

# Nomic Embed Vision: Expanding the Latent Space

Zach Nussbaum  
zach@nomic.ai

Brandon Duderstadt  
brandon@nomic.ai

Andriy Mulyar  
andriy@nomic.ai

## Abstract

This technical report describes the training of nomic-embed-vision, a highly performant, open-code, open-weights image embedding model that shares the same latent space as nomic-embed-text. Together, nomic-embed-vision and nomic-embed-text form the first unified latent space to achieve high performance across vision, language, and multimodal tasks.

## 1 Introduction

Beginning with CLIP (Radford et al., 2021) and ALIGN (Jia et al., 2021), unsupervised multimodal encoders trained on large amounts of noisy web crawled data have shown impressive zero-shot capabilities across retrieval and classification tasks. These self supervised models are competitive with, and sometimes outperform, supervised baselines. However, these models are only optimized for multimodal tasks, and the text encoders perform poorly on text-only benchmarks like MTEB (Muennighoff et al., 2023; Koukounas et al., 2024).

Recently, Jina CLIP v1 (Koukounas et al., 2024) was introduced to address this issue. Unfortunately Jina CLIP does not achieve state of the art performance, failing to exceed jina-embeddings-v2 (Günther et al., 2024) on MTEB and OpenAI CLIP ViT B/16 (Radford et al., 2021) on Datacomp (Gadre et al., 2023) and Imagenet Zero-Shot Classification.

In this technical report, we introduce nomic-embed-vision, a highly performant vision encoder that is aligned to the latent space of nomic-embed-text. To train nomic-embed-vision, we adopt a similar training style to Locked Image Tuning (LiT) (Zhai et al., 2022), but instead freeze a high-performing text embedder and train a vision encoder from a pretrained checkpoint. This enables us to maintain the performance of nomic-

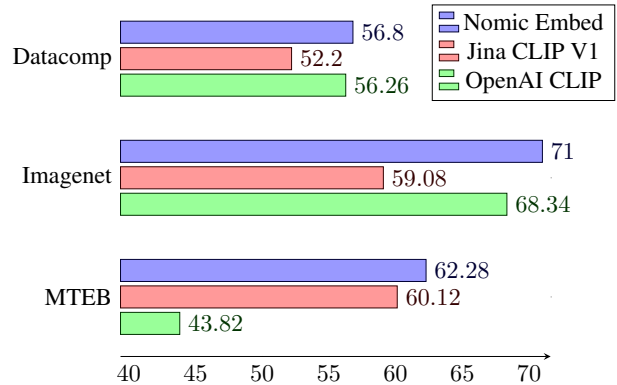


Figure 1: **Multimodal and Text Embedding Benchmark** Aggregate performance of Nomic Embed v1.5, OpenAI CLIP ViT B/16, and Jina CLIP v1 on text and multimodal benchmarks. Nomic Embed V1.5 is the only multimodal encoder to outperform OpenAI CLIP on multimodal and text benchmarks. X-axis units vary per benchmark suite. Imagenet is Imagenet Zero-Shot, Datacomp is a suite of 38 zero-shot multimodal evaluations, and MTEB evaluates performance of text embedding models.

embed-text as well as unlock new multimodal latent space capabilities. Together, nomic-embed-vision and nomic-embed-text form the first unified latent space to achieve high performance across vision, language, and multimodal tasks.

## 2 Related Work

Large scale noisy contrastive pretraining of image and text encoders was pioneered by Radford et al. (2021); Jia et al. (2021) using a large batch size and InfoNCE loss (van den Oord et al., 2019).

CLIP-style models are trained across a large noisy dataset created by crawling the web and extracting image-text pairs from webpages. These models are generally trained on billions of image-text pairs with a large batch size, which results in a massive pretraining compute requirement.

