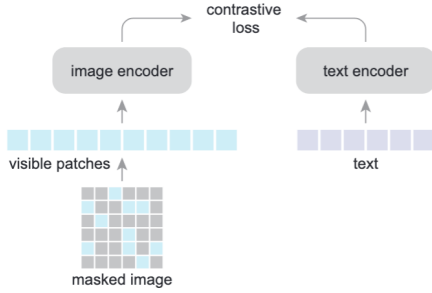


FLIP

A Preprint

Authors introduce FLIP, which is a simple method for efficient CLIP training. Idea behind it is to randomly remove patches/mask. FLIP trains $>3\times$ faster in wall-clock time for reaching similar accuracy as its CLIP counterpart; with the same number of epochs, FLIP reaches higher accuracy than its CLIP counterpart. The whole idea was inspired by MAE(Masked Autoencoder) which sparsely applies ViT encoder to visible content.



The intuition behind method is reducing computation so throughout training more image-text pairs can be used under the same time and also have a contrastive objective over a larger batch under the same memory constraint. As for Image masking authors adapt ViT; they mask 50-75 % of image and feed it into ViT encoder (by it meant the visible part). It is also possible to encode text the same way as was presented with image masking. however, as the text encoder is smaller, speeding it up does not lead to a better overall trade-off.

Unmasking: The simplest strategy is using the the same encoder pre-trained on masked images with just one simple setting: masking ratio 0 %

The training is performed on LAION-400M and evaluate zeroshot accuracy on ImageNet-1K validation. As results show, batchsize has a major impact on accuracy:

mask	batch	FLOPs	time	acc.	batch	mask 50%	mask 75%	text mask	text len	time	acc.
0%	16k	1.00×	1.00×	68.6	16k	68.5	65.8	baseline, 0%	32	1.00×	68.2
50%	32k	0.52×	0.50×	69.6	32k	69.6	67.3	random, 50%	16	0.92×	66.0
75%	64k	0.28×	0.33×	68.2	64k	70.4	68.2	prioritized, 50%	16	0.92×	67.8

(a) **Image masking** yields higher or comparable accuracy and speeds up training. Entries are subject to the same memory limit.

(b) **Batch size.** A large batch has big gains over smaller batches.

(c) **Text masking** performs decently, but the speed gain is marginal as its encoder is smaller. Here the image masking ratio is 75%.

mask 50%	mask 75%	mask 50%	mask 75%	mask 50%	mask 75%
w/ mask	66.4	60.9	baseline	69.6	68.2
w/ mask, ensemble	68.1	65.1	+ tuning	70.1	69.5
w/o mask	69.6	68.2	+ MAE	69.4	67.9

(d) **Inference unmasking.** Inference on intact images performs strongly even without tuning.

(e) **Unmasked tuning.** The distribution shift by masking is reduced by a short tuning.

(f) **Reconstruction.** Adding the MAE reconstruction loss has no gain.

Table 1 **Zeroshot ablation experiments.** Distribution is on LAION-400M for 8 epochs. Evaluated by zeroshot classification accuracy

FLIP and CLIP both use the same contrastive loss:

CLIP uses a symmetric contrastive loss based on InfoNCE, defined as:

$$\mathcal{L} = \frac{1}{2N} \sum_{i=1}^N \left[-\log \frac{\exp(\text{sim}(I_i, T_i)/\tau)}{\sum_{j=1}^N \exp(\text{sim}(I_i, T_j)/\tau)} - \log \frac{\exp(\text{sim}(T_i, I_i)/\tau)}{\sum_{j=1}^N \exp(\text{sim}(T_i, I_j)/\tau)} \right]$$

where:

- N — batch size;
- I_i — image embedding;
- T_i — text embedding;
- $\text{sim}(x, y) = \frac{x \cdot y}{\|x\| \|y\|}$ — cosine similarity between vectors;
- τ — learnable temperature parameter.