, and fine-tuned to VQAv2. In addition to CLIPPO and baselines produced in this work, we also compare to Pythia and MCAN models with ViT vision encoders from [61], and with comparably sized METER [16] and ViLT [34] models. “ViT B/16 JFT” is the model trained on JFT-300M from [15].

**Language understanding** Table 8 shows additional results for our models and baselines on the GLUE benchmark. We discuss a number of observations that were not discussed in the main paper.

First, it can be seen that a randomly initialized ViT performs much worse than all the other models, including the vision encoders of the different CLIP\* and 1T-CLIP variants, which all perform similarly, independently on the precise training setup.

We further present results for models that were trained with multilingual image/alt-text pairs (note that GLUE contains only English tasks). When trained for 100k steps, CLIP\*, 1T-CLIP and CLIPPO obtain a lower GLUE score than their counterparts trained on English-only alt-texts. The GLUE scores of these multilingual models improve when training for 250k steps. In particular, CLIPPO almost matches its English-only counterpart, whereas CLIP\* and 1T-CLIP still lag a few points behind their English-only counterparts.

	lan.	training dataset	steps	MNLI-M/MM	QQP	QNLI	SST-2	COLA	STS-B	MRPC	RTE	avg
BERT-Base	EN	Wiki + BC		84.0 / 84.0	87.6	91.0	92.6	60.3	88.8	90.2	69.5	83.1
PIXEL	EN	Wiki + BC		78.1 / 78.1	84.5	87.8	89.6	38.4	81.1	88.2	60.5	76.3
BiLSTM	EN			66.7 / 66.7	82.0	77.0	87.5	17.6	72.0	85.1	58.5	68.1
BiLSTM+Attn+ELMo	EN			72.4 / 72.4	83.6	75.2	91.5	44.1	56.1	82.1	52.7	70.0
ViT from scratch	EN			33.4 / 33.4	51.2	56.4	53.9	0.0	5.1	81.2	52.7	40.8
CLIP* img. enc.	EN	WebLI	100k	65.2 / 65.2	75.7	68.0	77.8	0.0	6.9	81.5	52.3	54.9
CLIP* text enc.	EN	WebLI	100k	70.6 / 70.6	80.6	71.1	85.9	0.0	62.4	82.1	54.9	64.3
1T-CLIP img. enc.	EN	WebLI	100k	64.4 / 64.4	74.2	65.8	74.5	0.0	12.0	81.6	53.8	54.7
1T-CLIP text enc.	EN	WebLI	100k	71.6 / 71.6	83.5	80.5	85.0	0.0	74.1	82.8	54.2	67.0
CLIPPO unt. img. enc.	EN	WebLI	100k	64.8 / 64.8	76.4	67.0	77.1	0.0	7.0	81.4	51.6	54.5
CLIPPO unt. text enc.	EN	WebLI	100k	65.2 / 65.2	83.7	74.8	86.6	3.1	56.1	81.8	54.9	63.5
CLIPPO	EN	WebLI	100k	72.2 / 72.2	84.0	81.2	86.7	0.0	81.0	84.0	57.8	68.8
CLIPPO	EN	WebLI + 25%C4	133k	77.0 / 77.0	85.4	82.8	90.9	20.1	83.1	83.6	54.5	72.7
CLIPPO	EN	WebLI + 50%C4	250k	78.8 / 78.8	86.0	84.8	92.0	34.4	83.1	84.2	58.8	75.6
