

# IN SEARCH OF FORGOTTEN DOMAIN GENERALIZATION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Out-of-Domain (OOD) generalization is the ability of a model trained on one or more domains to generalize to unseen domains. In the ImageNet era of computer vision, evaluation sets for measuring a model’s OOD performance were designed to be strictly OOD with respect to style. However, the emergence of foundation models and expansive web-scale datasets has obfuscated this evaluation process, as datasets cover a broad range of domains and risk test domain contamination. In search of the forgotten domain generalization, we create large-scale datasets subsampled from LAION—LAION-Natural and LAION-Rendition—that are strictly OOD to corresponding ImageNet and DomainNet test sets in terms of style. Training CLIP models on these datasets reveals that a significant portion of their performance is explained by in-domain examples. This indicates that the OOD generalization challenges from the ImageNet era still prevail and that training on web-scale data merely creates the illusion of OOD generalization. Furthermore, through a systematic exploration of combining natural and rendition datasets in varying proportions, we identify optimal mixing ratios for model generalization across these domains. Our datasets and results re-enable meaningful assessment of OOD robustness at scale—a crucial prerequisite for improving model robustness.

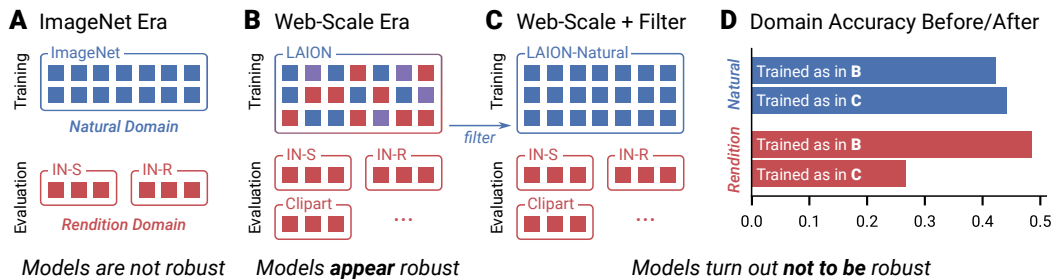


Figure 1: **Evaluated correctly, CLIP does not generalize across domains.** **A** Models used to be trained on a single domain like *natural* images from ImageNet (Russakovsky et al., 2015) and evaluated for out-of-domain (OOD) generalization on a different domain like *renditions* from test sets such as ImageNet-R (Hendrycks et al., 2021a), ImageNet-Sketch (Wang et al., 2019). **B** Today, large foundation models like CLIP (Radford et al., 2021) are trained on web-scale datasets such as LAION-400M (Schuhmann et al., 2021) containing images from many domains. Tested on a specific domain like renditions, CLIP exhibits unprecedented performance and appears robust. **C** We subsample from a deduplicated LAION-400M (Abbas et al., 2023) to obtain LAION-Natural, web-scale data set containing only natural images, which re-enables a meaningful assessment of CLIP’s generalization performance to renditions. **D** CLIP trained on LAION-Natural performs noticeably poorer on renditions, demonstrating that CLIP does not solve OOD generalization. The models are evaluated on refined test datasets containing samples only from their intended domains.

## 1 INTRODUCTION

Foundation models have revolutionized our world, demonstrating remarkable capabilities in solving grade school math problems, writing creative essays, generating stunning images, and comprehending visual content. One notable example is CLIP (Radford et al., 2021), a vision-language model

pre-trained on a vast dataset of image-text pairs, which forms the backbone of numerous other foundation models. CLIP has achieved unprecedented performance in various benchmarks across many domains—a stark difference to models in the ImageNet era, which struggled to generalize to unseen domains. This raises an important question: *Does CLIP solve out-of-domain generalization?* Out-of-domain (OOD) generalization refers to a model’s ability to perform well on data from domains other than its training (or *source*) domain. A *domain* is usually not rigorously defined and rather arises from collecting data in different contexts or environments. Nevertheless, some domains like the domain of *natural* images or the domain of *renditions* are delineated sufficiently clearly to enable the collection of datasets like ImageNet-Sketch (Wang et al., 2019), ImageNet-R (Hendrycks et al., 2021a), or DomainNet (Peng et al., 2019) for rigorous evaluation. CLIP’s impressive performance and generalization ability is primarily attributed to its extensive web-scale training set (Fang et al., 2022). Despite the large diversity of *natural* images in the training set, CLIP is likely to learn robust representations through exposure to many test domains during training. Indeed, Mayilvahanan et al. (2023) showed that CLIP’s training distribution contains exact or near duplicates of all commonly used OOD datasets but were also able to demonstrate that CLIP’s generalization performance remains high when correcting for this contamination. However, their analysis was only concerned with contamination on a data set level and failed to account for entire data domains. For example, even after their correction many *rendition* images remain in the training distribution (refer to Tab. 9). It is therefore unclear if CLIP will generalize to domain shifts if all datapoints from that domain are removed. We address this question with the following contributions:

- We develop a domain classifier that effectively distinguishes between *natural* images and *renditions*. We achieve this by labeling 19 000 random data points from LAION-400M for training and 6000 datapoints each from ImageNet and DomainNet test sets for evaluation.
- By applying the domain classifier to a deduplicated version of LAION-400M, we create two datasets: LAION-Natural, containing 57 million natural images, and LAION-Rendition, with 16 million renditions of objects and scenes. Additionally, we use the domain classifier to refine common OOD benchmarks by removing a small number of samples from an incorrect domain.
- Via our proposed LAION-Natural dataset, we demonstrate that CLIP trained on a single domain performs significantly worse on naturally-occurring domain shifts (see 1 for a summary). This indicates that CLIP’s strong performance is due to domain-contamination of the training data, rather than an inherent ability to generalize OOD.

## 2 RELATED WORK

**Measuring the OOD Generalization of CLIP Models** We aim to understand the OOD generalization capabilities of CLIP from a data-centric viewpoint. While multi-modal training with rich language captions does seem to contribute to robustness against distribution shifts (Xue et al., 2024), Fang et al. (2022) demonstrated that the nature of CLIP’s training distribution (as opposed to its mere size, its specific training objective, or natural language supervision) causes strong performance on various distribution shifts. However, it is unclear what aspects of the data distribution drive the robustness gains. Mayilvahanan et al. (2023) remove images highly similar to the test sets to show that data contamination and high perceptual similarity between training and test data do not explain generalization performance. While their data pruning technique removes some samples from LAION-400M that lie outside the natural image domain, they do not address domain generalization: They only account for the part of a domain covered by existing test sets and give no guarantee that all images of a given domain were removed. In another line of work, Nguyen et al. (2022) discover that a model’s effective robustness (Fang et al., 2022; Taori et al., 2020) on a test set interpolates when training data is compiled from various sources. However, they only consider mixing datasets that each cover multiple domains. In this work, we take their analysis further and show how mixing two data sources from distinct domains interpolates the effective robustness on those domains. Our study’s title is inspired by Gulrajani and Lopez-Paz (2021), who studied generalization from multiple distinct source domains. In contrast, we focus on generalization from single or mixed source domains to unseen domains. Overall, we aim for our work to be a valuable addition to the literature on OOD generalization (Liu et al., 2023; Koh et al., 2021; Madan et al., 2021; Gulrajani and Lopez-Paz, 2021; Madan et al., 2022; Arjovsky et al., 2019; Arjovsky, 2021).

### 3 BUILDING A DOMAIN CLASSIFIER

Our work hinges on filtering out datapoints that belong to specific domains from web-scale datasets. There is no precise definition for what constitutes a *domain* in general. Still, the community has come to agree on an implicit demarcation of the *natural* image and *renditions* domains by virtue of ImageNet compared to ImageNet-Sketch and ImageNet-R as well as DomainNet-Real compared to DomainNet-Sketch, -Quickdraw, -Infograph, -Clipart, and -Painting. Derived from the overall quality of an image, there is an intuitive, texture-centric notion of style us humans use which we adopt in this work. Further, we borrow methods from prior work that successfully classify images into different domains. We defer the reader to Sec. A for a thorough description of how we train and test our domain classifiers.

