

Summary is provided at the end of the paper. It is recommended to go through it before delving into the details.

# Sigmoid Loss for Language Image Pre-Training

Xiaohua Zhai\* Basil Mustafa Alexander Kolesnikov Lucas Beyer\*  
Google Research, Brain Team, Zürich

{xzhai, basilm, akolesnikov, lbeyer}@google.com

## Abstract

We propose a simple *pairwise sigmoid loss* for image-text pre-training. Unlike standard contrastive learning with softmax normalization, the sigmoid loss operates solely on image-text pairs and does not require a global view of the pairwise similarities for normalization. The sigmoid loss simultaneously allows further scaling up the batch size, while also performing better at smaller batch sizes. With only four TPUv4 chips, we can train a Base CLIP model at 4k batch size and a Large LiT model at 20k batch size, the latter achieves 84.5% ImageNet zero-shot accuracy in two days. This disentanglement of the batch size from the loss further allows us to study the impact of examples vs pairs and negative to positive ratio. Finally, we push the batch size to the extreme, up to one million, and find that the benefits of growing batch size quickly diminish, with a more reasonable batch size of 32k being sufficient. We hope our research motivates further explorations in improving the quality and efficiency of language-image pre-training.

## 1. Introduction

Contrastive pre-training using weak supervision from image-text pairs found on the web is becoming the go-to method for obtaining generic computer vision backbones, slowly replacing pre-training on large labelled multi-class datasets. The high-level idea is to simultaneously learn an aligned representation space for images and texts using paired data. Seminal works CLIP [31] and ALIGN [20] established the viability of this approach at a large scale, and following their success, many large image-text datasets became available privately [47, 11, 19, 39] and publicly [33, 5, 13, 6, 34].

The standard recipe to pre-train such models leverages the image-text contrastive objective. It aligns the image and text embeddings for matching (positive) image-text pairs while making sure that unrelated (negative) image-text pairs

\*equal contribution

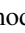




Table 1: **SigLiT and SigLiP results.** Sigmoid loss is memory efficient, allows larger batch sizes (BS) that unlocks language image pre-training with a small number of chips. SigLiT model with a frozen public  L/16 checkpoint [35], trained on the LiT image-text dataset [47] using four TPU-v4 chips for one day, achieves 79.7% 0-shot accuracy on ImageNet. The same setup with a g/14 checkpoint [46] leads to 84.5% accuracy, trained for two days. With a public unlocked  B/16 image checkpoint [35], trained on the WebLI dataset [11], SigLiP achieves 71.0% 0-shot accuracy using 16 TPU-v4 chips for three days. The last two rows show results with randomly initialized models.

	Image	Text	BS	#TPUv4	Days	INet-0
SigLiT	 B/8	L*	32 k	4	1	79.7
SigLiT	 g/14	L	20 k	4	2	84.5
SigLiP	 B/16	B	16 k	16	3	71.0
SigLiP	B/16	B	32 k	32	2	72.1
SigLiP	B/16	B	32 k	32	5	73.4

\* We use a variant of the L model with 12 layers.

Codes  
 

are dissimilar in the embedding space. This is achieved via a batch-level softmax-based contrastive loss, applied twice to normalize the pairwise similarity scores across all images, then all texts. A naive implementation of the softmax is numerically unstable; it is usually stabilized by subtracting the maximum input value before applying the softmax [16], which requires another pass over the full batch.

In this paper, we propose a simpler alternative: the sigmoid loss. It does not require any operation across the full batch and hence greatly simplifies the distributed loss implementation and boosts efficiency. Additionally, it conceptually decouples the batch size from the definition of the task. We compare the proposed sigmoid loss with the standard softmax loss across multiple setups. In particular, we investigate sigmoid-based loss with two prominent approaches for image-text learning: CLIP [31] and LiT [47], which we call sigmoid language image pre-

## 1. Advantages of sigmoid loss

## 2. Batch size scale

training (SigLIP) and sigmoid LiT (SigLiT), respectively. We find that the sigmoid loss performs significantly better than the softmax loss when the batch size is smaller than 16 k. As the train batch size grows, the gap closes. Importantly, the sigmoid loss is symmetric, requires just a single pass, and a typical implementation requires less memory than the softmax loss. This enables successful training of a SigLiT model at a batch size of *one million*. However, we find that the performance saturates with growing batch size, both for softmax and sigmoid. The good news is that a reasonable batch size, i.e. 32 k, is sufficient for image-text pre-training. This conclusion also holds for multilingual SigLIP training on over 100 languages.

In Table 1, we present setups for image-text pre-training that require a moderate amount of TPUv4 chips for training. SigLiT is surprisingly efficient, reaching 79.7% zero-shot accuracy on ImageNet in just a single day on four chips. SigLIP’s more demanding from-scratch training reaches 73.4% zero-shot accuracy in 5 days with 32 TPUv4 chips. This compares favorably to prior works such as FLIP [26] and CLIP [31], which require approximately 5 and 10 days respectively on 256 TPUv3 cores. When fine-tuning a pre-trained vision backbone in SigLIP, denoted as  $\square$  in Table 1, we found that disabling the weight decay on the pre-trained backbone leads to better results (see Figure 4 for details). We hope our work paves the way for making the nascent language-image pre-training field more accessible.

## 2. Related Work

**Contrastive learning with the sigmoid loss.** One prior work proposes a similar sigmoid loss for the task of unsupervised dimensionality reduction [17]; in the scope of contrastive image-text learning, the vast majority of works rely on the softmax-based InfoNCE loss as popularized by [37].

**Contrastive language-image pre-training** has become popular since CLIP [31] and ALIGN [20] applied softmax contrastive learning [48, 37, 9, 21] to large-scale image-text datasets. Both models perform very well on zero-shot transfer tasks, including classification and retrieval. Follow-up works show that contrastively pre-trained models produce good representations for fine-tuning [41, 14], linear regression [20], object detection [27], semantic segmentation [28] and video tasks [45].

**Generative language-image pre-training** Besides softmax contrastive pre-training, various alternatives have been proposed. GIT [39], SimVLM [40], and LEMON [19] successfully pre-train models using a generative text decoder instead, while CoCa [44] adds such a decoder to the discriminative CLIP/ALIGN setup, thus combining the pros and cons of both approaches into a single very capable model. BLIP [25] further proposes CapFit which uses the generative decoder to create better captions and the discrim-

I hope more people will provide explanations like this for the pseudo code in their paper! 🙏

### Algorithm 1 Sigmoid loss pseudo-implementation.

```
1 # img_emb      : image model embedding [n, dim]
2 # txt_emb      : text model embedding [n, dim]
3 # t_prime, b   : learnable temperature and bias
4 # n            : mini-batch size
5
6 t = exp(t_prime)
7 zimg = l2_normalize(img_emb)
8 ztxt = l2_normalize(txt_emb)
9 logits = dot(zimg, ztxt.T) * t + b
10 labels = 2 * eye(n) - ones(n) # -1 with diagonal 1
11 l = -sum(log_sigmoid(labels * logits)) / n
```

inative part of the model to filter pairs. Language-Image pre-training is a very active field and surveys [7] rapidly become outdated.

**Efficient language-image pre-training** On the other hand, few works have tried making language image pre-training more efficient. LiT [47] and FLIP [26] are notable attempts, the former requires a pre-trained and locked backbone, and the latter sacrifices quality by randomly dropping visual tokens. BASIC [30] and LAION [1] look at scaling batch-size but only go up to 16 k and 160 k respectively, by using many hundreds of chips, and for the former also mixing in a large private classification dataset [30, 43]. The recent Lion optimizer [10] claims to be able to reduce the training cost to reach similar quality.

## 3. Method

In this section, we first review the widely-used softmax-based contrastive loss. We then introduce the pairwise sigmoid loss and discuss its efficient implementation.

Given a mini-batch  $\mathcal{B} = \{(I_1, T_1), (I_2, T_2), \dots\}$  of image-text pairs, the contrastive learning objective encourages embeddings of matching pairs  $(I_i, T_i)$  to align with each other, while pushing embeddings of unmatched pairs  $(I_i, T_{j \neq i})$  apart. For practical purposes, it is assumed that for all images  $i$ , the text associated with a different image  $j$  is not related to  $i$ , and vice-versa. This assumption is usually noisy and imperfect.