CLIPPO	EN	C4	100k	79.3 / 79.3	86.4	85.4	93.2	47.7	84.2	83.7	59.6	77.6
CLIPPO	ML	WMT19	100k	72.9 / 72.9	80.8	74.5	88.6	4.0	19.6	81.9	55.6	61.2
CLIPPO	EN	WMT19 BT	100k	70.0 / 70.0	80.5	80.1	84.6	10.8	65.7	81.6	56.0	66.6
1T-CLIP L/16	EN	WebLI	100k	72.8 / 72.8	84.3	81.4	88.5	0.0	79.1	82.3	53.4	68.3
CLIPPO L/16	EN	WebLI	100k	67.4 / 67.4	84.9	76.7	86.5	0.0	81.5	82.9	53.1	66.6
CLIP* img. enc.	ML	WebLI	100k	63.3 / 63.3	73.8	65.9	75.6	0.0	7.0	81.7	54.5	54.0
CLIP* text enc.	ML	WebLI	100k	63.1 / 63.1	79.2	70.6	75.6	4.4	34.8	81.2	49.8	58.0
1T-CLIP img. enc.	ML	WebLI	100k	62.9 / 62.9	73.5	63.8	71.9	0.0	6.5	81.3	53.1	53.0
1T-CLIP text enc.	ML	WebLI	100k	64.9 / 64.9	80.5	74.7	78.6	4.2	66.0	81.5	50.2	62.8
CLIPPO	ML	WebLI	100k	72.0 / 72.0	82.1	80.4	85.0	0.0	16.1	81.6	50.9	60.0
1T-CLIP img. enc.	EN	LAION	100k	66.8 / 66.8	77.9	73.3	78.8	0.0	12.9	81.7	55.2	57.1
1T-CLIP text enc.	EN	LAION	100k	72.2 / 72.2	84.1	79.8	86.9	0.0	38.0	81.4	54.2	63.3
CLIPPO	EN	LAION	100k	73.2 / 73.2	84.2	80.9	86.5	0.0	75.3	82.2	53.8	67.7
CLIPPO	EN	LAION + 25%C4	133k	77.0 / 77.0	85.5	83.3	91.1	22.0	83.3	84.6	57.0	73.4
CLIPPO	EN	LAION + 50%C4	250k	78.8 / 78.8	86.1	84.3	92.2	38.3	83.7	83.9	55.2	75.7
CLIP* img enc.	EN	WebLI	250k	66.4 / 66.4	78.6	69.4	78.6	0.0	5.2	81.2	52.7	55.5
CLIP* text enc.	EN	WebLI	250k	71.8 / 71.8	82.7	73.0	86.2	6.6	65.0	81.4	53.8	65.9
1T-CLIP text enc.	EN	WebLI	250k	72.6 / 72.6	83.8	80.7	84.9	0.0	79.6	83.3	57.0	68.3
CLIPPO	EN	WebLI	250k	73.0 / 73.0	84.3	81.2	86.8	1.8	80.5	84.1	53.4	68.6
CLIPPO	EN	WebLI + 25%C4	333k	77.7 / 77.7	85.3	83.1	90.9	28.2	83.4	84.5	59.2	74.4
CLIPPO	EN	WebLI + 50%C4	500k	79.2 / 79.2	86.4	84.2	92.9	38.9	83.4	84.8	59.9	76.6
CLIPPO	EN	C4	250k	79.9 / 79.9	86.7	85.2	93.3	50.9	84.7	86.3	58.5	78.4
1T-CLIP L/16 text enc.	EN	WebLI	250k	74.3 / 74.3	85.1	81.6	86.6	8.0	82.5	83.1	57.4	70.4
CLIPPO L/16	EN	WebLI	250k	68.4 / 68.4	85.1	77.2	87.6	0.0	81.0	84.3	52.7	67.1
CLIPPO L/16	EN	WebLI + 25%C4	500k	76.6 / 76.6	87.1	79.9	93.2	48.2	84.1	84.6	56.0	76.1
CLIPPO L/16	EN	WebLI + 50%C4	500k	82.3 / 82.3	87.9	86.7	94.2	55.3	85.8	85.9	59.2	80.0
CLIPPO L/16	EN	C4	250k	83.9 / 83.9	87.9	89.1	94.7	62.0	87.1	87.0	62.5	82.0
CLIP* text enc.	ML	WebLI	250k	64.3 / 64.3	80.8	75.7	78.6	11.2	70.7	81.9	49.8	64.2
1T-CLIP text enc.	ML	WebLI	250k	65.8 / 65.8	80.9	75.0	80.7	0.0	71.1	81.9	51.6	63.6
CLIPPO	ML	WebLI	250k	71.1 / 71.1	82.8	79.6	85.2	0.0	78.3	83.1	53.1	67.1

