

# HQ-CLIP: Leveraging Large Vision-Language Models to Create High-Quality Image-Text Datasets and CLIP Models

## Supplementary Material

### 1. Experiments

DataComp Scale	Small	Medium	Large
CommonPool size	12.8M	128M	1.28B
Original DFN size	-	19.2M	192M
Reproduced DFN size	1.47M	14.7M	147M
Model	ViT-B/32	ViT-B/32	ViT-B/16
Batch size	4096	4096	8192

Table 1. Training setup and dataset scale.

#### 1.1. Setup

Our experimental setup primarily follows the configuration established in DFN [1]. The original DFN methodology processes CommonPool datasets (12.8M/128M/1.28B) to derive filtered subsets of 1.92M (small), 19.2M (medium), and 192M (large) image-text pairs. Due to partial URL inaccessibility, we obtained reduced subsets of **1.47M (small)**, **14.7M (medium)**, and **147M (large)** pairs for our implementation. In model training, we strictly adhere to DFN’s architectural specifications and batch size configurations. Notably, for the XLarge-scale model training, we employed CLIPA [4] to optimize computational efficiency and accelerate training convergence.

#### 1.2. Ablation study

**Ablation study on hard-negative sample quantity.** We investigate the optimal number of hard-negative variants per image for identification tasks. As Table 8 demonstrates, empirical evidence suggests the single-sample configuration emerges as optimal. Although increasing the number of samples initially appears to benefit performance metrics, practical constraints such as prohibitive GPU memory demands and computational overhead prevent further scaling. Consequently, we select one hard-negative instance as the computationally efficient yet effective solution.

**Ablation on the number of classes.** Our framework employs a frequency-based selection of the top-K most prevalent tags from the VLM-generated tag repository. As empirically validated in Table 3, we systematically determine the optimal class quantity parameter  $K$ .

**Ablation study on loss hyperparameters  $\alpha$  and  $\beta$ .** Performance sensitivity to the hard-negative identification loss weight ( $\alpha$ ) and short tag classification loss weight ( $\beta$ ) is quantified in Tables 5 and 4. The optimal configuration is observed at  $\alpha = 0.5$  and  $\beta = 10$ , where both loss components contribute maximally to model effectiveness.

Scale	Methods	Attribution	Relation
Medium	DFN <sup>†</sup>	54.2	53.2
	Ours	<b>61.1</b>	<b>54.4</b>
Large	DFN <sup>†</sup>	55.1	47.2
	Ours	<b>65.1</b>	<b>61.3</b>

Table 2. Comparison of attribution and relation metrics in the ARO benchmark [5].

Number of classes	3000	10000	30000	90000
ImageNet	39.2	<b>40.8</b>	40.5	40.6
ImageNet-Shifts	31.1	<b>32.9</b>	32.8	32.8
VTAB	40.3	40.5	40.5	<b>42.3</b>
Retrieval	36.6	37.4	<b>37.7</b>	37.3
Average over 38 datasets	39.9	40.2	40.1	<b>40.5</b>

Table 3. Ablation study on the number of classes.

$\beta$	1	10	100	1000
ImageNet	40.1	40.0	<b>40.7</b>	40.6
ImageNet-Shifts	32.4	32.4	<b>33.2</b>	32.1
VTAB	44.1	44.3	<b>44.7</b>	44.1
Retrieval	37.2	36.9	<b>37.5</b>	36.8
Average over 38 datasets	40.1	39.9	<b>40.5</b>	40.2

Table 4. Ablation study on the weight of  $\mathcal{L}_{STC}$ .

$\alpha$	0.1	0.2	0.5	1
ImageNet	40.6	<b>40.7</b>	40.1	40.2
ImageNet-Shifts	<b>32.7</b>	32.5	32.5	32.2
VTAB	40.8	41.2	<b>41.6</b>	41.3
Retrieval	38.7	37.7	<b>38.1</b>	<b>38.1</b>
Average over 38 datasets	40.1	39.9	<b>40.7</b>	40.0

Table 5. Ablation study on the weight of  $\mathcal{L}_{HNI}$ .