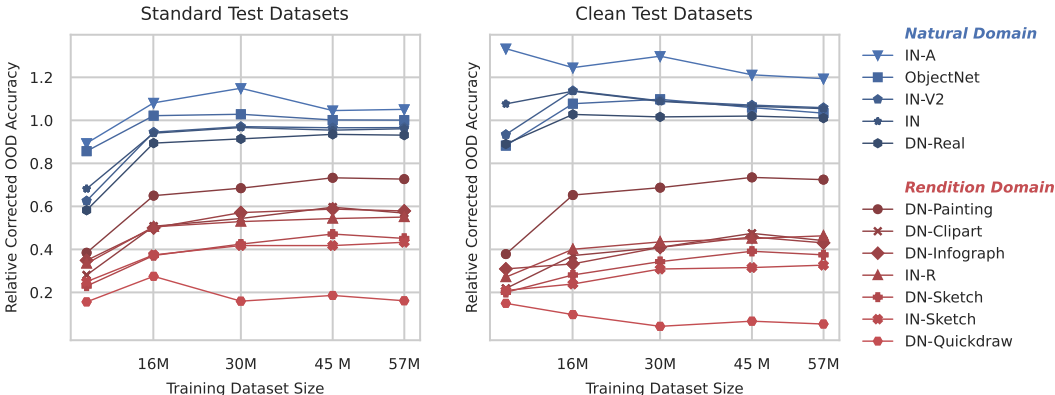


Figure 2: **Across scales, CLIP fails to generalize to unseen domains.** The *relative corrected OOD accuracy* shows performance losses or gains of a CLIP model trained exclusively on the *natural domain* via LAION-Natural to a CLIP model trained on a domain-contaminated dataset like LAION-200M. We evaluate on the original ImageNet and DomainNet test sets (left) and our cleaned versions of them (right, see Sec. 3.2 Without samples from the *rendition* domain, CLIP’s domain generalization ability suffers significantly and consistently across scales.

#### 3.1 DOMAIN COMPOSITION OF LAION-200M

We now deploy the chosen classifiers from Sec. 3 and label each sample in LAION-200M as *natural*, *rendition*, or *ambiguous*. We apply the classifiers with their strict thresholds at 98 % validation precision which yields a strong lower bound for the number of samples in each domain, as well as with their default thresholds which yields a more rounded estimate. From Tab. 9, it is clear that the LAION-200M contains a considerable portion of strictly stylistic images (with a lower bound of 7.90 % corresponding to 16 million images), and potentially many more images with some rendition elements are contained in the ambiguous group. We list further details in Sec. A.6.

Table 1: **Domain composition of LAION-200M.** We apply our *natural* and *rendition* domain classifiers with their strict thresholds at 98 % validation precision to get a lower bound of samples from each domain and with their default thresholds to obtain a more balanced estimate. Irrespective of the thresholding, LAION-200M still contains a large amount of renditions.

Classifier Precision		% Samples		
<i>Natural</i>	<i>Rendition</i>	<i>Natural</i>	<i>Ambiguous</i>	<i>Rendition</i>
0.79	0.77	60.74	25.41	13.86
0.98	0.98	28.40	63.70	7.90

### 3.2 CREATING SINGLE-DOMAIN DATASETS

To measure the true OOD performance of CLIP, we need to create a large dataset with only natural examples. We now use our trained domain classifiers at 98% validation-precision to subsample LAION-200M. We obtain LAION-Natural with roughly 57 million samples and LAION-Rendition with roughly 16 million samples. Figure 6 shows random samples from both datasets, more samples are shown in Figs. 19 and 20. We also deploy the domain classifiers on the ImageNet and DomainNet test sets to remove the domain-contamination reported above. The exact number of datapoints and the number of classes for each test set are detailed in Tab. 11. These datasets enable us to fairly assess CLIP’s domain generalization performance in the following sections.

Table 2: **Performance on the *rendition* domain is driven by renditions in the training data.** We compare CLIP trained without renditions on LAION-Natural to CLIP trained on datasets of the same size with renditions: LAION-Mix- $n$ M contains  $n$  million renditions, LAION-Rand is a random subset of LAION-200M with an estimated fraction of 7.9-13.86% renditions (see Tab. 9). Training with renditions greatly impacts performance on the *rendition* domain.

Dataset	Standard Datasets top-1 Acc.		Clean Datasets top-1 Acc.	
	<i>Natural</i>	<i>Rendition</i>	<i>Natural</i>	<i>Rendition</i>
LAION-Natural	36.88 %	21.98 %	39.72 %	18.75 %
LAION-Mix-12M	37.28 %	40.48 %	38.97 %	43.09 %
LAION-Mix-16M	36.92 %	41.46 %	38.58 %	41.46 %
LAION-Rand-57M	37.63 %	40.66 %	36.99 %	41.32 %

## 4 MEASURING CLIP’S OOD PERFORMANCE

For all our experiments, we train CLIP ViT-B/32 (Dosovitskiy et al., 2020) from scratch for 32 epochs with a batch size of 16384 on one node with either four or eight A100 GPUs (training takes several days, depending on dataset size). We use the implementation provided by Ilharco et al. (2021) and stick to their hyperparameters. We first train CLIP on the 57 M LAION-Natural and random subsets of it with 45 M, 30 M, and 16 M samples. We compare the classification accuracy of these models to that of CLIP models trained on random subsets of LAION-200M of the same sizes by reporting the accuracy ratio, which we refer to as *relative corrected OOD accuracy*. We measure this quantity on the original ImageNet and DomainNet test sets and their cleaned versions (see Sec. 3.2). Fig. 2 summarizes the results.