Radford et al. (2021) originally proposed evaluating CLIP models using zero-shot accuracy

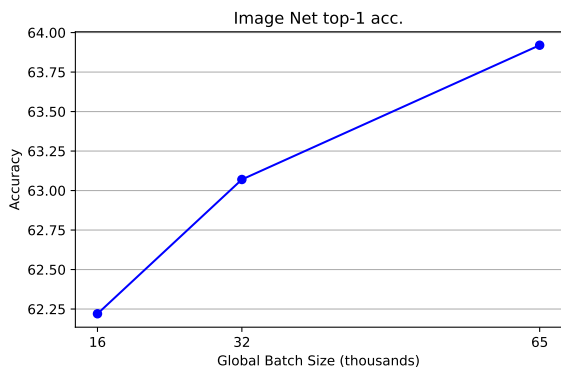


Figure 2: Imagenet Zero-Shot Top 1 Accuracy improves as we increase batch size in small scale experiments

across 27 datasets. Unfortunately, the lack of public information regarding the composition of the original web scale train set complicates this evaluation. To remedy this, [Gadre et al. \(2023\)](#) introduced Datacomp, an open benchmark to evaluate both CLIP-style models and their constituent training data mixes.

Taking inspiration from transfer learning, LiT ([Zhai et al., 2022](#)) and aligns a text encoder to a frozen pretrained image encoder, reducing the compute required to train a quality multimodal encoder. Three Towers ([Kossen et al., 2023](#)) improved upon LiT by introducing a third frozen pretrained image encoder and allowing the image and text encoders to take advantage of contrastive training as well as pretrained embeddings.

Imagebind ([Girdhar et al., 2023](#)) learns a joint embedding across many modalities by aligning modalities (e.g. audio) utilizing only image-paired data starting with a ViT-H from OpenCLIP ([Ilharco et al., 2021](#)).

Text embedding models are similarly trained contrastively on a large collection of text pairs and initializing with a pretrained transformer. [Reimers and Gurevych \(2019\)](#) train a pretrained BERT model contrastively for sentence similarity tasks. Since then, models such as E5 ([Wang et al., 2024](#)), GTE ([Li et al., 2023](#)), BGE ([Xiao et al., 2024](#)), InstructOR ([Su et al., 2023](#)), Jina ([Günther et al., 2024](#)), and Nomic ([Nussbaum et al., 2024](#)) train dual encoders in multiple stages.

MTEB ([Muennighoff et al., 2023](#)) aims to evaluate text embedding models across a suite of tasks including classification, retrieval, and semantic similarity.

### 3 Methods

Our goal is to learn a unified embedding space that performs well on multimodal tasks as well as unimodal text and image tasks. Contrastive Image Text Pretraining as introduced by [Radford et al. \(2021\)](#) leads to high performing multimodal models ([Radford et al., 2021](#); [Jia et al., 2021](#)). However, as shown in Figure 1 and noted by [Koukounas et al. \(2024\)](#), training only on these large scale datasets leads to poor general text embedding performance.

#### 3.1 Image Text Contrastive Training

Training CLIP-style models from scratch is expensive and requires large amounts of compute and data. [Zhai et al. \(2022\)](#) investigated ways to train CLIP models in a more efficient manner by freezing a pretrained vision encoder and training the text encoder from scratch. This methodology, which they named LiT, extends any pretrained vision encoder multimodal and zero-shot capabilities.

However, one downside of LiT is that freezing the image encoder prevents the vision encoder’s representations from being updated with signal from the text data. To remedy this, [Kossen et al. \(2023\)](#) proposes using a third frozen image tower to transfer representations to the main image and text encoders that are trained from scratch. This approach allows the encoders to be updated during training while also benefiting from the pretrained representations of the vision encoder. Three Towers outperforms LiT and CLIP-style models on retrieval tasks across initializations and pretraining datasets.

[Koukounas et al. \(2024\)](#) proposes a three stage contrastive training strategy to learn multimodal and text representations. In the first stage, they train the image and text encoders, initializing from EVA02 ([Fang et al., 2023b](#)) and a pretrained JinaBERT model, similar to ([Günther et al., 2024](#)) and optimize the image-text and text-text alignment. The second stage uses longer synthetic captions for further image-text alignment. The third stage introduces hard negatives to the text-text alignment to improve text embedding performance.

