

CLIPA-v2: Scaling CLIP Training with 81.1% Zero-shot ImageNet Accuracy within a \$10,000 Budget

Xianhang Li* Zeyu Wang* Cihang Xie

*equal technical contribution

UC Santa Cruz

<https://github.com/UCSC-VLAA/CLIPA>

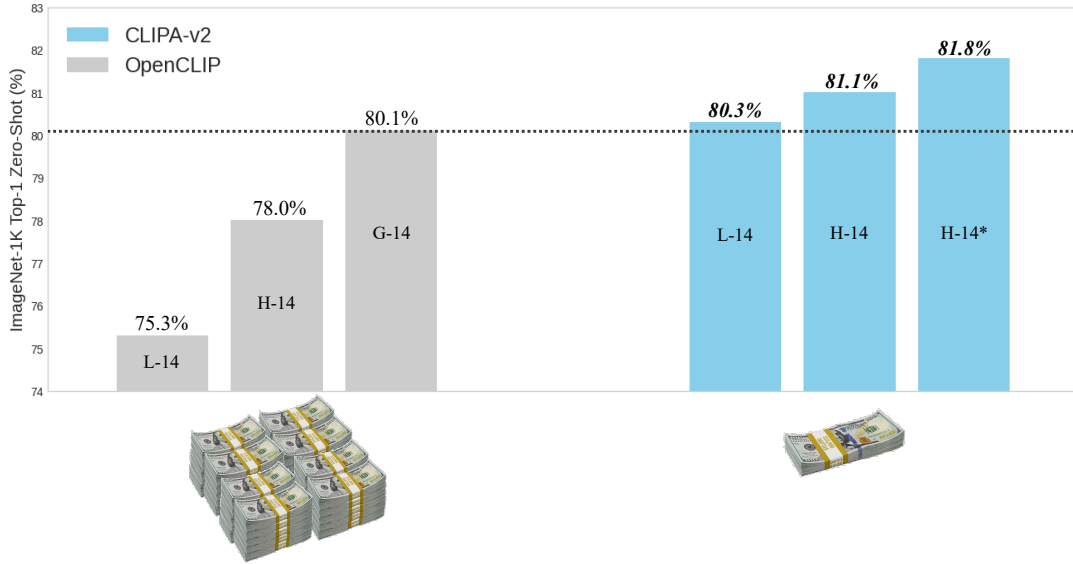


Figure 1: Compared to OpenCLIP [7], our CLIPA-v2 models achieve higher performance with lower training cost.

Abstract

The recent work CLIPA [9] presents an inverse scaling law for CLIP training — whereby the larger the image/text encoders used, the shorter the sequence length of image/text tokens that can be applied in training. This finding enables us to train high-performance CLIP models with significantly reduced computations. Building upon this work, we hereby present CLIPA-v2 with two key contributions. Technically, we find this inverse scaling law is also applicable in the finetuning stage, enabling further reduction in computational needs. Empirically, we explore CLIPA at scale, extending the experiments up to the H/14 model with $\sim 13B$ image-text pairs seen during training.

Our results are exciting — by only allocating a budget of \$10,000, our CLIP model achieves an impressive zero-shot ImageNet accuracy of 81.1%, surpassing the prior best CLIP model (from OpenCLIP, 80.1%) by 1.0% and meanwhile reducing the computational cost by $\sim 39\times$. Moreover, with an additional investment of \$4,000, we can further elevate the zero-shot ImageNet accuracy to 81.8%.

1. Introduction

CLIP [12] has emerged as the pioneering foundation model that bridges the gap between text and images, ushering computer vision research into the “post-ImageNet” era [7, 10, 20, 1, 13, 14, 16, 18, 3]. However, the demanding computational requirements of CLIP hinder its widespread exploration. The recent work CLIPA [9] offers a computationally efficient solution — with the introduction of an inverse scaling law for CLIP training, it reveals that larger models can be trained with fewer input tokens. Building upon this observation, CLIPA demonstrates its efficacy in scenarios with limited computational resources, resulting in a substantial reduction in the training cost of CLIP.

