# CLIP - Contrastive Language-Image Pretraining

General Info:

- OpenAI - 2021
- Multimodal (Image → Text)
  - One of (if not) the first model
  - Late fusion approach
- (very) Good paper

Claimed Advantages (compare to ResNet-50 - CNN):

- Match ResNet50 performance
- + Generality - zero-shot 30 different datasets
- + Scalability
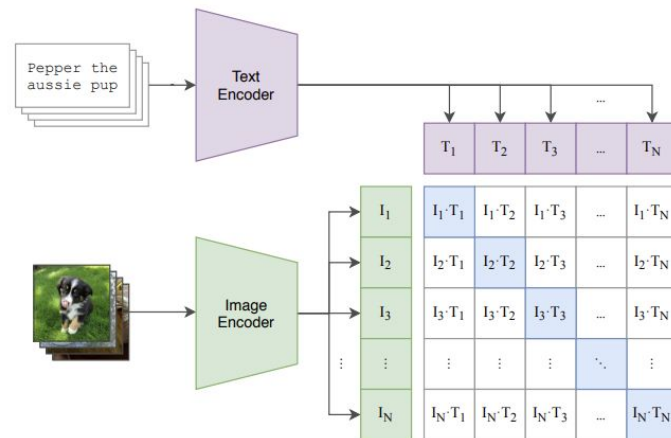
# CLIP - Contrastive Language-Image Pretraining

- ## Dataset:
  - WebImageNet (WIT): open-source internet-search with queries - **400Mil**
  - Reference - ImageNet has ~15Mil
  - *Train from scratch*

- ## Architecture (for pretrain):

- ## Loss function for batch of N pairs:
  - CE log(cosine (dis)similarity of NxN)