Similarly to ([Koukounas et al., 2024](#)), we aim to train general, high performing multimodal encoders. In this work, we adapt the LiT ([Zhai et al., 2022](#)) training recipe, and instead freeze the text

Vision Encoder	Pretrain	Supervised	IN-ZS	$I \rightarrow T$	$T \rightarrow I$	Mean R@1
Randomly Initialized	N/A	N/A	41.20	35.50	28.48	31.99
ViT (Dosovitskiy et al., 2021)	IN21k	Y	62.64	49.60	40.32	44.96
AugReg (Steiner et al., 2022)	IN21k	Y	57.56	50.80	42.88	46.84
ViT RoPE (Wightman, 2019)	IN1k	Y	61.25	52.50	42.32	47.41
Eva02 (Fang et al., 2023b)	IN21K	N	<b>65.19</b>	<b>59.90</b>	<b>48.32</b>	<b>54.11</b>

Table 1: Effect of initialization of vision backbone on Imagenet Zero-shot and Flickr 30k Image to Text Recall@1, Text to Image Recall@1, and mean Recall@1. The pretrain column refers to the dataset the vision encoder was pretrained on and Supervised is whether the vision encoder used a supervised task to pretrain.

encoder. Our early work in adapting Three Towers style (Kossen et al., 2023) training resulted in poor general text embedding models, so we focused our effort on Locked Text Tuning.

#### 4 Image Text Datasets

Radford et al. (2021) describes curating a dataset of 400 million image-text pairs by searching for images that overlap with 500,000 popular phrases. This dataset was never released publicly.

Subsequent works by Schuhmann et al. (2021) and Schuhmann et al. (2022) openly released Laion 400M and Laion 5B to facilitate the training of open source multimodal models. Xu et al. (2024) aim to reproduce the data curated in Radford et al. (2021) and outperforms the proprietary dataset without any reliance on an external model.

Gadre et al. (2023) also released Datacomp 1B, a top performing dataset on the Datacomp X-Large benchmark. Fang et al. (2023a) improves upon the dataset released in (Gadre et al., 2023) by learning a data filtering network that can be used to curate high quality image-text datasets. In this work, we use Data Filtering Networks 2B (DFN-2B), the curated dataset for the Datacomp X-Large track. At the time of curation, we were only able to obtain 1.5B of the 2B links.

#### 5 Experiments

Nomic Embed Vision v1 and v1.5 were trained with identical hyperparameters and recipes except for the initialization of their text encoders. We train on DFN-2B for 3 epochs with a batch size of 65,536, resulting in training on 5B samples. We initialize the text encoders for Nomic Embed Vision v1 and v1.5 from Nomic Embed Text v1 and v1.5 respectively (Nussbaum et al., 2024), and the vision encoder as EVA02-ViT B/16 (Fang et al., 2023b). We use the AdamW opti-

mizer (Loshchilov and Hutter, 2019) and a peak learning rate of 1e-3, 2000 warmup steps, and cosine decay. As noted in Zhai et al. (2023), we set weight decay to 0 for the pretrained vision encoder. We employ multi-head attention pooling (Kossen et al., 2023; Beyer et al., 2022; Wightman, 2019). We train on 224x224 pixel images and use the same image preprocessing as (Radford et al., 2021). We additionally employ small augmentations using random crops (Ilharco et al., 2021) and do not clamp the learnable logit scale unlike Radford et al. (2021).

##### 5.1 Evaluation of Design Decisions

Due to the high compute and time cost to training the full model, we explored different design decisions at smaller scales. We present evidence in favor of some of our design decisions.

For our small scale experiments, we train for 1 epoch and perform small hyperparameter searches over learning rate and weight decay. We employ the Locked Text Tuning strategy outlined above and freeze Nomic Embed Text v1.

##### 5.2 Evaluating Batch Size

As noted in (Zhai et al., 2022; Chen et al., 2020; Radford et al., 2021), large batch sizes can improve the performance of contrastively trained models. We initialize the vision encoder with a ViT B/16 from (Dosovitskiy et al., 2021) and train on 300M image-text pairs over Data Filtering Networks from the Datacomp Large track (DFN-Large) (Fang et al., 2023a; Gadre et al., 2023).

As shown in Figure 2, increasing the batch size leads to sizable improvements on ImageNet 0-shot accuracy. We choose to use 65,536 as this is the biggest batch size we can accommodate given our compute limitation. We leave it to future work to investigate whether performance increases from increased batch size plateau.

### 5.3 Evaluating Pretrained Vision Encoders

To evaluate pretrained vision encoders, we train for one epoch on DFN-Large (Fang et al., 2023a). For each encoder, we perform a sweep over learning rate and weight decay. We investigate vision encoders released in Dosovitskiy et al. (2021), Steiner et al. (2022), Wightman (2019), and Fang et al. (2023b).