### 3.1. Softmax loss for language image pre-training

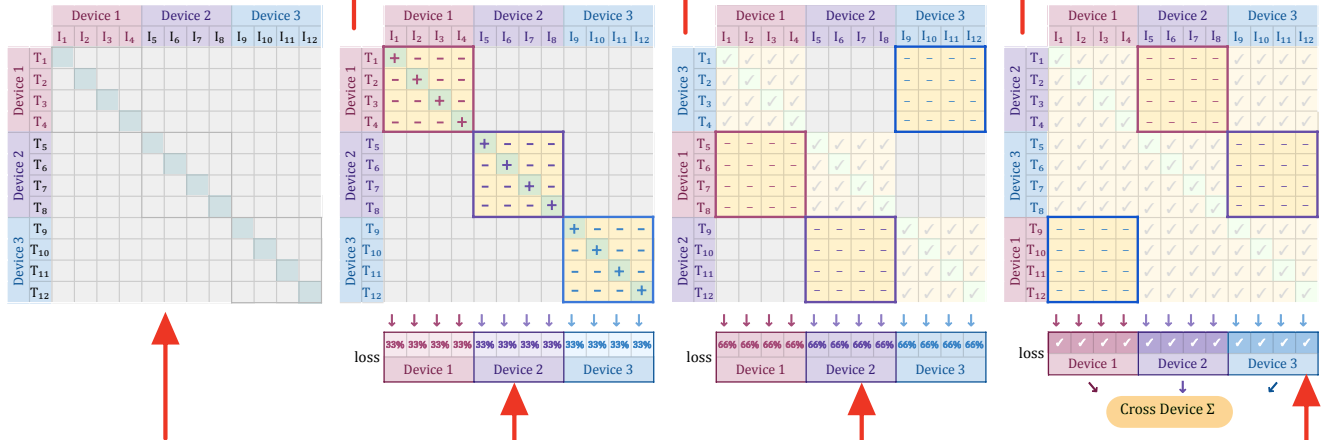
When using the softmax loss to formalize this objective, an image model  $f(\cdot)$  and a text model  $g(\cdot)$  are trained to minimize the following objective:

$$-\frac{1}{2|\mathcal{B}|} \sum_{i=1}^{|\mathcal{B}|} \left( \overbrace{\log \frac{e^{t\mathbf{x}_i \cdot \mathbf{y}_i}}{\sum_{j=1}^{|\mathcal{B}|} e^{t\mathbf{x}_i \cdot \mathbf{y}_j}}}^{\text{image} \rightarrow \text{text softmax}} + \overbrace{\log \frac{e^{t\mathbf{x}_i \cdot \mathbf{y}_i}}{\sum_{j=1}^{|\mathcal{B}|} e^{t\mathbf{x}_j \cdot \mathbf{y}_i}}}^{\text{text} \rightarrow \text{image softmax}} \right)$$

where  $\mathbf{x}_i = \frac{f(I_i)}{\|f(I_i)\|_2}$  and  $\mathbf{y}_i = \frac{g(T_i)}{\|g(T_i)\|_2}$ . In this paper, we adopt the vision transformer architecture [15] for images and the transformer architecture [38] for texts. Note that

Do you know why?

Think why we need to normalize the two parts separately



- (a) Initially each device holds 4 image and 4 text representations. Each device needs to see the representations from other devices to calculate the full loss.
- (b) They each compute the component of the loss (highlighted) for their representations, which includes the positives.
- (c) Texts are swapped across the devices, so device 1 now has  $I_{1:4}$  and  $T_{5:8}$  etc. The new loss is computed and accumulated with the previous.
- (d) This repeats till every image & text pair have interacted, e.g. device 1 has the loss of  $I_{1:4}$  and  $T_{1:12}$ . A final cross-device sum brings everything together.

Figure 1: **Efficient loss implementation** demonstrated via a mock setup with 3 devices and a global batch size of 12. There are no all-gathers, and at any point in time only the bright yellow square (size  $4 \times 4$ ) is materialized in memory.

due to the asymmetry of the softmax loss, the normalization is independently performed two times: across images and across texts [31]. The scalar  $t$  is parametrized as  $\exp(t')$ , where  $t'$  is a global freely learnable parameter.

### 3.2. Sigmoid loss for language image pre-training

Instead of the softmax-based contrastive loss, we propose a simpler alternative that **does not require computing global normalization factors**. The sigmoid-based loss processes every image-text pair independently, **effectively turning the learning problem into the standard binary classification** on the dataset of all pair combinations, with a positive labels for the matching pairs  $(I_i, T_i)$  and negative labels for all other pairs  $(I_i, T_{j \neq i})$ . It is defined as follows:

$$-\frac{1}{|\mathcal{B}|} \sum_{i=1}^{|\mathcal{B}|} \sum_{j=1}^{|\mathcal{B}|} \underbrace{\log \frac{1}{1 + e^{z_{ij}(-t\mathbf{x}_i \cdot \mathbf{y}_j + b)}}}_{\mathcal{L}_{ij}}$$

where  $z_{ij}$  is the label for a given image and text input, which equals 1 if they are paired and  $-1$  otherwise. Note that at initialization, the heavy imbalance coming from the many negatives dominates the loss, leading to large initial optimization steps attempting to correct this bias. To alleviate this, we introduce an additional learnable bias term  $b$  similar to the temperature  $t$ . We initialize  $t'$  and  $b$  to 10 and -10 respectively. This makes sure the training starts roughly close to the prior and does not require massive over-correction. Algorithm 1 presents a pseudocode implementation of the proposed sigmoid loss for language image pre-training.

### 3.3. Efficient “chunked” implementation

Contrastive training typically utilizes data parallelism. Computing the loss when data is split across  $D$  devices necessitates gathering all embeddings [47] with expensive all-gathers and, more importantly, the materialization of a memory-intensive  $|\mathcal{B}| \times |\mathcal{B}|$  matrix of pairwise similarities.

The sigmoid loss, however, is particularly amenable to a memory efficient, fast, and numerically stable implementation that ameliorates both these issues. Denoting the per-device batch size as  $b = \frac{|\mathcal{B}|}{D}$ , the loss is reformulated as:

$$-\frac{1}{|\mathcal{B}|} \underbrace{\sum_{d_i=1}^D}_{\text{A: } \forall \text{ device } d_i} \underbrace{\sum_{d_j=1}^D}_{\text{B: swap negs across devices}} \underbrace{\left( \sum_{i=bd_i}^{b(d_i+1)} \sum_{j=bd_j}^{b(d_j+1)} \mathcal{L}_{ij} \right)}_{\text{C: per device loss}} \underbrace{\substack{\text{all local positives} \\ \text{negs from next device}}}$$

This is particularly simple for the sigmoid loss as each pair is an independent term in the loss. Figure 1 illustrates this method. In words, we first compute the component of the loss corresponding to the positive pairs, and  $b-1$  negative pairs. We then permute representations across devices, so each device takes negatives from its neighbouring device (next iteration of sum B). The loss is then calculated with respect to this chunk (sum C). This is done independently in each device, such that each device computes the loss with respect to its local batch  $b$ . Losses can then simply be summed across all devices (sum A). Individual collective permutes (for sum B) are fast (and indeed  $D$  collective

Another proof that proper initialisation is much better than assuming “the model will learn to correct”

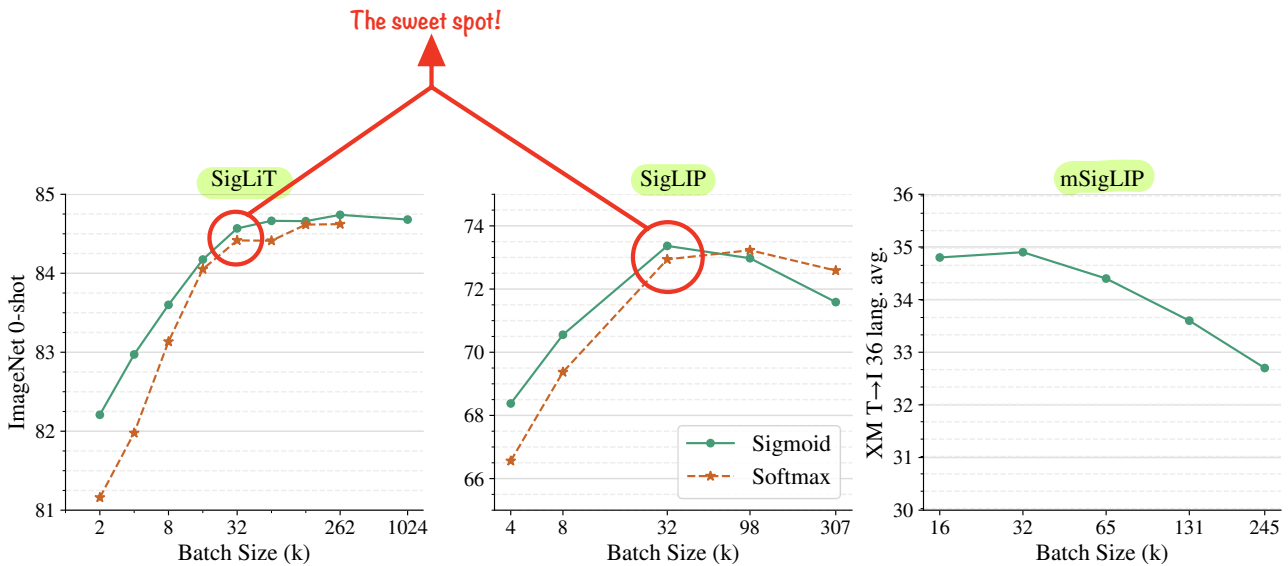


Figure 2: The effect of pre-training batch size. **Left: SigLiT results**, trained for 18B seen examples. Sigmoid loss outperforms the softmax loss significantly with small batch sizes, and performs similarly at larger batch sizes. We successfully trained an SigLiT model with up to *one million* batch size. However, performance for both sigmoid and softmax saturate at around 32 k batch size. **Middle: SigLiP results**, trained for 9B seen examples. Both sigmoid loss and softmax loss saturate at a reasonable batch size, while the peak of the sigmoid loss comes earlier and slightly outperforms the peak of the softmax loss. A very large batch size hurts both losses. **Right: mSigLiP results**, trained for 30B seen examples. With a multilingual setup using over 100 languages, 32 k batch size is surprisingly sufficient and scaling beyond that hurts performance on a 36-language cross-modal retrieval task.

permutes is typically faster than two all-gathers between  $D$  devices), and the memory cost at any given moment is reduced from  $|\mathcal{B}|^2$  to  $b^2$  (for sum C). Usually  $b$  is constant as scaling  $|\mathcal{B}|$  is achieved by increasing the number of accelerators. Due to being quadratic with respect to the batch size, the vanilla loss computation rapidly bottlenecks scaling up. This chunked approach enabled training with batch sizes over 1 million on relatively few devices.