Across the board, we find that the relative corrected OOD accuracy on the clean datasets is around or above 1.0 for *natural* test sets, but drops to around 0.4 for most *rendition* test sets. This demonstrates that without domain-contamination of the training distribution, CLIP does not generalize across domains nearly as effectively as previously assumed. Notably, the relative corrected OOD accuracy is very consistent across dataset scales, allowing us to conjecture that this result holds also for CLIP models trained on much larger data sizes. To further reinforce this observation, we build LAION-Mix- $n$ M by replacing  $n$  million samples from LAION-Natural with samples from LAION-Rendition. We show in Tab. 2 that adding 13 or 16 million renditions has little effect on performance on the *natural* domain, but greatly improves performance on the *rendition* domain, highlighting the effect of domain-contamination.

To put the corrected OOD accuracy in context, we evaluate effective robustness (Fang et al., 2022; Taori et al., 2020) on the *natural* and *rendition* domain. To this end, Fig. 15 shows the top-1 classification accuracy of multiple CLIP models trained on LAION-200M, LAION-Natural, LAION-Rendition, LAION-Mix-13M, and ResNets trained on ImageNet (see Sec. C for details). We include LAION-Mix-13M as opposed to LAION-Mix-16M since it matches the effective robustness results for LAION-200M most closely. As usual, models with the same training regimen lie on a line and the  $y$ -distance of a model to the ImageNet line indicates its effective robustness. While all LAION-trained models achieve a similar effective robustness on the *natural* domain (Fig. 15 left), effective robustness on the *rendition* domain varies greatly and is notably lowest for LAION-Natural-trained models. Effective robustness plots on the individual datasets can be found in App. D. Together, the

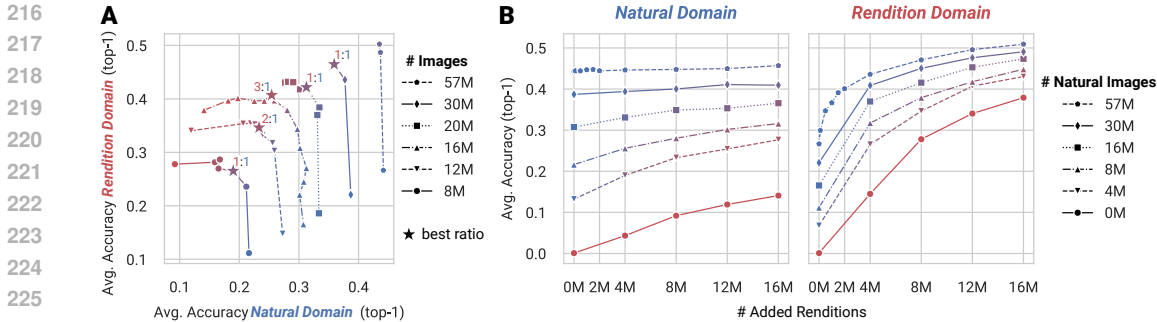


Figure 3: **A: Optimal data mixture.** We show the average accuracy on the *natural* and *rendition* domain for models trained with LAION-Mix of different absolute sizes and rendition-to-natural ratios (red indicates only renditions and blue only natural images). The best overall performance (corresponding to the point furthest from the origin) is achieved with a rendition-to-natural ratio between 1:1 and 3:1, which is consistent across scales. **B: Effect of adding renditions.** We also analyze model performance with increasingly more renditions added to a fixed-size training set of natural images (which increases overall dataset size). The amount of additional rendition samples required to reach a specific performance on the rendition domain depends on the number of natural samples included in the training set. While natural training samples give some performance boost on the rendition domain, rendition samples do this much more efficiently.

findings in this section demonstrate that CLIP’s unprecedented OOD generalization performance is a direct result of the domain-contamination of its training distribution. We defer a detailed discussion of Comparison of LAION training to ImageNet-training, Short-cut Learning, Domain Classification and Ambiguous Datapoints to Sec. 6.

## 5 INVESTIGATING DOMAIN MIXING AND SCALING EFFECTS

In the previous section, we explored training on single-domain datasets. Equipped with these clean datasets, we can now, for the first time, conduct a controlled investigation on what happens when large-scale datasets from different domains are mixed.