This report provides a follow-up on CLIPA. Firstly, we validate that the inverse scaling law is also applicable when finetuning models with input tokens at full resolution. This further reduces the training cost of CLIPA. Secondly, we investigate the performance of CLIPA at scale across various aspects, including model size (up to H/14), data (up to DataComp-1B [5] and LAION-2B [16] datasets), and training schedule (up to $\sim 13B$ samples seen).

model	# image token	# text token	data source	# seen samples	total compute ($\times 1e11$)	IN-1K
CLIPA-L/16	36	8	LAION-400M	2.56B + 128M	0.5	69.3
			LAION-400M	2.56B + 128M	0.8	72.8
CLIPA H/14	36	8	LAION-2B	2.56B + 128M	0.8	74.1
			LAION-2B	12.8B + 128M	4	77.9

Table 1: **Scaling up CLIPA-v1 [9]**. We employ a joint scaling strategy along the data, model, and schedule axes. We pretrain the H/14 model with 36 image tokens (84×84) and 8 text tokens. For finetuning, we use 256 (224×224) image tokens and 32 text tokens, following [9].

With these two contributions, we can train CLIP models with strong zero-shot performance on ImageNet [4], meanwhile significantly reducing training costs. For instance, we can train a H/14 model with 81.1% accuracy within a \$10,000 budget. We stress that, compared to the best publicly available CLIP model from OpenCLIP [7], ours is both better (**+1.0%**) and faster (by $\sim 39\times$). Moreover, we can further boost this accuracy to 81.8%, with an additional \$4,000 investment. These results are exciting as no prior work has thus far reached a similar performance within this small budget limitation. By open-sourcing our training code and models, we hope to contribute to the broader advancement and adoption of advanced CLIP models.

2. Background

CLIP has been a prominent foundation model due to its exceptional zero-shot capability and remarkable versatility [12, 8]. The tremendous success of CLIP can be attributed to the extensive scale of both the data [12, 15, 8, 2, 20, 21] and the model [19, 11, 17] it is built upon. Nevertheless, it also poses a significant cost barrier to researchers who wish to train a strong CLIP model. To reduce the computational burden, the recent work by Li et al. [9] presents an inverse scaling law, which reveals that larger models can effectively utilize fewer tokens for training without severe performance drop, therefore enabling highly efficient CLIP training. As a byproduct of this discovery, the CLIPA models are introduced, which attains a zero-shot top-1 ImageNet accuracy of 69.3% and can be trained on an 8 A100-GPU machine in just 4 days.

Our work is built upon CLIPA [9], but focuses on furthering its efficiency and scaling it up.

3. Experiments

Our experiments contain three parts. Firstly, we check the applicability of *inverse scaling law* during the finetuning stage with full-resolution tokens. Next, we scale up CLIPA in terms of data, model, and schedule. Lastly, we compare with other advanced CLIP models in terms of performance and computation cost.

Settings. Our pretraining setup strictly follows CLIPA [9]. We report the corresponding zero-shot top-1 accuracy on ImageNet [4].

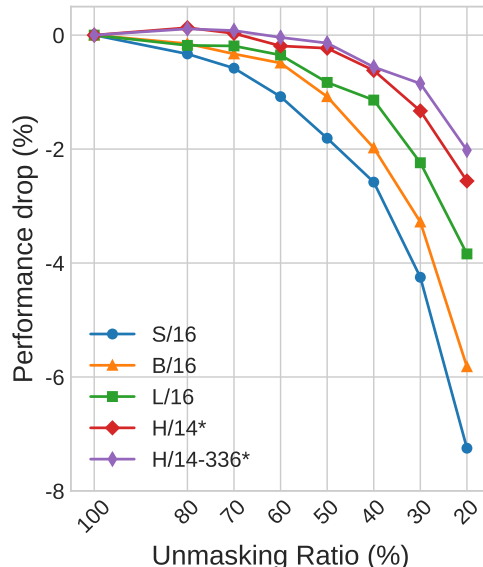


Figure 2: **The inverse scaling law on finetuning.** All models are finetuned with 128M samples, where we employ random masking for token reduction. *: H/14 model are pre-trained on LAION-2B, unlike others on LAION-400M.

Inverse scaling law in the finetuning stage. We first validate the inverse scaling law in the finetuning stage. Following [9], we choose four different scales of models: S/16, B/16, L/16, and H/14. Random masking [10, 6] is used as the token reduction strategy. As shown in Figure 2, we can observe that larger models can be finetuned with fewer tokens while still keeping the performance drop small. For instance, retaining 50% of the tokens merely results in a performance drop of 0.2% for the H/14 model, compared to larger drops of 0.8% for L/16, 1.1% for B/16, and 1.8% for S/16. More interestingly, we observed that when using larger image inputs, such as 336×336 , we could apply even higher masking ratios, as demonstrated in the comparison between H/14 and H/14-336.