Table 8. Complete results for the GLUE benchmark (dev set). The metric is accuracy except for the performance on QQP and MRPC, which is measured using the  $F_1$  score, CoLA which uses Matthew’s correlation, and STS-B which evaluated based on Spearman’s correlation coefficient. “avg” corresponds to the average across all metrics. The results for BERT-Base and PIXEL are from [54, Table 3], and BiLSTM and BiLSTM+Attn, ELMo from [66, Table 6]. All encoders considered here have a transformer architecture comparable to BERT-Base (up to the text embedding layer), except for CLIPPO L/16 which uses a ViT L/16, and the two BiLSTM model variants. Wiki and BC stand for (English) Wikipedia and Bookcorpus [78] data, respectively. “ViT from scratch” is a randomly initialized, untrained ViT B/16. “CLIPPO unt.” is a two tower model where two separate ViT B/16 models (i.e. with separate parameters) are used to encode images and rendered alt-texts. All models process rendered text except for “CLIP\* text enc.” and “1T-CLIP text enc.” which process tokenized text.

### C.3. The impact of weight sharing

To better understand whether a modality-shared patch embedding or modality-shared heads are degrading the performance of CLIPPO we train different models with separate embeddings and/or heads. The results in Table 9 show that neither of these variants leads to a consistent improvement in image classification or retrieval metrics compared to the default CLIPPO variant where both the embedding and head are shared for image and (rendered) text inputs.

	#param.	shared	separated	I1k 10s.	I1k 0s.	C I→T	C T→I	F I→T	F T→I
CLIP*	203M	-	all	52.9	62.8	47.2	29.7	76.8	57.2
CLIPPO untied	186M	-	all	52.4	61.8	47.2	29.5	76.6	55.0
1T-CLIP	118M	encoder, heads	embeddings	50.9	60.1	46.2	28.2	76.1	55.2
CLIPPO	93M	all	-	49.7	58.0	44.9	29.0	73.1	55.4
CLIPPO	94M	encoder, embeddings	heads	49.2	58.1	45.0	28.7	71.8	56.5
CLIPPO	94M	encoder, heads	embeddings	49.8	58.4	44.5	28.6	73.7	56.4
CLIPPO	94M	encoder	embeddings, heads	48.9	57.6	44.5	26.8	72.9	53.7

Table 9. We report ImageNet-1k 10-shot linear transfer validation accuracy (I1k 10s.), ImageNet-1k zero-shot transfer validation accuracy (I1k 0s.), image-to-text and text-to-image retrieval recall@1 on MS-COCO (C I→T and C T→I) and on Flickr30k (F T→I and F I→T). All models are trained for 100k iterations. “CLIPPO untied” is a two tower model where two separate ViT B/16 models (i.e. with separate parameters) are used to encode images and rendered alt-texts.

### C.4. Multilingual vision-language understanding

**Multilingual image/text retrieval** Fig. 8 shows the per-language retrieval performance on Crossmodal3600 [64] of CLIP\*, 1T-CLIP, and CLIPPO. CLIP\* obtains a slightly better performance than the other two methods which is not surprising given it uses about double the trainable parameters of the other models and separate text and image encoders. CLIPPO matches or outperforms 1T-CLIP on average, despite having fewer trainable parameters. Overall, the performance per-language correlates strongly across all models, with Japanese and Korean showing the biggest differences between CLIPPO and the other models.

**Tokenizers** We use the following open-source tokenizers in our experiments:

- *T5-en* [51]: `gs://t5-data/vocabs/cc_all.32000/sentencepiece.model`
- *T5-all* [51]: `gs://t5-data/vocabs/cc_en.32000/sentencepiece.model`
- *mT5* [72]: `gs://t5-data/vocabs/mc4.250000.100extra/sentencepiece.model`

We take the first 32,000 pieces of the mc4 vocabulary to create a vocabulary of equal size to the others.

**Tokenizer efficiency** Fig. 9 shows the average sequence length on 20,000 samples of different languages from C4. CLIPPO obtains a balanced average performance across the selected languages.