## 4. Results

In this section, we evaluate the proposed SigLiT and SigLiP models across a wide range of batch sizes. We discuss what can be achieved with a small number of accelerator chips, using both SigLiT and SigLiP recipes. We also briefly discuss the impact of batch size on multilingual language image pre-training. We ablate the importance of our large-batch stabilization modification and the introduced learned bias term and present a study on the effect of positive and negative pairs ratio in the sigmoid loss. Lastly, we explore SigLiP’s data noise robustness.

To validate our models, we report zero-shot transfer results on the ImageNet dataset [12] and zero-shot retrieval results across 36 languages on the XM3600 dataset [36]. We use the ScalingViT-Adafactor optimizer [46] by default for all our experiments.

### 4.1. SigLiT: Scaling batch size to the limit

Following [47], we use the same precomputed embeddings for the images using a ViT-g vision model, and train

a base size text tower from scratch with the same hyperparameters using the LiT image-text dataset [47].

We perform a study over a wide range of batch sizes, from 512 to 1 M, demonstrating the impact of batch size for contrastive learning. Results are presented in Figure 2 (left). When the batch size is smaller than 16 k, sigmoid loss outperforms softmax loss by a large margin. With growing batch sizes, we observe that softmax loss quickly catches up and potentially slightly underperforms sigmoid loss with a large enough batch size. Overall, we recommend using the SigLiP recipe for large batch sizes as well, due to the simplicity, compute savings, and straightforward memory efficient implementation.

There is a consensus that contrastive learning benefits from large batch sizes, while most of the existing studies stop at 64 k batch size [47, 30, 9]. We successfully trained an SigLiT model at one million batch size, to explore the limit of contrastive learning. To our surprise, the performance saturates at 32 k batch size, further scaling up the batch size only gives a minor boost, and the model peaks at 256 k batch size. Our best SigLiT with a  $B$ -sized text mode achieves 84.7% zero-shot transfer accuracy on ImageNet, while the original LiT paper reports a slightly better 85.2% score with a 10 times larger  $g$ -sized text model. Figure 3 presents the impact of training duration for different batch sizes. It demonstrates that large, 262 k batch size significantly outperforms smaller 8 k batch size when trained for a sufficiently long time. Note, that for short training durations, large batch size leads to the fewer absolute number of update steps and thus needs more time to ramp up.

Theory only works to a scale

Saying it again, it is always about the “sensible updates”, and not just about the number of updates!



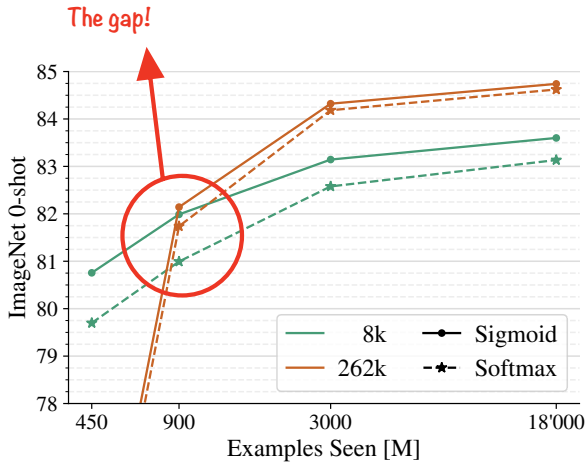


Figure 3: **SigLiT ImageNet 0-shot transfer results with different training durations.** Large batch size results in a big performance boost, but needs a sufficiently long schedule to ramp up, as for short schedules, very large batch size results in a small number of gradient update steps.

#### 4.2. SigLIP: Sigmoid loss is beneficial for language-image pre-training

We pre-train SigLIP models on the WebLI dataset [11], using only English image and text pairs. We use moderately-sized models: B/16 ViT for image embeddings and B-sized transformer for text embeddings. The input images are resized to  $224 \times 224$  resolution. The text is tokenized by a 32k vocabulary sentencepiece tokenizer [24] trained on the English C4 dataset [32], and a maximum of 16 text tokens are kept. Figure 2 middle plot shows SigLIP results. With less than 32k batch size, SigLIP outperforms CLIP baselines with the standard softmax loss. On the other end of the scale, the memory efficiency of the sigmoid loss enabled much larger batch sizes. For example, with four TPU-v4 chips, we could fit a batch size of 4096 with a Base SigLIP but only 2048 with a corresponding CLIP model. The two advantages together demonstrate significant benefits of the sigmoid loss for language image pre-training with fixed resources, which will be discussed in Section 4.5.

As batch size increases, the gap between the sigmoid and the softmax losses diminish. SigLIP performs best at batch size 32k, whereas the softmax loss required 98k for optimal performance and still didn’t outperform the sigmoid based variant. Scaling further, a larger batch size like 307k hurts both losses.

#### 4.3. mSigLIP: Multi-lingual pre-training

We further scale up the training data by keeping all the 100 languages from the WebLI dataset [11]. With multi-lingual data, one usually needs to use a larger international vocabulary. We first verify the impact of two tokenizers: a small multilingual vocabulary with 32k tokens [32], and a

	16 k	32 k	64 k	128 k	240 k
INet-0	71.6	73.2	73.2	73.2	73.1
XM avg	34.8	34.9	34.4	33.6	32.7
XM de	54.7	54.8	55.4	54.3	54.7
XM en	46.5	46.2	46.5	46.6	46.6
XM hi	9.1	8.5	7.9	8.1	7.3
XM ru	50.1	49.9	49.7	48.6	49.3
XM zh	30.7	32.5	32.0	30.6	23.7

Table 2: Multilingual SigLIP results with various batch sizes, pre-trained for 30 billion seen examples. We report zero-shot transfer results on ImageNet (INet-0) and averaged text to image retrieval results across 36 languages on the crossmodal 3600 dataset (XM). The full table on 36 languages can be found in Appendix.

large multilingual vocabulary with 250k tokens [42]. We train B-sized ViT and text models for 900M total examples seen, and observe slightly more than 1% improvement when using a larger vocabulary.

However, the token embeddings become huge for very large vocabulary sizes. Following the standard setup, we would need to store a  $N \times W$  token embedding lookup table to train the multilingual model, where  $N$  is the vocabulary size mentioned above and  $W$  is the embedding dimension of the text model. To save memory, we propose to use a “bottlenecked” token embedding. We use  $N \times K$  embedding matrix and additional  $K \times W$  projection, where the bottleneck  $K$  is much smaller than  $W$ .

In our experiments, we observed that using a large multilingual vocabulary with a bottleneck can be scaled up as efficiently as using a small multilingual vocabulary. Specifically, by enabling the bottleneck of size  $K = 96$  for Base architecture with  $W = 768$ , we only see about a half percent quality drop on ImageNet zero-shot transfer, compared to using the full 250k vocabulary.

With the memory improvements, we train mSigLIP models for various batch sizes, for a total of 30 billion examples seen. Table 2 and Figure 2 (right plot) show the results. A batch size of 32k is sufficient for a multilingual setup as well. On the XM3600 cross-modal retrieval tasks, we found that going beyond 32k batch size leads to worse results on average while on ImageNet zero-shot transfer it stays flat. mSigLIP sets the new state-of-the-art on XM3600 text to image retrieval task, with only a Base size model. Our best result is 34.9%, which is more than 6% higher than the previously reported result 28.5% [11] with a standard LiT model [47] using a much larger four billion ViT-e model.

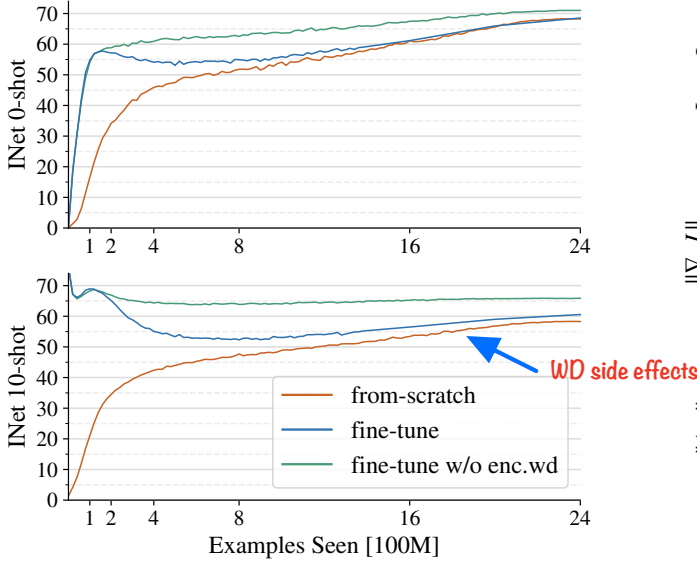


Figure 4: **Top:** SigLIP with pre-trained encoders ramps up quickly. However, only disabling weight decay on the pre-trained encoder weights leads to stable behavior and good ImageNet 0-shot transfer results. **Bottom:** ImageNet 10-shot transfer results, where decaying the pre-trained weights leads to deterioration of the pre-trained model visual representation quality. Disabling weight decay makes the curve flatter.