(1) Contrastive pre-training

# CLIP - Contrastive Language-Image Pretraining
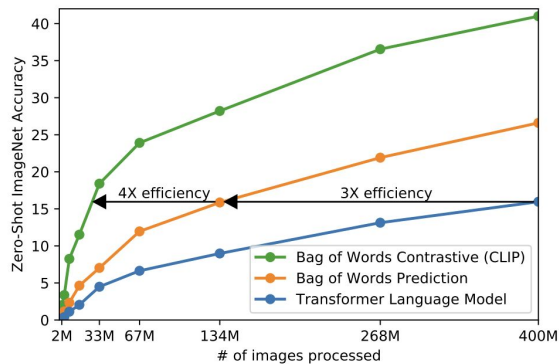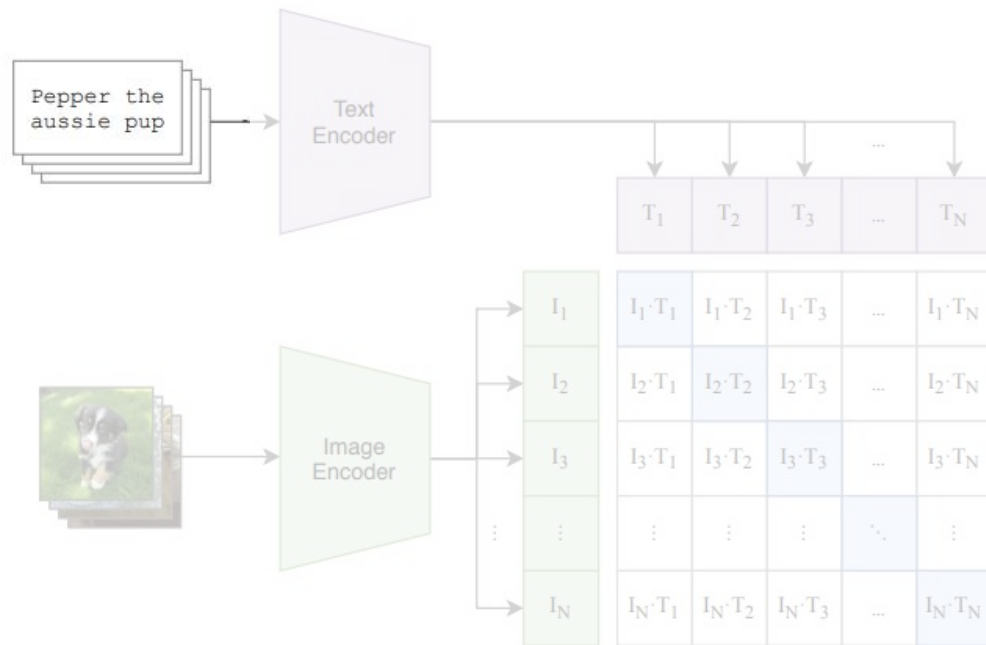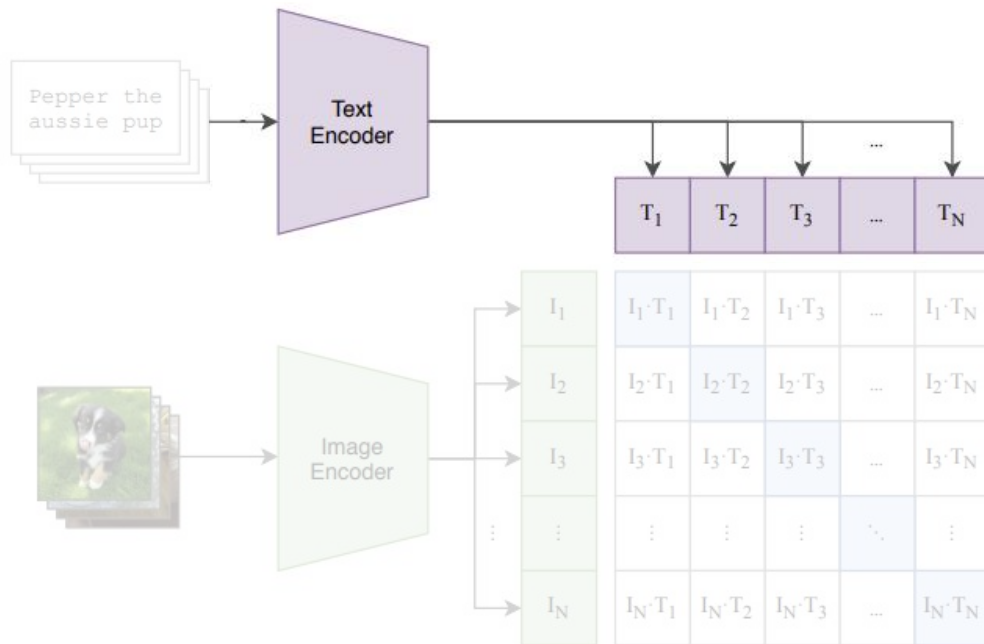
1. Text representation



*Figure 2.* **CLIP is much more efficient at zero-shot transfer than our image caption baseline.** Although highly expressive,

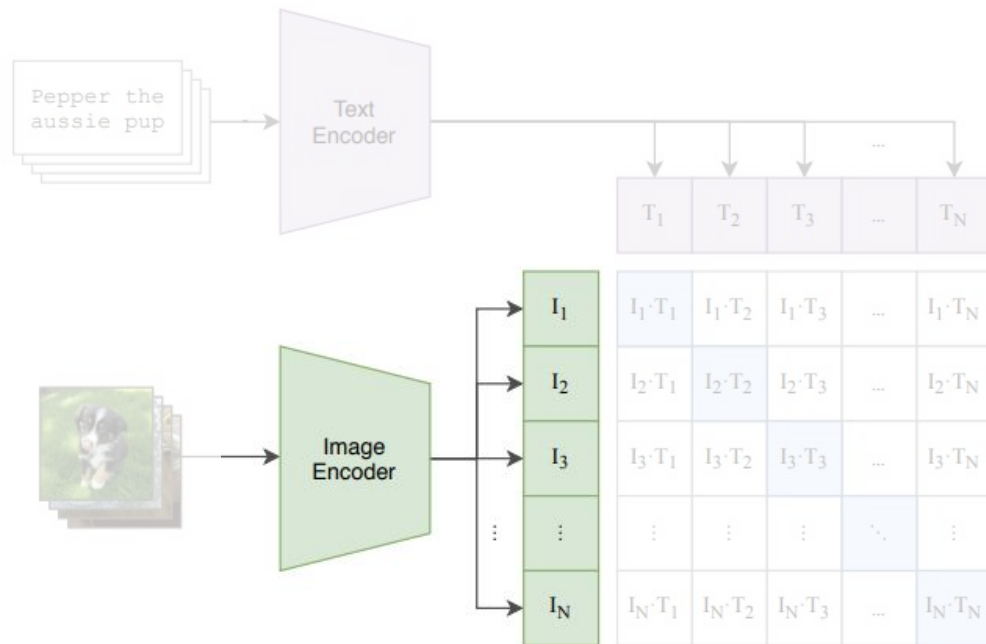# CLIP - Contrastive Language-Image Pretraining