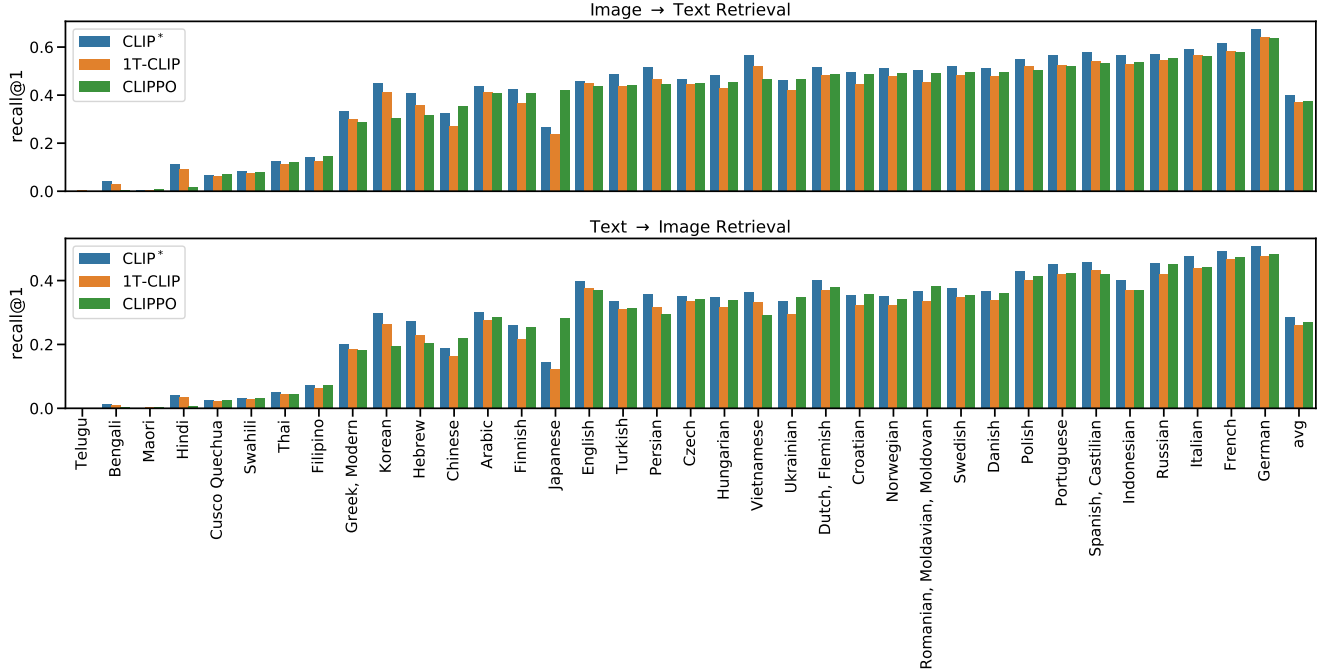


Figure 8. Per-language and average image-to-text and text-to-image recall@1 on the Crossmodal3600 data set. All the models are trained for 250k iterations on WebLI with multilingual alt-texts. CLIP\* and 1T-CLIP use a SentencePiece tokenizer with vocabulary size 32,000 built from 300M randomly sampled WebLI alt-texts, whereas CLIPPO is tokenizer-free by design.