#### 4.4. SigLiT with four TPU-v4 chips

For many practitioners, the important question usually is “what can be trained with a limited amount of resources?” We explore the usage of SigLiT models in this section with only four TPU-v4 chips, as the memory efficient sigmoid loss is suitable for this application scenario.

We follow the same setup as in section 4.1. We use the publicly available ViT-Augreg-B/8 [35] model as the frozen vision tower, and precompute embeddings to accelerate the training [47]. The text model is a Large Transformer, but with a depth of only 12 layers (instead of 24). It is trained using the LION [10] optimizer with decoupled weight decay  $1 \times 10^{-7}$ , linearly warm-up of learning rate over 6.5k steps up to a peak of  $1 \times 10^{-4}$ , followed by a cosine decay to 0. We train for a total of 65 000 steps with a batch size of 32k – this leads to just under one day of training. Table 1 shows the results when training a model on four chips for one day, achieving 79.7% 0-shot ImageNet classification accuracy; very competitive in this limited resource regime. With a ViT-g/14 [46] model as the vision tower and a Large text tower, we can train at 20k batch size on four chips for 107k steps in under two days. This further pushes the 0-shot ImageNet classification accuracy up to 84.5%.

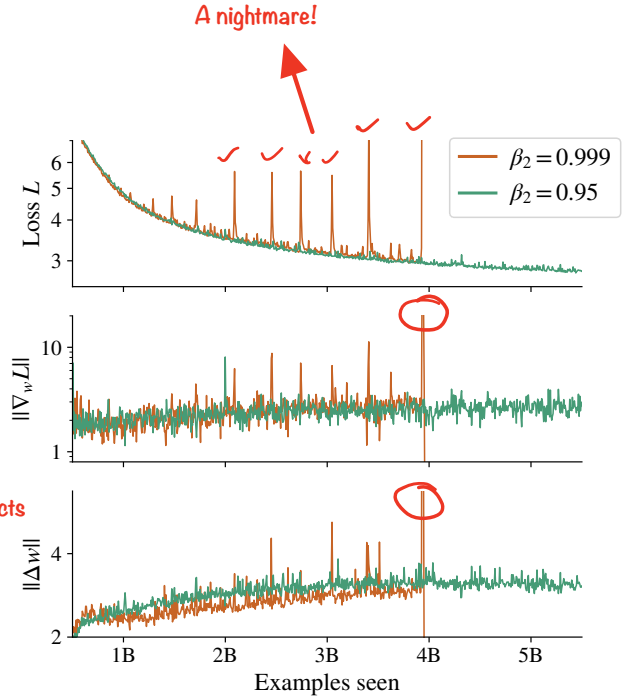


Figure 5: **The effect of Adam and AdaFactor’s  $\beta_2$ .** As we increase batch-size, we observe more frequent training instability. This instability can mainly be seen in the loss curves (top) and is caused by spikes in the gradient norm (middle) which results in large parameter updates (bottom). Decreasing the  $\beta_2$  momentum value stabilizes the training. Even though occasional gradient spikes still happen (see step at 2B), they do not destabilize the training process.

#### 4.5. SigLIP with a small amount of TPU-v4 chips

It’s resource demanding to train a CLIP model from-scratch in general, with SigLIP it’s possible to fit a larger train batch size with fewer amount of chips. In this section, we explore ways to train SigLIP models efficiently with pre-trained weights. We use pre-trained weights to initialize the image model to accelerate the pre-training, which was originally discussed in [47]. We use the public and unlocked ViT-Augreg-B/16 [35] model to initialize our vision tower and fine-tune on the same WebLI English data as used for SigLIP. In all the experiments, we apply a 0.1 learning rate multiplier to the pre-trained image tower to make it suitable for fine-tuning.

Figure 4 presents unlocked ViT fine-tuning results alongside from-scratch randomly initialized baselines. We used 16 TPU-v4 chips and train at 16k batch size for 2.4B examples seen. We found that the fine-tuning setup doesn’t perform well out-of-the-box; this is consistent with prior works [47] where finetuning image models degraded visual representation quality. This is evidenced by ImageNet 10-shot linear classification, where in Figure 4 the fine-tuned setup is barely better than the from-scratch baseline.

We hypothesize that the default weight decay applied to the pre-trained weights reduces their effectiveness. Moti-

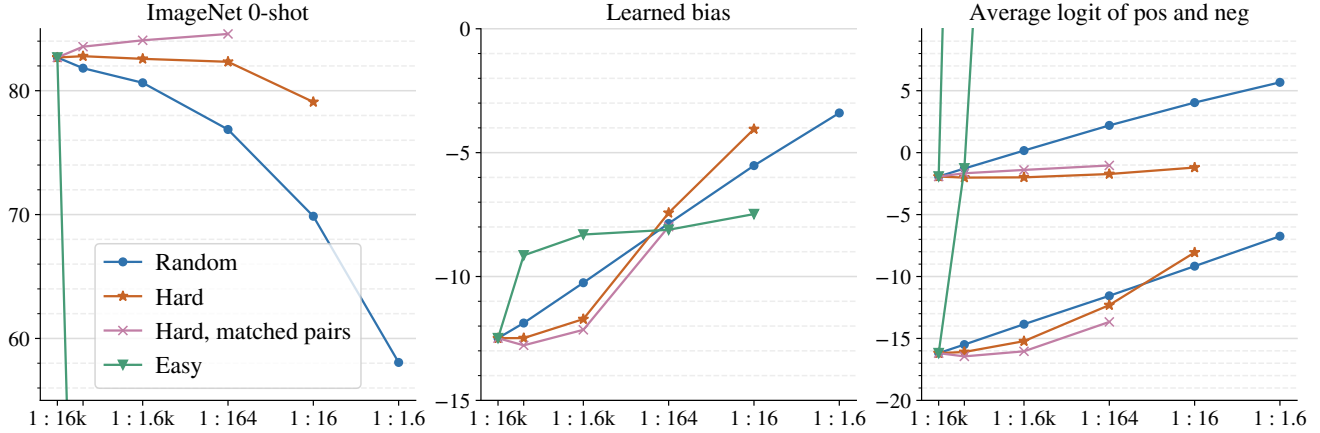


Figure 6: **The effect of batch composition.** We simulate various batch compositions by masking out negatives, either randomly, keeping only the hardest, or the easiest. With no masking, we have 16k negatives for each positive in the batch (1:16k) and the strongest masking we apply (1:1.6) results in almost balanced minibatches. In one setting we *match total pairs* seen by training for significantly longer. We observe ImageNet 0-shot score, the final value of the learned bias, and the average logits of positive and negative pairs. Overall, the imbalance does not seem to be detrimental, but finding an *efficient* way of mining negatives might be beneficial.

vated by the fine-tuning recipe from [15, 46, 22], that uses no weight decay, we also propose disabling weight decay on the pre-trained weights for SigLIP training. Weight decay is therefore only applied to the randomly initialized weights in the text model. This simple modification significantly improved SigLIP results. Figure 4 shows that with our improved recipe, SigLIP reaches 71% 0-shot accuracy on ImageNet, using 16k batch size, trained on 16 chips for three days. We also present from-scratch results in the bottom rows of Table 1: with 32 TPUv4 chips for only two days, SigLIP achieves 72.1% 0-shot accuracy. This presents a significant training cost reduction e.g. compared to CLIP (approx. 2500 TPUv3-days for 72.6%) reported in [26].

#### 4.6. Stabilizing large-batch training

As we move to large batch sizes, the language image pre-training using transformers becomes increasingly more unstable, even when using a modestly-sized model (e.g. Base size). The reason for these instabilities is large spikes in the gradient norms, which translate to large-magnitude changes in the weights that may destabilize the training process, see Figure 5. We observe that reducing  $\beta_2$  in Adam and AdaFactor from its default 0.999 to 0.95 (which was suggested in [18, 8]) is enough to stabilize the training.<sup>2</sup> Intuitively, this allows recovering from gradient spikes quicker. We opt for setting  $\beta_2 = 0.95$  by default for all our experiments.

<sup>2</sup>Lucas thanks Kaiming He and Xinlei Chen for discussion of  $\beta_2$ .

#### 4.7. Negative ratio in sigmoid loss

One question which arises when shifting the perspective from the softmax’s “pick the right class” view to the sigmoid’s “rate this pair” view, is the imbalance in positive versus negative pairs. For a batch size  $|\mathcal{B}|$ , the batch contains  $|\mathcal{B}|$  positive pairs, but  $|\mathcal{B}|^2 - |\mathcal{B}|$  negative examples. In the modest batch-size of 16k, there are actually 268M negative examples for only 16k positive ones. At the same time, because the sigmoid loss decomposes into a sum of per-example losses, we can perform controlled experiments to study the effect of the mini-batch composition and distribution of examples visited. We run experiments in the SigLiT setup at batch-size 16k for 900M steps and vary the composition of the batch by masking out (*i.e.* ignoring) enough negative examples to reach a target “positive : negative” ratio, masking in the following ways:

- **Random:** Randomly choose negative pairs to mask.
- **Hard:** Keep hardest negative pairs (highest loss).
- **Easy:** Keep easiest negatives pairs (lowest loss).
- **Hard + matching total pairs seen:** Masking examples while training for a fixed number of steps does decrease the total number of *pairs* seen during training. Hence in the *matched pairs* setting, we increase the number of training steps by the masking ratio in order to keep the number of pairs seen constant.