Similar to (Zhai et al., 2022), we found that the pretrained vision encoder backbone had a large effect on the quality of the final model, particularly in regards to multimodal retrieval. As shown in Table 1, we find that more broadly pretrained vision encoders lead to better multimodal retrieval and Imagenet zero-shot results. For example, a supervised vision encoder like the ViT B/16 released in (Dosovitskiy et al., 2021) performs well on Imagenet zero-shot but poorly on Flickr 30k retrieval (Young et al., 2014). Recently released ViTs from (Wightman, 2019) using techniques like global registers (Darcet et al., 2023) and rotary positional embeddings (Fang et al., 2023b) show promise even though they are trained on a small dataset like Imagenet.

From Table 1, we notice that even though models released by Wightman (2019) and Steiner et al. (2022) are trained in a supervised manner, they outperform the ViT B/16 released by Dosovitskiy et al. (2021). As noted by Zhai et al. (2022), training on large amounts of data leads to better and more general visual representations, even when training without a supervised objective. We hypothesize the high performance of the ViT B/16 released in Wightman (2019) is due to using heavy augmentation, like RandAugment (?), and training for many epochs.

Ultimately, the ViT B/16 released in (Fang et al., 2023b) performed the best across Imagenet zero-shot and Flickr retrieval, which leverages the unsupervised Masked Image Modeling (MIM) objective (Bao et al., 2022).

### 5.4 Evaluating Pooling Strategies

We evaluate different pooling layers for the vision encoder. We compare using the class token pooling, mean pooling, and multihead attention pooling (Kossen et al., 2023; Beyler et al., 2022). Again, our small scale experiments consist of training for 1 epoch over DFN Large (Fang et al., 2023a). We initialize the pretrained vision encoder from Fang et al. (2023b). We find that

Pooling Layer Comparison

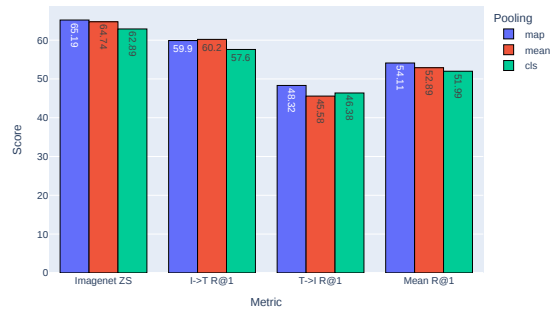


Figure 3: Effect of Pooling Layer on Performance in various retrieval and classification setups.

multihead attention pooling (MAP) performed the best over Imagenet Zero-Shot and Flickr 30k as shown in Figure 3.

## 6 Training Resources

Nomic Embed Vision v1 and v1.5 were trained on 2 8xH100s over 3.5 days. DFN-2B requires 62TB to store and several thousand dollars to preprocess.

## 7 Discussion

We present a recipe for enhancing a high quality text embedder with multimodal capabilities. While this model outperforms other unified embedding spaces, there are several important caveats. Consistent with prior CLIP literature, Nomic Embed exhibits bag of words like behavior on some tasks. (Yuksekgonul et al., 2023; Paiss et al., 2023). We also find that the retrieval scores resulting from Locked Text Tuning tend to skew low compared to similarly performing CLIP-style models as shown in Table 2. Three towers training (Kossen et al., 2023) presents a promising direction for remedying this.

Future work can investigate if similar strategies to those shown in Tschannen et al. (2023) can be adapted with a strong general purpose text encoder. However, some modifications may have to be made as the text encoder used in this work is bidirectional.

Moreover, recent work on multimodal embedding space geometry suggests that CLIP style training is not sufficient for closing the modality gap present in multimodal embedding spaces. Liang et al. (2022); Zhang et al. (2024b) As a result, we refer to the embedding spaces of Nomic Embed Text and Nomic Embed Image as unified

Model	ImageNet	ImageNet dist. shifts	VTAB	Retrieval	Average
Nomic Embed v1.5	0.710	0.551	0.561	0.469	0.568
Nomic Embed v1	0.707	0.551	0.565	0.457	0.567
CLIP ViT B-16	0.684	0.559	0.546	0.527	0.563
Jina CLIP v1	0.591	0.464	0.520	0.604	0.522

Table 2: Model Performance on DataComp Classification and Retrieval Tasks

and not aligned in this work. We leave it to future work to investigate whether closing the modality gap improves downstream performance.