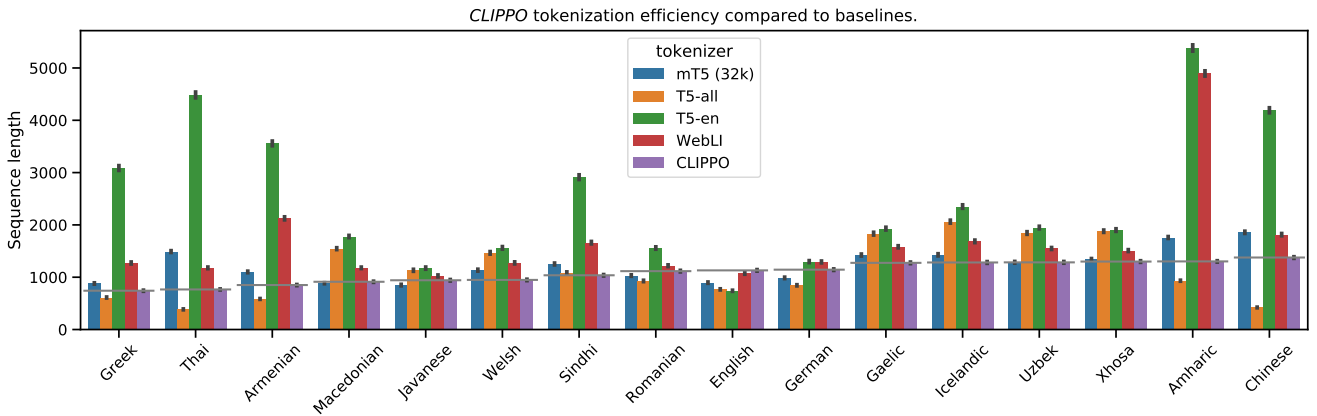


Figure 9. Sequence length of SentencePiece tokenizers derived from different corpora. All non-CLIPPO tokenizers have a vocabulary size of 32,000.

## C.5. Modality gap analysis

Fig. 10 shows additional modality gap visualizations, complementing those in the main paper (Sec. 4.6). In addition to a visualization for the WebLI validation set, we also show results on the MS-COCO validation set. The qualitative and quantitative trend across model variants on MS-COCO is similar to that observed for WebLI, except that the modality gap is somewhat larger for a given model variant (we use the formula from [41, Sec. 4.2] to compute the modality gap). This might be due to the fact that image/caption pairs from MS-COCO have a different distribution than the image/alt-text pairs from WebLI. We further observe that 1T-CLIP and CLIPPO models have a comparable modality gap, and adding more C4 data to the training data mix does not necessarily lead to a reduction in modality gap (going from 25% to 50% C4 data increases the modality gap for MS-COCO).



Since the modality gap measures the Euclidean distance between the image and alt-text mean embeddings it does not fully reflect how the pairwise Euclidean distance between embeddings of corresponding images and alt-texts changes. We plot histograms of the latter in Fig. 11 and observe that the average pairwise distance across models roughly follows the trend of the modality gap. However, the average pairwise distance remains larger than 0.5 even when the modality gap is smaller than 0.1, hence corresponding images and alt-text are not mapped to the same embedding.