2. Text Encoder
→ Transformer architecture + Linear
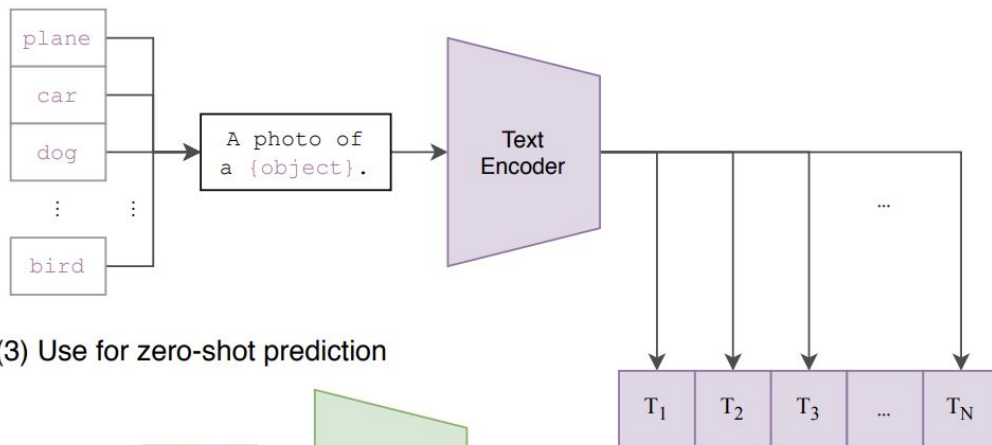
# CLIP - Contrastive Language-Image Pretraining

3. Image Encoder

- Candidate 1: ResNet50 with attention-pooling
  - 3 sizes: 4x, 16x, and 64x
- Candidate 2: ViT (2020)
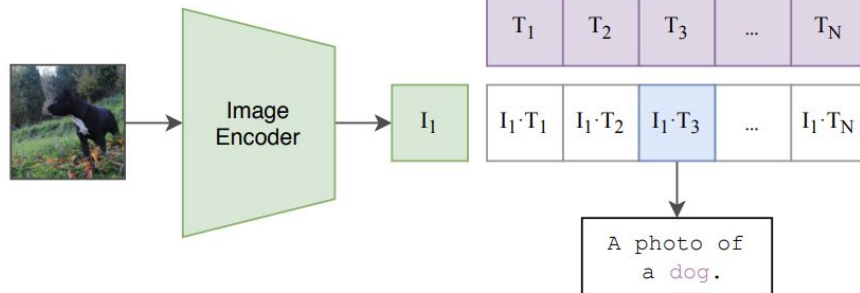  - 3 sizes: ViT-B/32, ViT-B/16, ViT-L/14