Figure 6 shows the effect of the various masking strategies. Randomly removing negatives to rebalance does deteriorate performance. Keeping the easiest examples does not work at all, while keeping the hardest negatives does almost

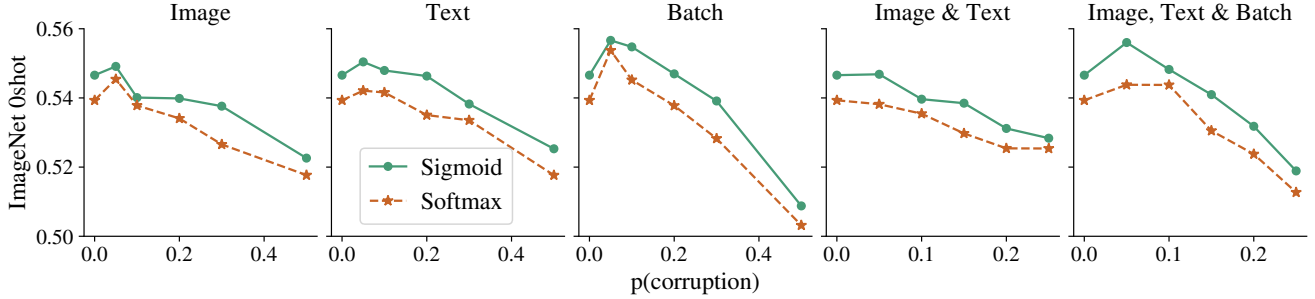


Figure 7: **Sigmoid-training increases robustness** to data noise. Titles show the type of corruption applied, and x-axes show the probability with which they are applied. With increasing corruption severity, M-scale models trained with sigmoid loss for 3.6 billion examples retain superiority over corresponding softmax baseline.

maintain the quality, indicating that, as could be expected, a lot of the learning on the negative side comes from the harder examples. This is further confirmed by the slightly increased performance of training longer on the hardest examples in order to match the total pairs seen.

We also look at the value of the learned bias at the end of training as well as the average logit value for positive and negative examples across these settings, and find the result mostly follows what one would expect: as fewer negatives are present, the bias and logits become more positive overall. Interestingly, when training with more hard negative pairs, the average logits of positive pairs stays mostly flat.

This study confirms that (1) the imbalance does not seem to be a major reason for concern, while at the same time (2) coming up with an *efficient* way of including more negative examples can be promising but is not trivial.

#### 4.8. Bias term in sigmoid loss

We ablate the bias term in the loss function, using the Base architecture with an 8k batch size, trained for 900M examples with the SigLIP setup. Zero-shot transfer results are reported on ImageNet [12], Oxford-iiit pet [29] and Cifar100 [23]. Table 3 presents results with and without a bias term in the sigmoid loss.

Table 3: **Bias (b) and temperature (t') initialization.** Results are reported using Base architecture, 8k batch size, trained for 900M examples. Enabling the bias term b with  $-10$  initialization improves results consistently.

b	t'	INet-0	Pet-0	C100-0
n/a	log 10	62.0	81.8	59.9
<b>-10</b>	<b>log 10</b>	<b>63.0</b>	<b>82.4</b>	<b>61.0</b>
-10	log 1	61.0	80.0	60.4
0	log 10	61.7	79.9	59.0
<b>0</b>	<b>log 1</b>	<b>53.7</b>	<b>73.2</b>	<b>53.8</b>

Enabling the bias term with a  $-10$  initialization consistently improves performance across all tasks. This is because the bias term ensures that the training starts close to the prior, preventing dramatic over-correction in early optimization. In contrast, a randomly chosen bias term initialization, such as the 0 initialization in Table 3, fails to address the over-correction issue, leading to significantly worse results. This effect is particularly noticeable when using a small temperature  $t'$  initialization. We set the bias and temperature initialization to  $b = -10$  and  $t' = \log 10$  (hence  $t = 10$ ) as the default for all experiments.

#### 4.9. Label noise robustness

Prior works demonstrated improved robustness against label noise when using the sigmoid loss for classification models [2]. This property would be particularly useful here in the face of the famously noisy nature of popular large-scale image-text datasets. In order to study this for SigLIP, we train M/16 image models alongside an M text model at batch size 16384 for 3.6 billion seen examples. We corrupt the training data using one of the following methods:

- **Image:** With probability  $p$ , replace the image with uniform random noise.
- **Text:** With probability  $p$ , replace tokenized text with a new sequence of randomly sampled tokens, up to some (sampled) sequence length.
- **Batch alignment:** Randomly shuffle the ordering of  $p\%$  of the batch.
- **Image & text:** Independently apply (1) and (2), each with probability  $p$ .
- **Image, text & batch:** Alongside (4), also shuffle fraction  $p$  of alignments.

Results from varying the likelihood of the corruption are shown in Figure 7. Models trained with sigmoid loss are increasingly robust to all kinds of added noise.



## 5. Conclusion

We conducted a study on two language-image pre-training instances that used the sigmoid loss: SigLiT and SigLiP. Our results demonstrate that the sigmoid loss performs better than the softmax baseline, particularly for small train batch sizes. This loss function is also more memory efficient, which allows larger train batch sizes without requiring additional resources. We performed a thorough investigation of the batch size in contrastive learning. Surprisingly, we found that a relatively modest batch size of 32k yielded nearly optimal performance. Further studies have been performed to understand better the introduced bias term in the sigmoid loss, robustness to data noises and the impact of positive and negative pairs ratio in the sigmoid loss. We hope this work will facilitate language-image pre-training research with limited resources.

**Acknowledgements.** We thank Daniel Keysers, Ilya Tolstikhin, Olivier Bousquet and Michael Tschannen for their valuable feedback and discussions on this paper. We thank Joan Puigcerver, Josip Djolonga and Black Hechtman for discussions on efficient implementations of the chunked contrastive loss. We also thank Ross Wightman for spotting a mistake in the pseudocode in the first version of this paper. Similarly, we thank Boris Dayma for spotting a typo making  $t$  vs  $t'$  confusing, which we fixed in the third version of this paper. As always, we thank the Google Brain team at large for providing a supportive research environment. We use the `big_vision` codebase [4, 3] for all experiments in this project.

## References

- [1] Reaching 80% zero-shot accuracy with OpenCLIP: ViT-G/14 trained on LAION-2B. <https://web.archive.org/web/20230127012732/https://laion.ai/blog/giant-openclip/>. 2
- [2] Lucas Beyer, Olivier J. Hénaff, Alexander Kolesnikov, Xiaohua Zhai, and Aäron van den Oord. Are we done with imagenet? *CoRR*, abs/2006.07159, 2020. 8
- [3] Lucas Beyer, Xiaohua Zhai, and Alexander Kolesnikov. Better plain vit baselines for imagenet-1k, 2022. 9
- [4] Lucas Beyer, Xiaohua Zhai, and Alexander Kolesnikov. Big vision. [https://github.com/google-research/big\\_vision](https://github.com/google-research/big_vision), 2022. 9
- [5] Minwoo Byeon, Beomhee Park, Haecheon Kim, Sungjun Lee, Woonhyuk Baek, and Saehoon Kim. Coyo-700m: Image-text pair dataset. <https://github.com/kakaobrain/coyo-dataset>, 2022. 1
- [6] Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12M: Pushing web-scale image-text pre-training to recognize long-tail visual concepts. In *CVPR*, 2021. 1
- [7] Feilong Chen, Duzhen Zhang, Minglun Han, Xiu-Yi Chen, Jing Shi, Shuang Xu, and Bo Xu. VLP: A survey on vision-language pre-training. *Int. J. Autom. Comput.*, 20(1):38–56, 2023. 2
- [8] Mark Chen, Alec Radford, Rewon Child, Jeffrey Wu, Heewoo Jun, David Luan, and Ilya Sutskever. Generative pre-training from pixels. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages 1691–1703. PMLR, 2020. 7
- [9] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey E. Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, 2020. 2, 4
- [10] Xiangning Chen, Chen Liang, Da Huang, Esteban Real, Kaiyuan Wang, Yao Liu, Hieu Pham, Xuanyi Dong, Thang Luong, Cho-Jui Hsieh, Yifeng Lu, and Quoc V. Le. Symbolic discovery of optimization algorithms, 2023. 2, 6
- [11] 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, 5
- [12] 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, 8
- [13] Karan Desai, Gaurav Kaul, Zubin Aysola, and Justin Johnson. Redcaps: Web-curated image-text data created by the people, for the people. In Joaquin Vanschoren and Sai-Kit Yeung, editors, *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual*, 2021. 1
- [14] Xiaoyi Dong, Jianmin Bao, Ting Zhang, Dongdong Chen, Shuyang Gu, Weiming Zhang, Lu Yuan, Dong Chen, Fang Wen, and Nenghai Yu. Clip itself is a strong fine-tuner: Achieving 85.7% and 88.0% top-1 accuracy with vit-b and vit-l on imagenet. *CoRR*, abs/2212.06138, 2022. 2
- [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. 2, 7
- [16] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>. 1
- [17] Raia Hadsell, Sumit Chopra, and Yann LeCun. Dimensionality reduction by learning an invariant mapping. In *CVPR*, volume 2, 2006. 2
- [18] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross B. Girshick. Masked autoencoders