## 8 Conclusion

We adapt the contrastive tuning framework presented in (Zhai et al., 2022) to enhance a high performing text encoder with multimodal capabilities. We call this training paradigm Locked Text Tuning, and use it to train Nomic Embed Vision. Together, Nomic Embed Vision and Nomic Embed Text form the first unified latent space to achieve high performance across vision, language, and multimodal tasks.

## References

- Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. 2022. *Beit: Bert pre-training of image transformers*.
- Lucas Beyer, Xiaohua Zhai, and Alexander Kolesnikov. 2022. *Big vision*. [https://github.com/google-research/big\\_vision](https://github.com/google-research/big_vision).
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. *A simple framework for contrastive learning of visual representations*.
- Timothée Darcet, Maxime Oquab, Julien Mairal, and Piotr Bojanowski. 2023. *Vision transformers need registers*. *arXiv preprint arXiv:2309.16588*.
- 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. 2021. *An image is worth 16x16 words: Transformers for image recognition at scale*.
- Alex Fang, Albin Madappally Jose, Amit Jain, Ludwig Schmidt, Alexander Toshev, and Vaishaal Shankar. 2023a. *Data filtering networks*.
- Yuxin Fang, Quan Sun, Xinggang Wang, Tiejun Huang, Xinlong Wang, and Yue Cao. 2023b. *Eva-02: A visual representation for neon genesis*.
- Samir Yitzhak Gadre, Gabriel Ilharco, Alex Fang, Jonathan Hayase, Georgios Smyrnis, Thao Nguyen, Ryan Marten, Mitchell Wortsman, Dhruva Ghosh, Jieyu Zhang, Eyal Orgad, Rahim Entezari, Giannis Daras, Sarah Pratt, Vivek Ramanujan, Yonatan Bitton, Kalyani Marathe, Stephen Mussmann, Richard Vencu, Mehdi Cherti, Ranjay Krishna, Pang Wei Koh, Olga Saukh, Alexander Ratner, Shuran Song, Hannaneh Hajishirzi, Ali Farhadi, Romain Beaumont, Sewoong Oh, Alex Dimakis, Jenia Jitsev, Yair Carmon, Vaishaal Shankar, and Ludwig Schmidt. 2023. *Datacomp: In search of the next generation of multimodal datasets*.
- Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan Misra. 2023. *Imagebind: One embedding space to bind them all*.
- Michael Günther, Jackmin Ong, Isabelle Mohr, Alaeddine Abdesslem, Tanguy Abel, Mohammad Kalim Akram, Susana Guzman, Georgios Mastrapas, Saba Sturua, Bo Wang, Maximilian Werk, Nan Wang, and Han Xiao. 2024. *Jina embeddings 2: 8192-token general-purpose text embeddings for long documents*.
- Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali Farhadi, and Ludwig Schmidt. 2021. *Openclip*. If you use this software, please cite it as below.
- Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc V. Le, Yunhsuan Sung, Zhen Li, and Tom Duerig. 2021. *Scaling up visual and vision-language representation learning with noisy text supervision*.
- Jannik Kossen, Mark Collier, Basil Mustafa, Xiao Wang, Xiaohua Zhai, Lucas Beyer, Andreas Steiner, Jesse Berent, Rodolphe Jenatton, and Efi Kokiopoulou. 2023. *Three towers: Flexible contrastive learning with pretrained image models*.
- Andreas Koukounas, Georgios Mastrapas, Michael Günther, Bo Wang, Scott Martens, Isabelle Mohr, Saba Sturua, Mohammad Kalim Akram, Joan Fontanals Martínez, Saahil Ognawala, Susana Guzman, Maximilian Werk, Nan Wang, and Han Xiao. 2024. *Jina clip: Your clip model is also your text retriever*.