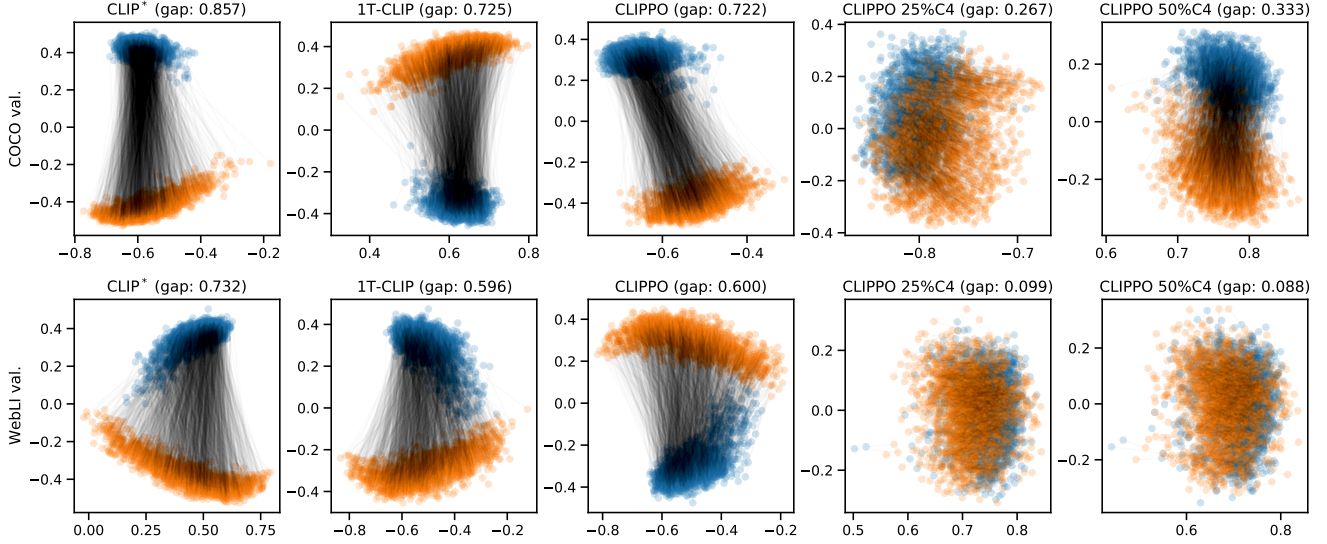


Figure 10. Visualization of the modality gap for examples from the WebLI and MS-COCO validation sets. The visualization follows the analysis from [41] and shows embedded images (blue dots) and corresponding alt-text (orange dots), projected to the first two principal components of the validation data matrix.

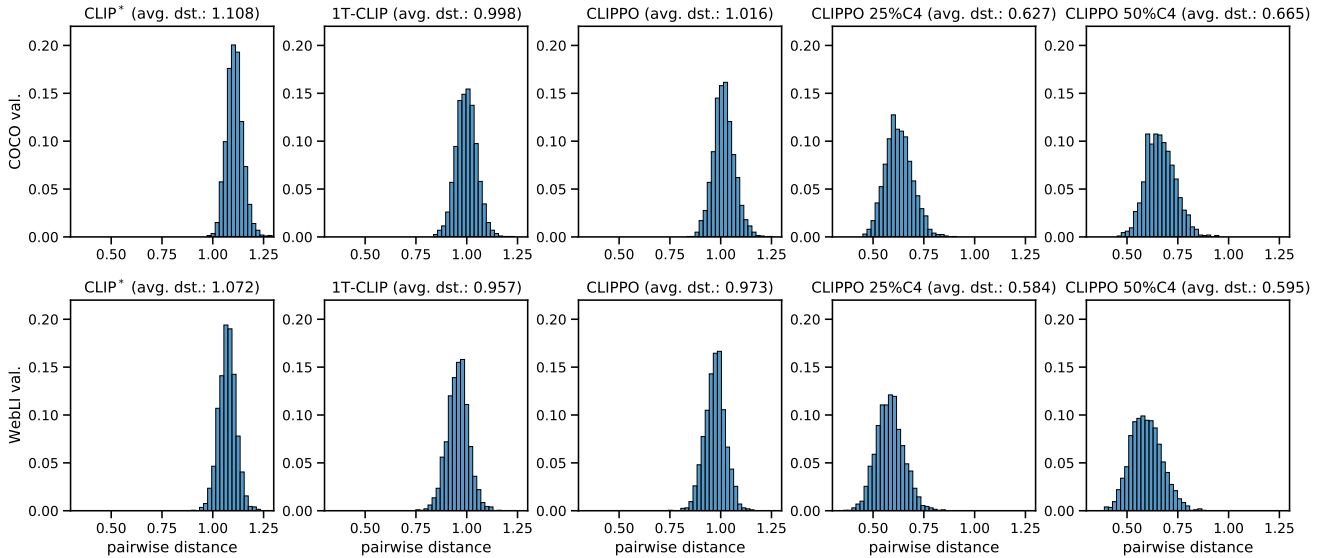


Figure 11. Histograms of the distribution of the Euclidean distance between corresponding image and alt-text embeddings. The average distance across models follows the trend of the modality gap, but the reduction in distance between embeddings when co-training with C4 is not as drastic as for the modality gap.

## C.6. Patch embedding analysis

Following [15], we inspect the patch embedding of different CLIPPO variants and baselines. Concretely, we visualize the top 30 principal components of the patch embedding kernel in Fig. 12. Qualitatively, the top components for CLIP\* and 1T-CLIP are similar to those for supervised ViT training in [15, Sec. 4.5], resembling a plausible basis for image patches. There seems to be no substantial visual difference between the patch embedding structure for English and multilingual variants of CLIP\* and 1T-CLIP. By contrast, the top components for CLIPPO appear to contain more horizontal, high-frequency visual features than the other models, with these features becoming more pronounced as the fraction of C4 data in the training mix increases, or when multilingual alt-text is used. We speculate that this structure might be useful to represent letters and subwords with varying horizontal position as prevalent in the rendered text images fed to CLIPPO.