# CLIP - Contrastive Language-Image Pretraining



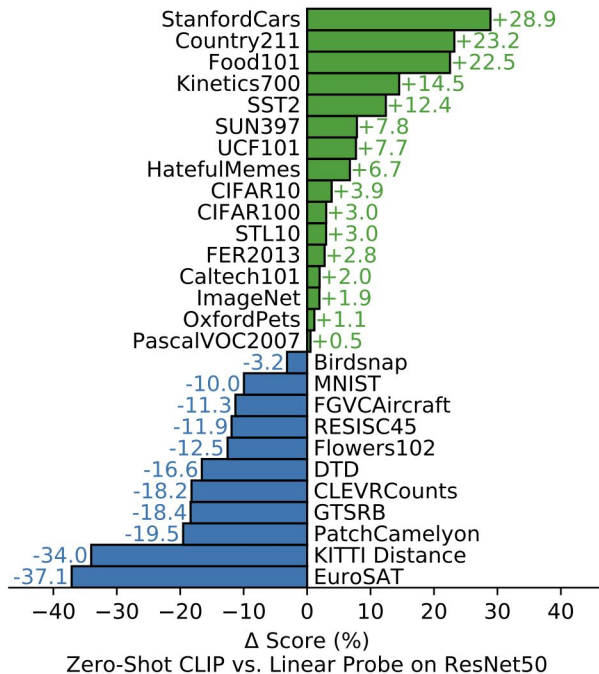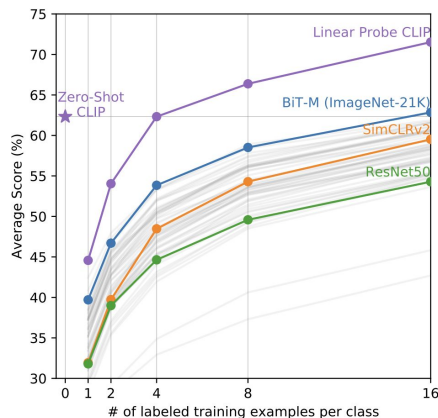(2) Create dataset classifier from label text

plane
car
dog
⋮
bird

A photo of a {object}.

Text Encoder

T₁  T₂  T₃  ...  T_N

(3) Use for zero-shot prediction

Image Encoder

I₁

I₁·T₁  I₁·T₂  I₁·T₃  ...  I₁·T_N

A photo of a dog.

# CLIP - Contrastive Language-Image Pretraining

Result (zero-shot)


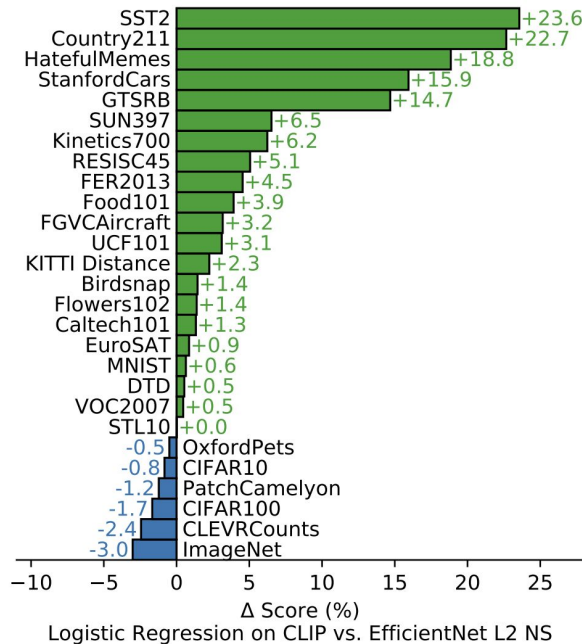
Δ Score (%)
Zero-Shot CLIP vs. Linear Probe on ResNet50

# CLIP - Contrastive Language-Image Pretraining

Result (both linear probe)



Figure 6. **Zero-shot CLIP outperforms few-shot linear probes.** Zero-shot CLIP matches the average performance of a 4-shot linear classifier trained on the same feature space and nearly matches the best results of a 16-shot linear classifier across publicly available models. For both BiT-M and SimCLRv2, the best performing model is highlighted. Light gray lines are other models in the eval suite. The 20 datasets with at least 16 examples per class were used in this analysis.



Figure 11. **CLIP's features outperform the features of the best ImageNet model on a wide variety of datasets.** Fitting a linear classifier on CLIP's features outperforms using the Noisy Student EfficientNet-L2 on 21 out of 27 datasets.