- Zehan Li, Xin Zhang, Yanzhao Zhang, Dingkun Long, Pengjun Xie, and Meishan Zhang. 2023. [Towards general text embeddings with multi-stage contrastive learning](#).
- Weixin Liang, Yuhui Zhang, Yongchan Kwon, Serena Yeung, and James Zou. 2022. [Mind the gap: Understanding the modality gap in multi-modal contrastive representation learning](#).
- Ilya Loshchilov and Frank Hutter. 2019. [Decoupled weight decay regularization](#).
- Niklas Muennighoff, Nouamane Tazi, Loïc Magne, and Nils Reimers. 2023. [Mteb: Massive text embedding benchmark](#).
- Zach Nussbaum, John X. Morris, Brandon Duderstadt, and Andriy Mulyar. 2024. [Nomic embed: Training a reproducible long context text embedder](#).
- Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2019. [Representation learning with contrastive predictive coding](#).
- Roni Paiss, Ariel Ephrat, Omer Tov, Shiran Zada, Inbar Mosseri, Michal Irani, and Tali Dekel. 2023. [Teaching clip to count to ten](#).
- 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. 2021. [Learning transferable visual models from natural language supervision](#).
- Nils Reimers and Iryna Gurevych. 2019. [Sentencebert: Sentence embeddings using siamese bert-networks](#).
- 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. 2022. [Laion-5b: An open large-scale dataset for training next generation image-text models](#).
- Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. 2021. [Laion-400m: Open dataset of clip-filtered 400 million image-text pairs](#).
- Andreas Steiner, Alexander Kolesnikov, Xiaohua Zhai, Ross Wightman, Jakob Uszkoreit, and Lucas Beyer. 2022. [How to train your vit? data, augmentation, and regularization in vision transformers](#).
- Hongjin Su, Weijia Shi, Jungo Kasai, Yizhong Wang, Yushi Hu, Mari Ostendorf, Wen tau Yih, Noah A. Smith, Luke Zettlemoyer, and Tao Yu. 2023. [One embedder, any task: Instruction-finetuned text embeddings](#).
- Michael Tschannen, Manoj Kumar, Andreas Steiner, Xiaohua Zhai, Neil Houlsby, and Lucas Beyer. 2023. [Image captioners are scalable vision learners too](#).
- Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. 2024. [Text embeddings by weakly-supervised contrastive pre-training](#).
- Ross Wightman. 2019. [Pytorch image models](#). <https://github.com/huggingface/pytorch-image-models>.
- Shitao Xiao, Zheng Liu, Peitian Zhang, Niklas Muennighoff, Defu Lian, and Jian-Yun Nie. 2024. [C-pack: Packaged resources to advance general chinese embedding](#).
- Hu Xu, Saining Xie, Xiaoqing Ellen Tan, Po-Yao Huang, Russell Howes, Vasu Sharma, Shang-Wen Li, Gargi Ghosh, Luke Zettlemoyer, and Christoph Feichtenhofer. 2024. [Demystifying clip data](#).
- Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. 2014. [From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions](#). *Transactions of the Association for Computational Linguistics*, 2:67–78.
- Mert Yuksekogonul, Federico Bianchi, Pratyusha Kalluri, Dan Jurafsky, and James Zou. 2023. [When and why vision-language models behave like bags-of-words, and what to do about it?](#)
- Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. 2023. [Sigmoid loss for language image pre-training](#).
- Xiaohua Zhai, Xiao Wang, Basil Mustafa, Andreas Steiner, Daniel Keysers, Alexander Kolesnikov, and Lucas Beyer. 2022. [Lit: Zero-shot transfer with locked-image text tuning](#).
- Beichen Zhang, Pan Zhang, Xiaoyi Dong, Yuhang Zang, and Jiaqi Wang. 2024a. [Long-clip: Unlocking the long-text capability of clip](#).
- Yuhui Zhang, Elaine Sui, and Serena Yeung-Levy. 2024b. [Connect, collapse, corrupt: Learning cross-modal tasks with uni-modal data](#).