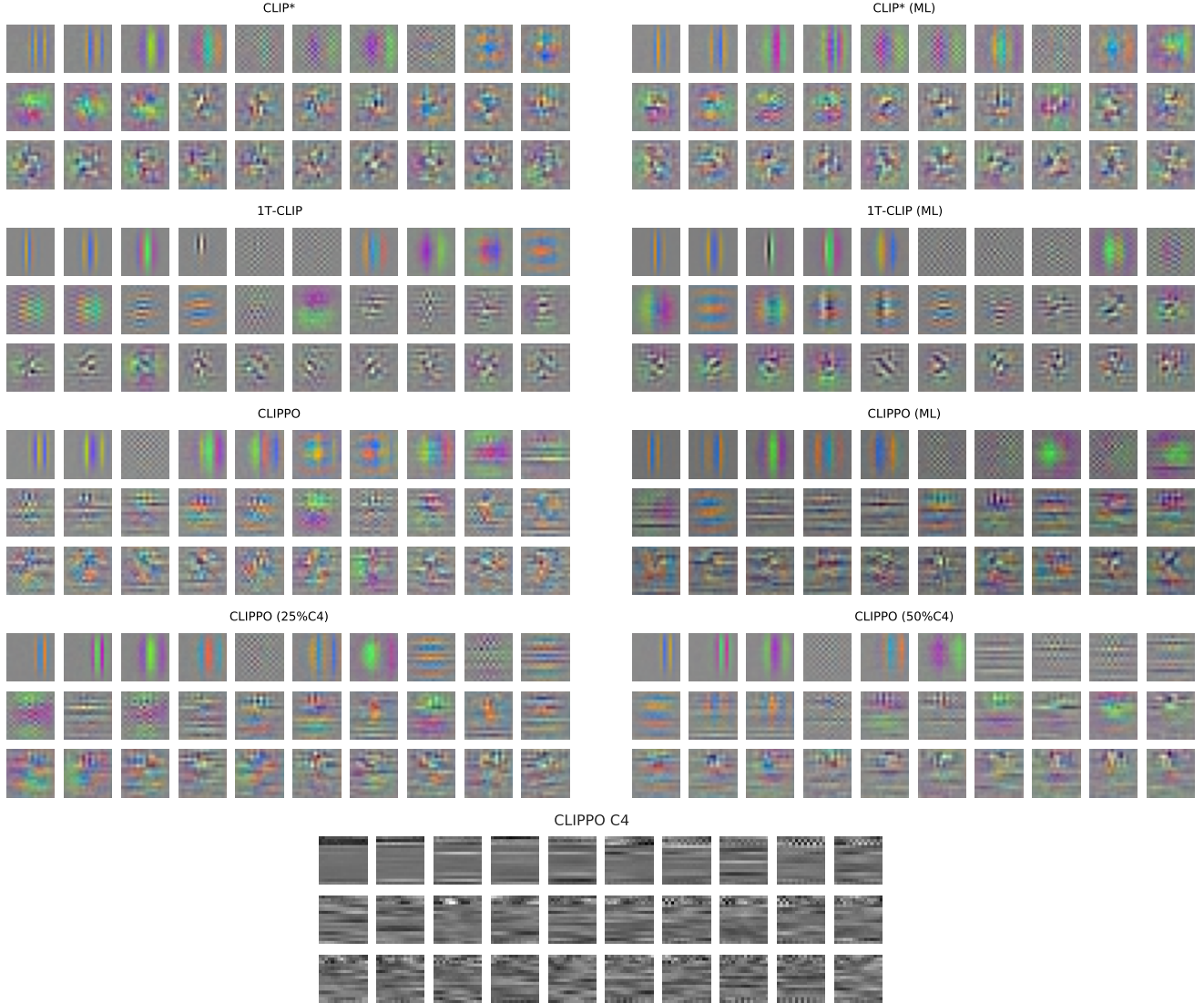


Figure 12. Visualization of the top 30 principal components of the patch embedding kernel for CLIPPO variants and baselines. The top components for CLIPPO appear to contain more horizontal, high-frequency visual features than the other models.