# CLIP - Contrastive Language-Image Pretraining

- Limitations:
  - Not yet beats SOTAs with just baseline model
  - Bad at fine-grained classification tasks, abstract and systematic tasks
  - Learning social bias from dataset
  - Out-of-distribution problem
  - No available dataset for image captioning at the time
- Personal critics:
  - Bag-of-words and Winogrounds (link)
  - Modality overpower

# Flamingo - Multimodal Few-shot Learning

**General Info:**

- DeepMind - NeurIPS 2022
- Multimodal
  - Image + Video (partial) ↔ Text
  - Frozen modules
- Designed for few-shot learning tasks

# Flamingo - Multimodal Few-shot Learning

pretrained and frozen
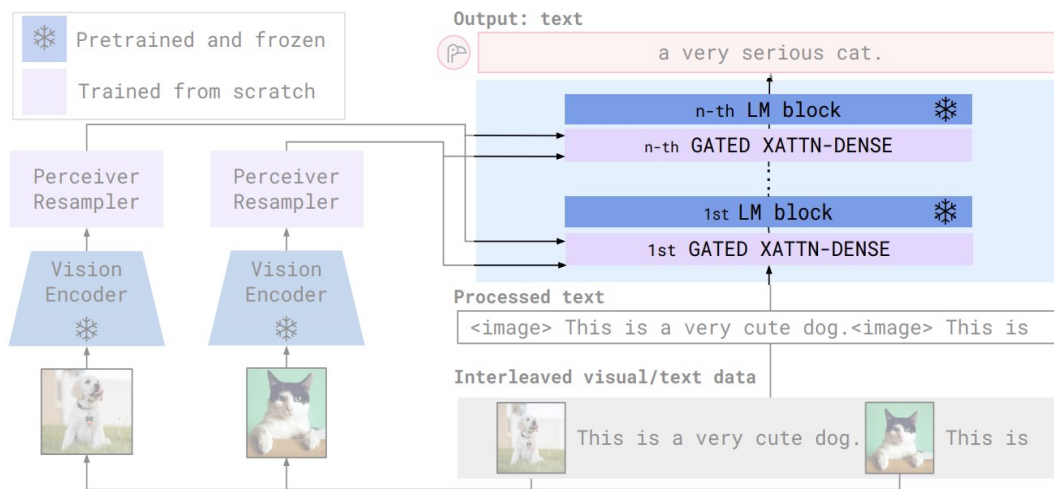Using contrastive text-image data
NormalizerFree ResNet (NFNet) - F6



Figure 3: **Flamingo architecture overview.** Flamingo is a family of visual language models (VLMs) that take as input visual data interleaved with text and produce free-form text as output.

# Flamingo - Multimodal Few-shot Learning

- pretrained and frozen: Transformer decoder
- From scratch: gated cross-attention dense - with tanh(a) gating mechanism



Figure 3: **Flamingo architecture overview.** Flamingo is a family of visual language models (VLMs) that take as input visual data interleaved with text and produce free-form text as output.