#### 1.3. Comparison with state-of-the-art method

**ARO benchmark evaluation.** As shown on Tab. 2, our approach exhibits superior comprehension of attribution and relation compared to the DFN<sup>†</sup> baseline. By benefiting from descriptions with enhanced semantic richness and the specialized hard-negative identification loss during training, our method achieves significant and scalable performance improvements on Visual Genome attribution metrics.

**Comparison with VeCLIP.** Given the exceptional performance claims of VeCLIP[3] in its original publication, com-

	Dataset size	IN	INv2	COCO	Flickr	Caltech101	CIFAR100	SVHN	DTD	OxPet	Flowers102	EuroSAT	RESISC45	Camelyon	Average
VeCLIP [3]	200M	64.6	57.7	<b>57.8</b>	<b>83.7</b>	83.1	68.1	44.9	<b>62.0</b>	72.6	68.5	47.4	55.1	<b>62.6</b>	62.7
Ours	148M	<b>70.6</b>	<b>63.1</b>	52.2	77.9	<b>93.1</b>	<b>81.0</b>	<b>45.9</b>	51.5	<b>89.5</b>	<b>69.0</b>	<b>47.6</b>	<b>60.6</b>	46.2	<b>64.9</b>

Table 6. Comparison of Our Method with VeCLIP. The metrics for VeCLIP are sourced from the original paper. Our method demonstrates superior average performance.

	IN	IN-Shifts	VTAB	Retrieval	Average over 38 datasets
VeCLIP*	52.5	45.9	46.8	55.2	48.2
Ours	<b>70.6</b>	<b>57.2</b>	<b>57.6</b>	<b>60.9</b>	<b>58.6</b>

Table 7. Comparison of Performance on the DataComp [2] Benchmark with VeCLIP. The metrics for VeCLIP were obtained by using the weights provided in its official GitHub repository, trained on the 100 Million dataset, and evaluated using the DataComp benchmark code and Hugging Face tools.

$N^-$	1	2	3	4
ImageNet	39.9	39.8	<b>40.0</b>	39.5
ImageNet-Shifts	<b>32.5</b>	<b>32.5</b>	32.1	32.3
VTAB	41.8	<b>42.0</b>	40.9	41.1
Retrieval	<b>38.1</b>	38.0	37.7	37.9
Average over 38 datasets	40.1	<b>40.2</b>	40.1	39.9

Table 8. Ablation study on number of hard negative samples  $M$ .

prehensive benchmarking becomes imperative. However, since VeCLIP did not include DataComp benchmark results in their work, a direct comparison in our main results table (Table 2) proves infeasible. We therefore provide supplementary comparisons with more performance metrics in the Supplementary Materials between our method and the ViT-B variant of VeCLIP trained on 200 million samples (as reported in their paper), where our approach demonstrates superior comprehensive performance (Table 6).

To facilitate rigorous benchmarking, we sought to evaluate VeCLIP under the DataComp [2] framework. While the authors provide clear instructions for loading their ViT-H weights, documentation gaps were identified regarding ViT-B weight implementation. Technical challenges emerged from (1) framework-specific implementation details in TensorFlow and (2) compatibility constraints with VeCLIP’s text encoder architecture in the DataComp library. To address these methodological challenges, we re-implemented a PyTorch version of VeCLIP’s data pipeline and modified the DataComp evaluation code.

Due to technical limitations in loading VeCLIP model weights trained on the 200M subset, our analysis employs the 100M variant for standardized DataComp benchmark comparisons (Table 7). HQ-CLIP significantly outperforms VeCLIP. We are actively seeking verification through direct communication with the authors’ team to ensure correct



Figure 1. Comparison of recognition results between our model and DFN.

comparison and sincerely welcome their insights.

## 1.4. Recognition Results

Figure 1 shows the classification results of our model compared to the DFN model. For each image, binary classification is performed using manually crafted text to demonstrate the fine-grained understanding capability of the models. Our model shows better recognition of detailed semantics in the images.

## 1.5. Details of other experiments

We showcase the full 38 dataset result for some experiments on main paper, as shown in Tab. ?? and 10.