## Appendix

Table 3: Detailed performance on the CLIP Benchmark. Numbers for JinaCLIP (Koukounas et al., 2024), OpenAI CLIP (Radford et al., 2021), EVA02-CLIP (Fang et al., 2023b), and Long CLIP (Zhang et al., 2024a) reported from Koukounas et al. (2024)

Model	JinaCLIP	Nomic Embed	OpenAI CLIP	EVA02-CLIP	LongCLIP
<b>Zero-shot Image Retrieval - Recall@5 [%]</b>					
Average	80.31	69.43	75.62	<b>82.15</b>	81.72
Flickr30k	89.02	77.98	85.60	<b>91.10</b>	90.46
Flickr8k	85.50	74.10	82.84	<b>88.50</b>	88.40
MSCOCO	66.42	56.21	58.42	<b>66.85</b>	66.31
<b>Zero-shot Text Retrieval - Recall@5 [%]</b>					
Average	89.91	80.44	88.12	90.59	<b>90.79</b>
Flickr30k	96.50	89.89	96.20	96.60	<b>98.00</b>
Flickr8k	94.20	84.50	91.40	<b>94.60</b>	94.00
MSCOCO	79.02	67.02	76.76	<b>80.58</b>	80.38
<b>Image Classification - Accuracy@1 [%]</b>					
Average	43.28	46.62	46.16	<b>48.70</b>	46.67
Cars	68.03	<b>87.60</b>	64.73	78.56	59.17
Country211	13.45	16.35	<b>22.85</b>	21.34	20.28
Fer2013	<b>49.07</b>	20.30	46.18	51.17	47.80
Fgvc-aircraft	11.49	23.64	24.27	<b>25.11</b>	22.56
Gtsrb	38.70	45.22	43.58	<b>46.33</b>	42.93
Imagenet-a	29.92	46.04	49.93	<b>53.89</b>	46.84
Imagenet-o	33.40	20.55	42.25	34.10	<b>42.65</b>
Imagenet-r	73.66	<b>82.46</b>	77.69	82.42	76.63
Imagenet1k	59.08	71.03	68.32	<b>74.75</b>	66.84
Imagenet-sketch	45.04	57.51	48.25	<b>57.70</b>	47.12
Imagenetv2	51.37	62.17	61.95	<b>66.98</b>	60.17
Mnist	48.07	59.42	65.51	47.16	<b>71.84</b>
Objectnet	45.41	62.02	55.35	<b>62.29</b>	50.79
Renderedsst2	59.14	55.29	<b>60.68</b>	54.15	59.31
Stl10	97.89	97.47	98.28	<b>99.49</b>	98.41
Sun397	65.92	65.12	64.37	<b>70.62</b>	68.73
Voc2007	72.83	61.75	78.34	<b>80.17</b>	75.35
Vtab/caltech101	82.68	<b>84.58</b>	82.19	82.78	82.63
Vtab/cifar10	93.49	96.82	90.78	<b>98.46</b>	91.22
Vtab/cifar100	72.08	83.62	66.94	<b>87.72</b>	69.17
Vtab/clevr-closest-object-distance	15.61	15.84	15.83	15.72	<b>15.90</b>
Vtab/clevr-count-all	<b>22.35</b>	21.62	21.09	21.27	20.71
Vtab/diabetic-retinopathy	2.82	4.51	3.44	14.19	10.99
Vtab/dmlab	<b>19.53</b>	13.97	15.49	14.67	15.45
Vtab/dsprites-label-orientation	2.44	1.63	2.34	1.94	1.12
Vtab/dsprites-label-x-position	3.07	2.95	2.95	3.11	<b>3.15</b>
Vtab/dsprites-label-y-position	3.17	2.87	3.11	<b>3.21</b>	3.16
Vtab/dtd	<b>55.43</b>	50.27	44.89	52.82	45.27
Vtab/eurosat	49.52	37.27	55.93	<b>66.33</b>	60.44
Vtab/flowers	59.62	68.23	71.13	<b>75.75</b>	69.85
Vtab/kitti-closest-vehicle-distance	22.93	<b>38.82</b>	26.44	22.08	34.60
Vtab/pcam	55.54	61.48	50.72	50.95	52.55
Vtab/pets	80.98	91.79	89.04	<b>92.10</b>	89.21
Vtab/resisc45	55.46	57.12	58.27	60.37	<b>60.63</b>
Vtab/smallnorb-label-azimuth	<b>5.40</b>	5.30	5.21	4.96	5.14
Vtab/smallnorb-label-elevation	11.31	9.62	<b>12.17</b>	9.79	10.59
Vtab/svhn	25.46	<b>42.70</b>	31.20	17.65	27.65