These results confirm the existence of the inverse scaling law in the finetuning stage, which enables us to reduce the required computations for CLIP training further.

Scaling up CLIPA [9]. We next investigate the scaling behavior beyond the largest case studied in CLIPA. Specifically, our scaling efforts cover three aspects: model, data, and training schedule. The results are reported in Table 1.

case	mask ratio	resolution	# seen samples	training FLOPs	IN-1K
CLIPA-v1	0%	224 ²	128M	177.0G	77.9
(1)	30%	224 ²	128M	135.9G	78.0
(2)	30%	224 ²	512M	135.9G	78.6
(3)	30%	224 ²	640M	135.9G	78.5
(4)	40%	336 ²	640M	237.8G	78.9
(5)	30%+40%	224 ² + 336 ²	512M+128M	156.3G	79.1

Table 2: **Ablation of CLIPA-v2.** We report zero-shot top-1 ImageNet accuracy. In case (5), we use 224×224 input with a mask ratio of 30% for the first 512M samples, and 336×336 input with a mask ratio of 40% for the rest 128M samples.

	model	data source	# seen samples@input size	GPU hours ¹	Est. cost ²	IN-1K
OpenCLIP	H/14	LAION-2B	32.0B@224 ²	216,712	\$190,706	78.0
CLIPA-v2			12.8B@84 ² + 512M@224 ² + 128M@336 ²	8,640	\$13,613	79.1
OpenCLIP	L/14	DataComp-1B	12.8B@224 ²	41,472	\$65,111	79.2
	G/14*	LAION-2B	32.0B@224 ² + 6.7B@224 ²	232,448	\$366,105	80.1
CLIPA-v2	H/14	DataComp-1B	12.8B@70 ² + 512M@224 ²	5,920	\$9,324	81.1
CLIPA-v2	L/14	DataComp-1B	12.8B@84 ² + 512M@224 ²	4,008	\$6,318	79.7
			+128M@336 ²	+512	+\$806	80.3
	H/14		12.8B@84 ² + 512M@224 ²	7,776	\$12,247	81.5
			+128M@336 ²	+864	+\$1,366	81.8

Table 3: **Comparison with OpenCLIP [7].** Our CLIPA-v2’s GPU hour is estimated using an 8-A100 80GB GPU machine on Google Cloud, while the OpenCLIP’s GPU hour is calculated based on their report¹. The corresponding training cost is estimated based on 80GB A100’s cloud pricing². * denotes this model is trained with FLIP at a masking ratio of 50%.

First, we can observe that scaling the model size from L/14 to H/14 boosts the performance from 69.3% to 72.8%. Furthermore, we note switching the training dataset from LAION-400M [16] to LAION-2B [15] yields another 1.3% improvement, suggesting the importance of data diversity. Lastly, by increasing the training schedule by a factor of 5, resulting in a total of ~ 13 B seen samples, we achieve an impressive performance of 77.9%. We stress that this scaled version of CLIPA H/14 model readily outperforms its counterpart in FLIP [10] by 0.3% while requiring only 1/3 of the training budget.

These results confirm the efficiency and effectiveness of training CLIPA at scale. Next, we set this CLIPA H/14 as our baseline for further ablation in the finetuning stage.

Ablation. The results of our ablation studies on different finetuning setups are summarized in Table 2. Interestingly, compared to finetuning input tokens at the full resolution, we observe that 30% masked finetuning even leads to a slight performance improvement (+0.1%). Additionally, this masking strategy enables a $1.7\times$ speedup of the finetuning process. Furthermore, adopting a $4\times$ finetuning schedule results in an additional improvement of 0.6%. However, further increasing the finetuning schedule does not lead to any substantial performance gains.

Following [7], we also investigate progressively finetuning with large resolutions. Initially, for the first 512 million samples, we finetune the model using a 224×224 input size with a masking ratio of 30%; subsequently, for the remaining 128 million samples, we adopt a larger 336×336 input size with a masking ratio of 40% and a smaller learning rate.

As shown in Table 2, progressive finetuning results in a performance improvement of 0.2% compared to direct finetuning with a 336×336 input size and meanwhile achieving a notable $1.5\times$ speedup of the finetuning process.