| Method | FT | Shot | OKVQA (I) | VQAv2 (I) | COCO (I) | MSVDQA (V) | VATEX (V) | VizWiz (I) | Flick30K (I) | MSRVTTQA (V) | iVQA (V) | YouCook2 (V) | STAR (V) | VisDial (I) | TextVQA (I) | NextQA (I) | HatefulMemes (I) | RareAct (V) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Zero/Few shot SOTA | ✗ | | [34] | [114] | [124] | [58] | - | - | - | [58] | [135] | - | [143] | [79] | - | - | [85] | [85] |
| | | | 43.3 | 38.2 | 32.2 | 35.2 | | | | 19.2 | 12.2 | | 39.4 | 11.6 | | | 66.1 | 40.7 |
| | | (X) | (16) | (4) | (0) | (0) | | | | (0) | (0) | | (0) | (0) | | | (0) | (0) |
| *Flamingo*-3B | ✗ | 0 | 41.2 | 49.2 | 73.0 | 27.5 | 40.1 | 28.9 | 60.6 | 11.0 | 32.7 | 55.8 | 39.6 | 46.1 | 30.1 | 21.3 | 53.7 | 58.4 |
| | ✗ | 4 | 43.3 | 53.2 | 85.0 | 33.0 | 50.0 | 34.0 | 72.0 | 14.9 | 35.7 | 64.6 | 41.3 | 47.3 | 32.7 | 22.4 | 53.6 | - |
| | ✗ | 32 | 45.9 | 57.1 | 99.0 | 42.6 | 59.2 | 45.5 | 71.2 | 25.6 | 37.7 | 76.7 | 41.6 | 47.3 | 30.6 | 26.1 | 56.3 | - |
| *Flamingo*-9B | ✗ | 0 | 44.7 | 51.8 | 79.4 | 30.2 | 39.5 | 28.8 | 61.5 | 13.7 | 35.2 | 55.0 | 41.8 | 48.0 | 31.8 | 23.0 | 57.0 | 57.9 |
| | ✗ | 4 | 49.3 | 56.3 | 93.1 | 36.2 | 51.7 | 34.9 | 72.6 | 18.2 | 37.7 | 70.8 | **42.8** | 50.4 | 33.6 | 24.7 | 62.7 | - |
| | ✗ | 32 | 51.0 | 60.4 | 106.3 | 47.2 | 57.4 | 44.0 | 72.8 | 29.4 | 40.7 | 77.3 | 41.2 | 50.4 | 32.6 | 28.4 | 63.5 | - |
| *Flamingo* | ✗ | 0 | 50.6 | 56.3 | 84.3 | 35.6 | 46.7 | 31.6 | 67.2 | 17.4 | 40.7 | | 39.7 | 52.0 | 35.0 | 26.7 | 46.4 | **60.8** |
| | ✗ | 4 | 57.4 | 63.1 | 103.2 | 41.7 | 56.0 | 39.6 | 75.1 | 23.9 | 44.1 | 74.5 | 42.4 | **55.6** | 36.5 | 30.8 | 68.6 | - |
| | ✗ | 32 | **57.8** | **67.6** | **113.8** | **52.3** | **65.1** | **49.8** | **75.4** | **31.0** | **45.3** | **86.8** | 42.2 | **55.6** | 37.9 | **33.5** | 70.0 | - |
| Pretrained FT SOTA | ✔ | | 54.4 | 80.2 | 143.3 | 47.9 | 76.3 | 57.2 | 67.4 | 46.8 | 35.4 | 138.7 | 36.7 | 75.2 | 54.7 | 25.2 | 79.1 | |
| | | | [34] | [140] | [124] | [28] | [153] | [65] | [150] | [51] | [135] | [132] | [128] | [79] | [137] | [129] | [62] | - |
| | | (X) | (10K) | (444K) | (500K) | (27K) | (500K) | (20K) | (30K) | (130K) | (6K) | (10K) | (46K) | (123K) | (20K) | (38K) | (9K) | |

Table 1: **Comparison to the state of the art.** A *single* Flamingo model reaches the state of the art

| Method | VQAV2 | | COCO | VATEX | VizWiz | | MSRVTTQA | VisDial | | YouCook2 | TextVQA | | HatefulMemes |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | test-dev | test-std | test | test | test-dev | test-std | test | valid | test-std | valid | valid | test-std | test seen |
| 🦩 32 shots | 67.6 | - | 113.8 | 65.1 | 49.8 | - | 31.0 | 56.8 | - | 86.8 | 36.0 | - | 70.0 |
| 🦩 Fine-tuned | **82.0** | **82.1** | 138.1 | **84.2** | **65.7** | 65.4 | **47.4** | 61.8 | 59.7 | 118.6 | **57.1** | 54.1 | **86.6** |
| SotA | 81.3† | 81.3† | **149.6†** | 81.4† | 57.2† | 60.6† | 46.8 | **75.2** | **75.4†** | **138.7** | 54.7 | **73.7** | 84.6† |
| | [133] | [133] | [119] | [153] | [65] | [65] | [51] | [79] | [123] | [132] | [137] | [84] | [152] |