- are scalable vision learners. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pages 15979–15988. IEEE, 2022. 7
- [19] Xiaowei Hu, Zhe Gan, Jianfeng Wang, Zhengyuan Yang, Zicheng Liu, Yumao Lu, and Lijuan Wang. Scaling up vision-language pre-training for image captioning. *CoRR*, abs/2111.12233, 2021. 1, 2
- [20] 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. 1, 2
- [21] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. 2
- [22] Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Joan Puigcerver, Jessica Yung, Sylvain Gelly, and Neil Houlsby. Big transfer (BiT): General visual representation learning. In *ECCV*, 2020. 7
- [23] Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, Univ. of Toronto, 2009. 8
- [24] Taku Kudo and John Richardson. SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In *EMNLP*, 2018. 5, 12
- [25] Junnan Li, Dongxu Li, Caiming Xiong, and Steven C. H. Hoi. BLIP: bootstrapping language-image pre-training for unified vision-language understanding and generation. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato, editors, *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pages 12888–12900. PMLR, 2022. 2
- [26] Yanghao Li, Haoqi Fan, Ronghang Hu, Christoph Feichtenhofer, and Kaiming He. Scaling language-image pre-training via masking. *CoRR*, abs/2212.00794, 2022. 2, 7
- [27] Matthias Minderer, Alexey A. Gritsenko, Austin Stone, Maxim Neumann, Dirk Weissenborn, Alexey Dosovitskiy, Aravindh Mahendran, Anurag Arnab, Mostafa Dehghani, Zhuoran Shen, Xiao Wang, Xiaohua Zhai, Thomas Kipf, and Neil Houlsby. Simple open-vocabulary object detection. In Shai Avidan, Gabriel J. Brostow, Moustapha Cissé, Giovanni Maria Farinella, and Tal Hassner, editors, *Computer Vision - ECCV 2022 - 17th European Conference, Tel Aviv, Israel, October 23-27, 2022, Proceedings, Part X*, volume 13670 of *Lecture Notes in Computer Science*, pages 728–755. Springer, 2022. 2
- [28] Jishnu Mukhoti, Tsung-Yu Lin, Omid Poursaeed, Rui Wang, Ashish Shah, Philip H. S. Torr, and Ser-Nam Lim. Open vocabulary semantic segmentation with patch aligned contrastive learning, 2022. 2
- [29] Omkar M. Parkhi, Andrea Vedaldi, Andrew Zisserman, and C. V. Jawahar. Cats and dogs. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2012. 8
- [30] Hieu Pham, Zihang Dai, Golnaz Ghiasi, Hanxiao Liu, Adams Wei Yu, Minh-Thang Luong, Mingxing Tan, and Quoc V. Le. Combined scaling for zero-shot transfer learning. *CoRR*, abs/2111.10050, 2021. 2, 4
- [31] 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
- [32] 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. *arXiv e-prints*, 2019. 5, 12
- [33] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. LAION-5B: an open large-scale dataset for training next generation image-text models. *CoRR*, abs/2210.08402, 2022. 1
- [34] Krishna Srinivasan, Karthik Raman, Jiecao Chen, Michael Bendersky, and Marc Najork. WIT: wikipedia-based image text dataset for multimodal multilingual machine learning. *CoRR*, abs/2103.01913, 2021. 1
- [35] Andreas Steiner, Alexander Kolesnikov, Xiaohua Zhai, Ross Wightman, Jakob Uszkoreit, and Lucas Beyer. How to train your ViT? Data, augmentation, and regularization in vision transformers. *CoRR*, abs/2106.10270, 2021. 1, 6
- [36] Ashish V. Thapliyal, Jordi Pont-Tuset, Xi Chen, and Radu Soricut. Crossmodal-3600: A massively multilingual multimodal evaluation dataset. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang, editors, *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022*, pages 715–729. Association for Computational Linguistics, 2022. 4
- [37] Aäron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *CoRR*, abs/1807.03748, 2018. 2
- [38] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017. 2
- [39] 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
- [40] Zirui Wang, Jiahui Yu, Adams Wei Yu, Zihang Dai, Yulia Tsvetkov, and Yuan Cao. Simvln: Simple visual language model pretraining with weak supervision. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net, 2022. 2

- [41] Mitchell Wortsman, Gabriel Ilharco, Jong Wook Kim, Mike Li, Simon Kornblith, Rebecca Roelofs, Raphael Gontijo Lopes, Hannaneh Hajishirzi, Ali Farhadi, Hongseok Namkoong, and Ludwig Schmidt. Robust fine-tuning of zero-shot models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pages 7949–7961. IEEE, 2022. 2
- [42] 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. 5
- [43] Jianwei Yang, Chunyuan Li, Pengchuan Zhang, Bin Xiao, Ce Liu, Lu Yuan, and Jianfeng Gao. Unified contrastive learning in image-text-label space. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pages 19141–19151. IEEE, 2022. 2
- [44] Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu. Coca: Contrastive captioners are image-text foundation models. *CoRR*, abs/2205.01917, 2022. 2
- [45] Lu Yuan, Dongdong Chen, Yi-Ling Chen, Noel Codella, Xiyang Dai, Jianfeng Gao, Houdong Hu, Xuedong Huang, Boxin Li, Chunyuan Li, Ce Liu, Mengchen Liu, Zicheng Liu, Yumao Lu, Yu Shi, Lijuan Wang, Jianfeng Wang, Bin Xiao, Zhen Xiao, Jianwei Yang, Michael Zeng, Luowei Zhou, and Pengchuan Zhang. Florence: A new foundation model for computer vision. *CoRR*, abs/2111.11432, 2021. 2
- [46] Xiaohua Zhai, Alexander Kolesnikov, Neil Houlsby, and Lucas Beyer. Scaling vision transformers. *CVPR*, 2022. 1, 4, 6, 7, 12
- [47] 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 *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pages 18102–18112. IEEE, 2022. 1, 2, 3, 4, 5, 6, 12
- [48] Yuhao Zhang, Hang Jiang, Yasuhide Miura, Christopher D. Manning, and Curtis P. Langlotz. Contrastive learning of medical visual representations from paired images and text. In Zachary C. Lipton, Rajesh Ranganath, Mark P. Sendak, Michael W. Sjoding, and Serena Yeung, editors, *Proceedings of the Machine Learning for Healthcare Conference, MLHC 2022, 5-6 August 2022, Durham, NC, USA*, volume 182 of *Proceedings of Machine Learning Research*, pages 2–25. PMLR, 2022. 2

## A. More results for SigLiT

In section 4.1, we use the same precomputed embeddings for the images using a ViT-g vision model from [47]. Only resize augmentation is applied, to a fixed  $288 \times 288$  resolution. We train a standard base size text tower, using the ScalingViT-Adafactor optimizer [46] with  $\beta_1 = 0.9$  and  $\beta_2 = 0.95$ . We use 0.001 learning rate with a linear warmup schedule for the first 200 M examples seen, and then the learning rate is decayed to zero with a cosine learning rate schedule. Weight decay is set to 0.0001 for all the experiments. The text is tokenized by a 32 k vocabulary sentence-piece tokenizer [24] trained on the English C4 dataset [32], and a maximum of 16 text tokens are kept. Table 5 shows results with multiple train examples seen and batch sizes, for both the sigmoid loss and the softmax loss baseline.

For training SigLiT in under a day with 4 chips (Section 4.4), we used the LION optimizer with peak learning rate  $1 \times 10^{-4}$  and weight decay  $1 \times 10^{-7}$ . The learning rate was warmed linearly to the peak in 6.5 k steps, then cosine decayed to zero for the remaining 58.5 k steps.

## B. More results for SigLIP

In Table 4, we present more results for SigLIP with multiple train examples seen: 3 billion examples and 9 billion examples respectively.

Batch Size	3 B		9 B	
	sigmoid	softmax	sigmoid	softmax
512	<b>51.5</b>	47.7	-	-
1 k	<b>57.3</b>	53.2	-	-
2 k	<b>62.1</b>	59.3	-	-
4 k	<b>65.3</b>	63.8	<b>68.4</b>	66.6
8 k	<b>68.6</b>	66.6	<b>70.6</b>	69.4
32 k	<b>69.9</b>	<b>69.9</b>	<b>73.4</b>	72.9
98 k	69.5	<b>69.7</b>	73.0	<b>73.2</b>
307 k	-	-	71.6	<b>72.6</b>

Table 4: **SigLIP zeor-shot accuracy (%) on the ImageNet benchmark.** Both the sigmoid loss and the softmax loss baseline are presented. Experiments are performed on multiple train examples seen (3 B, 9 B) and train batch sizes (from 512 to 307 k). When trained for 9 B examples, the peak of the sigmoid loss comes earlier at 32 k than the peak of the softmax loss at 98 k. Together with the memory efficient advantage for the sigmoid loss, it allows one to train the best language-image model with much fewer amount of accelerators.