Comparison with OpenCLIP [7]. We summarize the results in Table 3. Firstly, when trained on the LAION-2B dataset, our CLIPA-v2 H/14 model outperforms OpenCLIP’s version by 1.1% (79.1% vs. 78.0%) and meanwhile significantly reducing the training cost by $\sim 14\times$. Furthermore, when upgrading to the DataComp-1B dataset, our CLIPA-v2 H/14 (pretrained on images at 70×70) achieves an impressive zero-shot ImageNet accuracy of **81.1%**, while keeping the training cost within \$10,000. Notably, this 81.1% accuracy is 1.0% higher than the prior best CLIP model, which is OpenCLIP’s G/14 model with a zero-shot ImageNet accuracy of 80.1%.

With an additional investment of \$4000, we can further enhance CLIPA-v2’s training by 1) pretraining with a larger resolution (from 70 to 84) and 2) applying the progressive finetuning with a larger image resolution of 336. This leads to an additional 0.7% improvement, resulting in the *best-performing CLIP model to date with an 81.8% zero-shot ImageNet accuracy*. We have open-sourced these CLIP models to facilitate future research.

¹We measure OpenCLIP [7]’s training time based on <https://laion.ai/blog/large-openclip/> and <https://laion.ai/blog/giant-openclip/>.

²We estimate the total training cost based on <https://cloud.google.com/compute/gpus-pricing>, which is \$1.575 per GPU hour, and <https://lambdalabs.com/service/gpu-cloud/pricing>, which is \$1.5 per GPU hour.

References

- [1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. In *NeurIPS*, 2022.
- [2] 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.
- [3] Yuchen Cui, Scott Niekum, Abhinav Gupta, Vikash Kumar, and Aravind Rajeswaran. Can foundation models perform zero-shot task specification for robot manipulation? *arXiv preprint arXiv:2204.11134*, 2022.
- [4] 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.
- [5] Samir Yitzhak Gadre, Gabriel Ilharco, Alex Fang, Jonathan Hayase, Georgios Smyrnis, Thao Nguyen, Ryan Marten, Mitchell Wortsman, Dhruva Ghosh, Jieyu Zhang, et al. Datacomp: In search of the next generation of multimodal datasets. *arXiv preprint arXiv:2304.14108*, 2023.
- [6] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *CVPR*, 2022.
- [7] Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hananeh Hajishirzi, Ali Farhadi, and Ludwig Schmidt. Open-clip, July 2021.
- [8] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *ICML*, 2021.
- [9] Xianhang Li, Zeyu Wang, and Cihang Xie. An inverse scaling law for clip training. *arXiv preprint arXiv:2305.07017*, 2023.
- [10] Yanghao Li, Haoqi Fan, Ronghang Hu, Christoph Feichtenhofer, and Kaiming He. Scaling language-image pre-training via masking. In *CVPR*, 2023.
- [11] OpenAI. Gpt-4 technical report. 2023.
- [12] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, 2021.
- [13] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *ICML*, 2021.
- [14] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022.
- [15] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. In *NeurIPS*, 2022.
- [16] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114*, 2021.
- [17] Quan Sun, Yuxin Fang, Ledell Wu, Xinlong Wang, and Yue Cao. Eva-clip: Improved training techniques for clip at scale. *arXiv preprint arXiv:2303.15389*, 2023.
- [18] Hu Xu, Saining Xie, Po-Yao Huang, Licheng Yu, Russell Howes, Gargi Ghosh, Luke Zettlemoyer, and Christoph Feichtenhofer. Cit: Curation in training for effective vision-language data. *arXiv preprint arXiv:2301.02241*, 2023.
- [19] Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu. Coca: Contrastive captioners are image-text foundation models. *arXiv preprint arXiv:2205.01917*, 2022.
- [20] Lu Yuan, Dongdong Chen, Yi-Ling Chen, Noel Codella, Xiyang Dai, Jianfeng Gao, Houdong Hu, Xuedong Huang, Boxin Li, Chunyuan Li, et al. Florence: A new foundation model for computer vision. *arXiv preprint arXiv:2111.11432*, 2021.
- [21] Wanrong Zhu, Jack Hessel, Anas Awadalla, Samir Yitzhak Gadre, Jesse Dodge, Alex Fang, Youngjae Yu, Ludwig Schmidt, William Yang Wang, and Yejin Choi. Multimodal c4: An open, billion-scale corpus of images interleaved with text. *arXiv preprint arXiv:2304.06939*, 2023.