Table 2: **Comparison to SotA when fine-tuning *Flamingo*.** We fine-tune *Flamingo* on all nine tasks where *Flamingo* does not achieve SotA with few-shot learning. *Flamingo* sets a new SotA on five of them, outperfoming methods (marked with †) that use tricks such as model ensembling or domain-specific metric optimisation (e.g., CIDEr optimisation).

| Ablated setting | *Flamingo*-3B original value | Changed value | Param. count↓ | Step time↓ | COCO CIDEr↑ | OKVQA top1↑ | VQAv2 top1↑ | MSVDQA top1↑ | VATEX CIDEr↑ | Overall score↑ |
|---|---|---|---|---|---|---|---|---|---|---|
| | **_Flamingo_-3B model** | | 3.2B | 1.74s | 86.5 | 42.1 | 55.8 | 36.3 | 53.4 | **70.7** |
| **(i)** Training data | All data | w/o Video-Text pairs | 3.2B | 1.42s | 84.2 | 43.0 | 53.9 | 34.5 | 46.0 | 67.3 |
| | | w/o Image-Text pairs | 3.2B | 0.95s | 66.3 | 39.2 | 51.6 | 32.0 | 41.6 | 60.9 |
| | | Image-Text pairs→ LAION | 3.2B | 1.74s | 79.5 | 41.4 | 53.5 | 33.9 | 47.6 | 66.4 |
| | | w/o M3W | 3.2B | 1.02s | 54.1 | 36.5 | 52.7 | 31.4 | 23.5 | 53.4 |
| **(ii)** Optimisation | Accumulation | Round Robin | 3.2B | 1.68s | 76.1 | 39.8 | 52.1 | 33.2 | 40.8 | 62.9 |
| **(iii)** Tanh gating | ✓ | ✗ | 3.2B | 1.74s | 78.4 | 40.5 | 52.9 | 35.9 | 47.5 | 66.5 |
| **(iv)** Cross-attention architecture | GATED XATTN-DENSE | VANILLA XATTN | 2.4B | 1.16s | 80.6 | 41.5 | 53.4 | 32.9 | 50.7 | 66.9 |
| | | GRAFTING | 3.3B | 1.74s | 79.2 | 36.1 | 50.8 | 32.2 | 47.8 | 63.1 |
| **(v)** Cross-attention frequency | Every | Single in middle | 2.0B | 0.87s | 71.5 | 38.1 | 50.2 | 29.1 | 42.3 | 59.8 |
| | | Every 4th | 2.3B | 1.02s | 82.3 | 42.7 | 55.1 | 34.6 | 50.8 | 68.8 |
| | | Every 2nd | 2.6B | 1.24s | 83.7 | 41.0 | 55.8 | 34.5 | 49.7 | 68.2 |
| **(vi)** Resampler | Perceiver | MLP | 3.2B | 1.85s | 78.6 | 42.2 | 54.7 | 35.2 | 44.7 | 66.6 |
| | | Transformer | 3.2B | 1.81s | 83.2 | 41.7 | 55.6 | 31.5 | 48.3 | 66.7 |
| **(vii)** Vision encoder | NFNet-F6 | CLIP ViT-L/14 | 3.1B | 1.58s | 76.5 | 41.6 | 53.4 | 33.2 | 44.5 | 64.9 |
| | | NFNet-F0 | 2.9B | 1.45s | 73.8 | 40.5 | 52.8 | 31.1 | 42.9 | 62.7 |
| **(viii)** Freezing LM | ✓ | ✗ (random init) | 3.2B | 2.42s | 74.8 | 31.5 | 45.6 | 26.9 | 50.1 | 57.8 |
| | | ✗ (pretrained) | 3.2B | 2.42s | 81.2 | 33.7 | 47.4 | 31.0 | 53.9 | 62.7 |

Table 3: **Ablation studies.** Each row should be compared to the baseline Flamingo run (top row). Step time measures the time spent to perform gradient updates on all training datasets.