## C. More results for mSigLIP

We present the crossmodal retrieval results on the Crossmodal-3600 dataset, across all the 36 languages in Figure 8, Figure 9 and Table 6.

## D. Label noise experiments

All models had an M/16 image tower and a M text tower. They were trained from random initialisation for 3.6B examples seen, with a batch size of 16384. A cosine learning rate schedule was used, with an initial linear warmup for 10% of steps up to a peak learning rate of 0.001.



Batch Size	450 M		900 M		3 B		18 B	
	sigmoid	softmax	sigmoid	softmax	sigmoid	softmax	sigmoid	softmax
512	72.5	69.5	75.0	72.8	77.2	74.6	-	-
1 k	75.5	73.6	77.2	76.0	79.6	77.9	-	-
2 k	77.1	76.3	79.3	78.1	81.3	80.1	82.2	81.2
4 k	79.2	78.3	80.8	79.8	82.4	81.2	83.0	82.0
8 k	80.8	79.7	82.0	81.0	83.1	82.6	83.6	83.1
16 k	81.2	81.2	82.7	82.1	83.8	83.5	84.2	84.1
32 k	81.9	81.4	83.1	82.7	84.2	84.0	84.6	84.4
64 k	81.6	81.6	83.0	82.8	84.3	84.1	84.7	84.4
128 k	80.5	80.0	83.1	83.2	84.2	84.4	84.7	84.6
256 k	72.8	72.2	82.1	81.7	84.3	84.2	84.7	84.6
1024 k	-	-	-	-	-	-	84.7	-

Table 5: **SigLiT zero-shot accuracy (%) on the ImageNet benchmark.** Both the sigmoid loss and the softmax loss baseline are presented. Extensive experiments are performed on multiple train examples seen (450 M, 900 M, 3 B, 18 B) and train batch sizes (from 512 to 1 M).

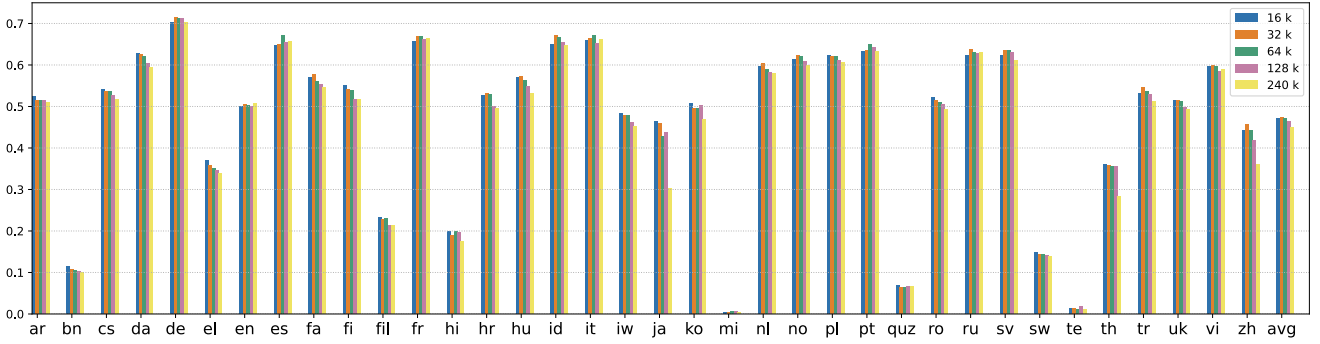


Figure 8: **Zero-shot image to text retrieval results on the Crossmodal-3600 benchmark.** Multiple batch sizes are reported.

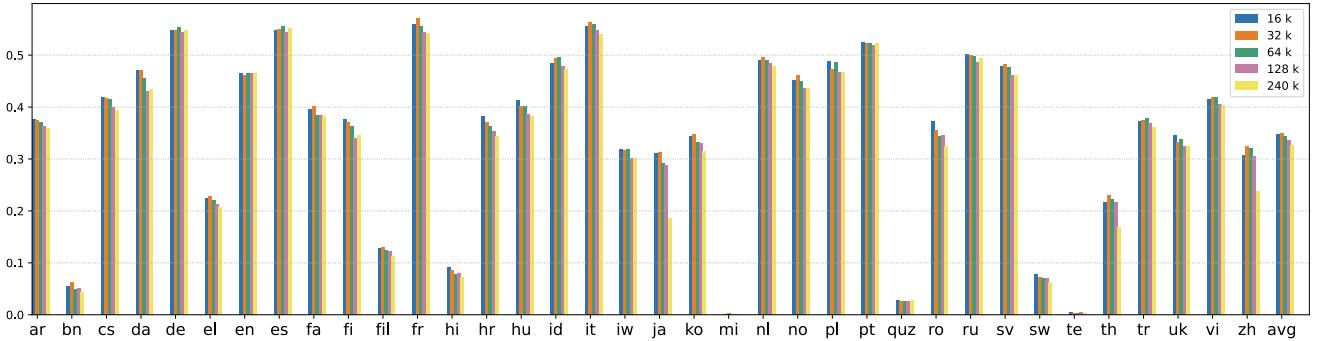


Figure 9: **Zero-shot text to image retrieval results on the Crossmodal-3600 benchmark.** Multiple batch sizes are reported.

Lang.	Image-to-text					Text-to-image				
	16 k	32 k	64 k	128 k	240 k	16 k	32 k	64 k	128 k	240 k
ar	52.39	51.33	51.50	51.53	51.06	37.61	37.37	37.14	36.31	35.98
bn	11.39	10.78	10.44	10.31	9.89	5.53	6.25	4.94	5.14	4.36
cs	54.08	53.69	53.67	52.75	51.78	41.82	41.64	41.46	39.93	39.38
da	62.72	62.42	62.00	60.42	59.33	47.01	47.01	45.55	43.03	43.47
de	70.33	71.42	71.19	71.11	70.25	54.70	54.83	55.36	54.31	54.71
el	36.94	35.81	35.06	34.53	33.81	22.42	22.78	22.00	21.25	20.79
en	50.11	50.53	50.22	49.94	50.67	46.46	46.21	46.55	46.60	46.60
es	64.69	64.94	67.19	65.31	65.61	54.81	55.04	55.51	54.47	55.24
fa	57.03	57.75	56.06	55.28	54.64	39.61	40.15	38.43	38.36	38.30
fi	54.94	54.08	53.78	51.69	51.67	37.70	37.14	36.38	33.98	34.50
fil	23.22	22.75	22.92	21.39	21.22	12.83	12.93	12.41	12.19	11.34
fr	65.69	66.92	66.97	66.14	66.47	55.92	57.08	55.50	54.39	54.29
hi	19.86	18.81	19.89	19.53	17.36	9.09	8.55	7.86	8.06	7.28
hr	52.67	53.03	52.97	49.92	49.58	38.16	37.09	36.37	35.25	34.33
hu	56.97	57.11	56.33	54.83	53.03	41.37	40.20	40.22	38.55	38.25
id	64.83	67.06	66.56	65.39	64.72	48.53	49.42	49.49	47.82	47.29
it	65.86	66.42	67.11	65.25	66.08	55.52	56.39	55.85	54.75	54.11
iw	48.36	47.86	47.72	46.06	45.25	31.78	31.76	31.89	30.08	30.12
ja	46.42	45.94	42.89	43.72	30.17	31.04	31.32	29.21	28.87	18.50
ko	50.78	49.53	49.44	50.22	46.78	34.44	34.72	33.15	33.06	31.54
mi	0.36	0.42	0.58	0.56	0.42	0.16	0.22	0.19	0.19	0.19
nl	59.56	60.36	58.94	58.31	57.86	48.95	49.55	48.95	48.37	47.88
no	61.36	62.39	61.97	60.89	59.86	45.25	46.21	45.04	43.53	43.71
pl	62.19	62.03	61.97	61.11	60.50	48.80	47.36	48.70	46.79	46.72
pt	63.14	63.61	64.89	64.31	63.25	52.41	52.34	52.30	51.93	52.37
quz	6.78	6.42	6.36	6.64	6.67	2.74	2.57	2.67	2.69	2.79
ro	52.06	51.44	50.97	50.58	49.31	37.20	35.60	34.34	34.52	32.50
ru	62.22	63.64	63.08	62.69	63.08	50.11	49.89	49.71	48.61	49.31
sv	62.33	63.53	63.53	63.06	61.19	47.89	48.18	47.64	46.17	46.16
sw	14.83	14.42	14.31	14.17	13.78	7.81	7.17	7.11	6.94	6.34
te	1.25	1.25	1.19	1.69	1.06	0.40	0.29	0.32	0.47	0.32
th	36.11	35.78	35.64	35.56	28.33	21.58	23.08	22.22	21.62	16.76
tr	53.08	54.50	53.72	52.94	51.25	37.33	37.38	37.81	36.97	36.08
uk	51.42	51.50	51.17	49.86	49.22	34.54	33.21	33.79	32.49	32.39
vi	59.58	59.78	59.53	58.53	58.83	41.43	41.92	41.85	40.63	40.26
zh	44.11	45.67	44.11	41.92	36.08	30.74	32.45	32.05	30.61	23.72
avg	47.21	47.36	47.11	46.34	45.00	34.82	34.87	34.44	33.58	32.72

Table 6: **Image-to-text and text-to-image zero-shot retrieval results on all 36 languages of Crossmodal-3600**, with mSigLIP models trained at different batch sizes for 30 B examples seen.