## 2. VLM-150M

### 2.1. Examples

We present some examples from the acquired dataset. As shown in Figure 2, we obtained more comprehensive annotations.



Figure 2. Examples of VLM-150M.

Model	XCom2	LLaVA	Qwen2-VL	Qwen2-VL	Qwen2-VL	Qwen2-VL
Parameters	7B	7B	7B	2B	72B	7B
GPT4o SFT	✓	✓	✓	✓		✓
Caption Input	✓	✓		✓	✓	✓
ImageNet 1k	41.1	39.9	37.6	40.8	<b>41.2</b>	40.2
ImageNet Sketch	30.9	31.1	26.9	31.9	<b>31.9</b>	31.7
ImageNet V2	<b>34.1</b>	33.3	30.6	33.8	34.1	33.4
ImageNet-A	7.1	<b>7.5</b>	6.2	6.8	7.2	7.2
ImageNet-O	<b>48.9</b>	47.8	46.0	<b>48.9</b>	48.1	48.0
ImageNet-R	<b>47.6</b>	47.5	42.5	47.4	47.5	47.5
Caltech-101	81.7	80.4	78.9	80.8	80.7	<b>83.8</b>
CIFAR-10	89.8	88.2	83.8	88.1	88.4	<b>89.8</b>
CIFAR-100	63.8	63.6	59.2	65.2	65.0	<b>65.5</b>
CLEVR Counts	14.9	<b>26.2</b>	13.1	24.3	17.1	25.0
CLEVR Distance	<b>21.2</b>	18.6	16.4	15.9	15.9	15.8
SVHN	<b>26.8</b>	10.6	20.4	21.9	9.8	23.1
DTD	28.0	26.1	22.0	27.7	27.8	<b>28.7</b>
EuroSAT	35.9	<b>40.9</b>	22.5	31.4	36.5	32.6
KITTI distance	20.5	28.7	16.7	27.1	<b>34.2</b>	32.1
Oxford Flowers-102	38.8	35.8	39.3	39.3	<b>39.7</b>	36.3
Oxford-IIIT Pet	59.5	60.0	57.0	58.8	<b>61.3</b>	55.4
PatchCamelyon	57.5	54.7	56.8	52.3	<b>58.7</b>	53.1
RESISC45	31.0	34.8	28.7	<b>36.7</b>	33.9	34.5
FGVC Aircraft	3.3	3.2	3.3	2.6	<b>3.5</b>	3.4
Food-101	56.1	54.5	52.8	<b>56.5</b>	55.4	56.1
GTSRB	15.5	18.6	13.9	17.1	17.1	<b>19.7</b>
MNIST	29.8	22.8	23.4	29.5	26.1	<b>31.8</b>
ObjectNet	28.5	<b>28.7</b>	24.3	28.6	28.0	28.4
Pascal VOC 2007	63.8	70.2	54.7	67.6	69.1	<b>71.0</b>
Rendered SST2	50.2	50.1	<b>50.4</b>	49.9	49.2	49.7
Stanford Cars	45.3	44.1	<b>48.9</b>	45.7	48.2	42.5
STL-10	89.9	89.9	87.1	89.8	90.0	<b>90.2</b>
SUN-397	48.7	47.4	44.5	48.8	48.7	<b>49.7</b>
Country211	5.0	4.8	4.5	5.3	<b>5.3</b>	5.3
iWildCam	2.9	2.2	2.3	2.5	<b>3.5</b>	2.6
Camelyon17	57.0	65.8	<b>67.8</b>	53.1	66.0	55.8
FMoW	<b>0.0</b>	<b>0.0</b>	<b>0.0</b>	<b>0.0</b>	<b>0.0</b>	<b>0.0</b>
Dollar Street	<b>49.3</b>	46.1	47.1	48.2	48.7	47.4
GeoDE	73.0	70.5	66.2	<b>74.0</b>	<b>74.0</b>	68.8
Flickr30k	40.8	42.3	29.5	42.0	39.6	<b>44.6</b>
MSCOCO	25.3	18.0	17.2	26.2	24.2	<b>26.7</b>
WinoGAViL	43.1	37.7	36.9	41.6	<b>46.4</b>	40.5
Avg. over 38 datasets	39.6	39.3	36.3	39.7	<b>40.1</b>	39.9

Table 9. Comparison of the performance of different data refinement pipelines. Compared to other LVLMS, Qwen2VL demonstrates superior performance. Despite a tenfold difference in parameter size, Qwen2VL-7B with GPT-4o SFT still exhibits performance comparable to the 72B model. Additionally, the inclusion of captions significantly enhances dataset quality.