# Summary

- Contrastive pre-training using weak supervision from image-text pairs found on the web is becoming the go-to method for obtaining generic computer vision backbones, slowly replacing pre-training on large labelled multi-class datasets. The high-level idea is to simultaneously learn an aligned representation space for images and texts using paired data. CLIP from OpenAI has shown that this works on a scale.
- The standard recipe to pre-train such models leverages the image-text contrastive objective. It aligns the image and text embeddings for matching (positive) image-text pairs while making sure that unrelated (negative) image-text pairs are dissimilar in the embedding space. This is achieved via a batch-level softmax based contrastive loss, applied twice to normalize the pairwise similarity scores across all images, then all texts as shown below

$$-\frac{1}{2|\mathcal{B}|} \sum_{i=1}^{|\mathcal{B}|} \left( \overbrace{\log \frac{e^{t\mathbf{x}_i \cdot \mathbf{y}_i}}{\sum_{j=1}^{|\mathcal{B}|} e^{t\mathbf{x}_i \cdot \mathbf{y}_j}}}^{\text{image} \rightarrow \text{text softmax}} + \overbrace{\log \frac{e^{t\mathbf{x}_i \cdot \mathbf{y}_i}}{\sum_{j=1}^{|\mathcal{B}|} e^{t\mathbf{x}_j \cdot \mathbf{y}_i}}}^{\text{text} \rightarrow \text{image softmax}} \right)$$

- The above loss works fine but has two major disadvantages when training models at scale:
  - Data parallelism is a typical setting used for training such models. Computing the loss when data is split across  $D$  devices necessitates gathering all embeddings with expensive all-gathers
  - The  $|\mathcal{B}| \times |\mathcal{B}|$  pairwise similarity matrix is memory intensive
- To overcome these limitations, the authors propose a simpler yet effective alternative: The sigmoid loss
- The proposed sigmoid loss does not require any operation across the full batch and hence greatly simplifies the distributed loss implementation and boosts efficiency. It also conceptually decouples the batch size from the definition of the task.
- The sigmoid-based loss processes every image-text pair independently, turning the problem statement into a standard binary classification problem on the dataset of all pair combinations, with a positive labels for the matching pairs  $(I_i, T_i)$  and negative labels for all other pairs  $(I_i, T_j)$  ( $j \neq i$ ). It is defined as follows:

$$-\frac{1}{|\mathcal{B}|} \sum_{i=1}^{|\mathcal{B}|} \sum_{j=1}^{|\mathcal{B}|} \underbrace{\log \frac{1}{1 + e^{z_{ij}(-t\mathbf{x}_i \cdot \mathbf{y}_j + b)}}}_{\mathcal{L}_{ij}}$$

where  $z_{ij}$  is the label for a given image and text input, which equals 1 if they are paired and -1 otherwise.

- The sigmoid loss is memory-efficient, fast, and numerically stable. If the model is replicated over  $D$  devices with a global batch size of  $|\mathcal{B}|$ , such that each device operates on a batch size of  $b = |\mathcal{B}| / D$ , the loss can be computed efficiently as shown below:

$$-\frac{1}{|\mathcal{B}|} \underbrace{\sum_{d_i=1}^D}_{\text{A: } \forall \text{ device } d_i} \underbrace{\sum_{d_j=1}^D}_{\text{B: swap negs across devices}} \underbrace{\sum_{i=bd_i}^{b(d_i+1)} \sum_{j=bd_j}^{b(d_j+1)} \mathcal{L}_{ij}}_{\text{C: per device loss}} \underbrace{\substack{\text{all local} \\ \text{positives}}}_{\substack{\text{negs from} \\ \text{next device}}}$$

- The above equation is an efficient chunked implementation, and can be broken down as:
  - Compute the loss corresponding to the positive pair, and  $(b-1)$  negative pairs
  - Then permute representations across devices, so each device takes negatives from its neighbouring device (next iteration of sum B).
  - The loss is then calculated with respect to this chunk (sum C). This is done independently in each device, such that each device computes the loss with respect to its local batch  $b$
  - Losses can then simply be summed across all devices (sum A)
  - Individual collective permutes (for sum B) are fast (and indeed  $D$  collective permutes is typically faster than two all-gathers between  $D$  devices), and the memory cost at any given moment is reduced from  $|\mathcal{B}|^2$  to  $b^2$  (for sum C).

## Results and other findings

- Scaling batch size to the limit

- When the batch size is smaller than 16 k, sigmoid loss outperforms softmax loss by a large margin. With growing batch sizes, the softmax loss quickly catches up and potentially slightly underperforms sigmoid loss with a large enough batch size.
- The authors successfully trained a model with a batch size of 1 M but found out that the performance saturates at a batch size of 32 K, and further scaling up gives only a minor boost to the performance peaking around a batch size of 256 K
- 262 k batch size significantly outperforms smaller 8 k batch size when trained for a sufficiently long time. This is mostly because large batch size leads to the fewer absolute number of update steps and thus needs more time to ramp up.

- SigLIP: Sigmoid loss is beneficial for language- image pre-training

- The authors pretrain SigLIP models on the WebLI dataset using only English image and text pairs. B/16 ViT for image embeddings and B-sized transformer for text embeddings.
- The text is tokenized by a 32 k vocabulary sentencepiece tokenizer trained on the English C4 dataset, and a maximum of 16 text tokens are kept.
- With less than 32 k batch size, SigLIP outperforms CLIP baselines with the standard softmax loss.
- Because the loss function is now more memory efficient, the authors were able to fit a batch size of 4096 with base SigLIP compared to a batch size of 2048 with a corresponding CLIP model
- As batch size increases, the gap between the sigmoid and the softmax losses diminish. SigLIP performs best at batch size 32 k, whereas the softmax loss required 98 k for optimal performance and still didn't outperform the sigmoid based variant. Scaling further, a larger batch size like 307 k hurts both losses.

- mSigLIP: Multi-lingual pre-training

- The authors scale up the training data by keeping all the 100 languages from the WebLI dataset
- Large vocab size proved to be more performative but the token embeddings become a bottleneck with large vocab as it requires a  $N \times W$  token embedding lookup table to train the multilingual model, where  $N$  is the vocabulary size mentioned above and  $W$  is the embedding dimension of the text model.
- To save memory, the authors proposed a bottlenecked embedding:  $N \times K$  embedding matrix and additional  $K \times W$  projection, where the bottleneck  $K$  is much smaller than  $W$ . (👉👉👉👉)
- mSigLIP sets the new SOTA on XM3600 text to image retrieval task, with only a Base size model at 34.9%, which is more than 6% higher than the previously reported result 28.5%

- Stabilizing large batch training

- Language-image pretraining becomes more unstable with large batch sizes. The reason for these instabilities is large spikes in the gradient norms, which translate to large magnitude changes in the weights that may destabilize the training process
- $\beta_2$  in Adam and Adafactor was the culprit (not surprised though 😂). The default value was changed from 0.999 to 0.95. This provided a good enough stable training, though some spikes were still observed in the loss curve



- Negative ratio in the sigmoid loss

- For a batch size  $|B|$ , the batch contains  $|B|$  positive pairs, but  $|B|^2 - |B|$  negative examples. In the modest batch size of 16 k, there are actually 268 M negative examples for only 16 k positive ones. To study the effects of mini-batch composition, and distribution, the authors performed some experiments by masking out negatives to vary the mini-batch composition:
  - **Random**: Randomly choose negative pairs to mask -> Performance deteriorates!
  - **Hard**: Keep hardest negative pairs (highest loss) -> Maintain the quality
  - **Easy**: Keep easiest negatives pairs (lowest loss). -> Doesn't work or meh!
  - **Hard + matching total pairs seen**: Masking examples while training for a fixed number of steps does decrease the total number of pairs seen during training -> Maintain the quality
- This study confirms that:
  - The imbalance does not seem to be a major reason for concern
  - Coming up with an efficient way of including more negative examples can be promising but is not trivial.

- Bias term in sigmoid loss

- The sigmoid loss is heavily imbalanced because of the ratio of positive-negative ratio in a mini batch
- To alleviate this, the authors introduced an additional learnable bias term  $b$  similar to the temperature  $t$ . They initialized  $t'$  and  $b$  to 10 and -10 respectively.
- Enabling the bias term with a -10 initialization consistently improves performance across all tasks. This is because the bias term ensures that the training starts close to the prior, preventing dramatic over-correction in early optimization.
- In contrast, a randomly chosen bias term initialization, such as the zero initialization fails to address the overcorrection issue, leading to significantly worse results.

- Label noise robustness

- The authors trained the training data using one of the following methods:
  - **Image**: With probability  $p$ , replace the image with uniform random noise.
  - **Text**: With probability  $p$ , replace tokenized text with a new sequence of randomly sampled tokens, up to some (sampled) sequence length.
  - **Batch alignment**: Randomly shuffle the ordering of  $p\%$  of the batch.
  - **Image & text**: Independently apply (1) and (2), each with probability  $p$ .
  - **Image, text & batch**: Alongside (4), also shuffle fraction  $p$  of alignments.
- Models trained with sigmoid loss are increasingly robust to all kinds of added noise.