Method	Ours	DFN
DataComp scale	Large	Large
Dataset size	146.6M	146.6M
ImageNet 1k	<b>70.6</b>	68.7
ImageNet Sketch	<b>57.3</b>	54.9
ImageNet V2	<b>63.1</b>	60.0
ImageNet-A	<b>39.1</b>	29.9
ImageNet-O	43.0	<b>53.5</b>
ImageNet-R	<b>80.1</b>	75.4
Caltech-101	<b>93.1</b>	91.2
CIFAR-10	<b>96.2</b>	94.8
CIFAR-100	<b>81.0</b>	79.1
CLEVR Counts	<b>27.5</b>	14.7
CLEVR Distance	<b>22.2</b>	20.0
SVHN	45.9	<b>48.5</b>
DTD	<b>51.5</b>	46.9
EuroSAT	47.6	<b>49.9</b>
KITTI distance	<b>43.0</b>	24.9
Oxford Flowers-102	69.0	<b>71.0</b>
Oxford-IIIT Pet	<b>89.5</b>	88.7
PatchCamelyon	47.5	<b>51.0</b>
RESISC45	<b>60.6</b>	56.0
FGVC Aircraft	11.3	<b>13.2</b>
Food-101	<b>87.8</b>	86.2
GTSRB	<b>54.4</b>	44.2
MNIST	<b>77.7</b>	61.5
ObjectNet	<b>60.6</b>	55.0
Pascal VOC 2007	<b>78.8</b>	75.0
Rendered SST2	<b>51.7</b>	51.2
Stanford Cars	<b>85.3</b>	85.1
STL-10	<b>98.1</b>	96.0
SUN-397	<b>69.7</b>	67.2
Country211	<b>15.9</b>	13.5
iWildCam	<b>12.2</b>	10.0
Camelyon17	46.2	<b>63.1</b>
FMoW	<b>15.1</b>	10.9
Dollar Street	<b>61.3</b>	60.3
GeoDE	<b>88.7</b>	87.3
Flickr30k	<b>77.9</b>	68.2
MSCOCO	<b>52.2</b>	43.7
WinoGAViL	<b>52.8</b>	51.8
Avg. over 38 datasets	<b>58.6</b>	55.9

Table 10. Training on VLM-150M yields state-of-the-art CLIP models. We evaluate these models using the DataComp evaluation protocol. For detailed comparisons on specific datasets, we also provide the reproduced results for DFN. The symbol † indicates the results that we reproduced. Due to some broken links in the dataset, the amount of data used in our reproduction is slightly lower than that in the original paper.

## References

- [1] Alex Fang, Albin Madappally Jose, Amit Jain, Ludwig Schmidt, Alexander Toshev, and Vaishaal Shankar. Data filtering networks, 2023. [1](#)
- [2] 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. Datacomp: In search of the next generation of multimodal datasets, 2023. [2](#)
- [3] Zhengfeng Lai, Haotian Zhang, Bowen Zhang, Wentao Wu, Haoping Bai, Aleksei Timofeev, Xianzhi Du, Zhe Gan, Jiulong Shan, Chen-Nee Chuah, Yinfei Yang, and Meng Cao. Veclip: Improving clip training via visual-enriched captions, 2024. [1](#), [2](#)
- [4] Xianhang Li, Zeyu Wang, and Cihang Xie. Clipa-v2: Scaling clip training with 81.1 [1](#)
- [5] Mert Yuksekgonul, Federico Bianchi, Pratyusha Kalluri, Dan Jurafsky, and James Zou. When and why vision-language models behave like bags-of-words, and what to do about it? In *International Conference on Learning Representations*, 2023.