First, we show performance on the *natural* and *rendition* domain for models trained on LAION-Mix of different sizes and mixing ratios in Fig. 3A. Varying the mixing ratio while keeping the overall training set size constant reveals that a rendition-to-natural ratio between 1:3 and 1:1 achieves the best overall performance. This optimal range is consistent across training set sizes, although insights on larger scales are limited by the availability of LAION-Rendition samples (in total 16 million images). We hope our results can help practitioners while mixing such domains.

In our second experiment, we progressively add more rendition samples to fixed-size training sets of natural images (Fig. 3B). We find that models starting with more natural images require far fewer renditions to achieve the same performance on the rendition domain. This suggests that large amounts of natural images help the model learn some features that can be useful for generalizing to renditions, and relatively few additional renditions suffice to reach good performance on the rendition domain. In addition to boosting the performance on rendition test sets, adding rendition samples to the training set marginally boosts the performance on natural test sets, albeit with quickly diminishing returns. While performance in the natural domain benefits from rendition samples, natural samples are much more helpful. Likewise, training on few rendition samples gives higher performance than training on substantially more natural samples (see Fig. 3B, Tab. 2)—echoing our conclusion in Sec. 4 that CLIP does slightly generalize but much less than previously assumed.

## 6 DISCUSSION

**Contextualizing our core result** The literature often assumes that CLIP is capable of generalizing OOD (Radford et al., 2021; Abbasi et al., 2024; Nguyen et al., 2024; Fang et al., 2022; Li et al., 2023; Shu et al., 2023). Our main result is that CLIP’s strong generalization to rendition domains is largely

270 due to the presence of samples from those domains in its training distribution. Fang et al. (2022)  
271 showed CLIP’s robustness is tied to its data distribution but do not mention any specific characteristic.  
272 In contrast, Mayilvahanan et al. (2023) indicate that other dataset properties, not train-test similarity  
273 on a per-sample level, influence robustness. We conclusively demonstrate that CLIP’s apparent OOD  
274 robustness on standard OOD benchmarks like ImageNet-Sketch or ImageNet-R is often an artifact of  
275 overlapping domain data, rather than genuine OOD generalization. This refines the conclusion of  
276 Fang et al. (2022) and directly challenges Mayilvahanan et al. (2023) (see Appx. A.6 and Sec. 4), and  
277 several other works (Radford et al., 2021; Abbasi et al., 2024; Nguyen et al., 2024; Fang et al., 2022;  
278 Li et al., 2023; Shu et al., 2023). To the best of our knowledge, no work exists that addresses OOD  
279 generalization without domain contamination at this paper’s scale (10s of millions).

280 **Validity of conclusions for larger datasets** Although our training sets are constrained by the  
281 availability of natural and rendition samples, we believe that the insights gained from analyzing  
282 datasets with sizes spanning over one order of magnitude will remain applicable to even larger  
283 datasets. Specifically, the disparity in ‘relative corrected accuracy’ shown in Fig. 2 remains stable  
284 across dataset sizes from 4M to 57M. Similarly, effective robustness illustrated in Fig. 15 is influenced  
285 by the training distribution rather than the dataset size, which is also supported by findings in previous  
286 works (Miller et al., 2021; Fang et al., 2022; Mayilvahanan et al., 2023). Lastly, CLIP’s performance  
287 scales predictably across domain mixtures as shown in Sec. 5. Overall, we see no indication that our  
288 results should not transfer to larger scales.

289 **Validity of conclusions for other architectures and loss functions** Prior work strongly supports  
290 the generalizability of our findings on data contamination and optimal ratios across architectures and  
291 training methods beyond CLIP (Miller et al., 2021; Fang et al., 2022). For instance, Fang et al. (2022)  
292 demonstrates that CLIP’s robustness is driven primarily by the training distribution, with factors like  
293 dataset size, language supervision, and contrastive loss playing minimal roles. They also show that  
294 models trained on identical data distributions, regardless of loss functions (e.g., SimCLR+FT, CLIP,  
295 Supervised) or architectures (e.g., varying backbones and parameter sizes), exhibit similar effective  
296 robustness. This indicates that our conclusions are likely to hold across model types. We further  
297 address their validity across dataset sizes in Sec. 6.

298 **Choice of domain and validity of conclusions for other domains** For models to align with human  
299 perception, it is essential that they generalize to rendition domains, particularly in out-of-distribution  
300 (OOD) scenarios. Humans are adept at interpreting abstract visual renditions (Hendrycks et al.,  
301 2021a), while machines often depend primarily on textural cues (Geirhos et al., 2019). Consequently,  
302 we focus on natural images vs. renditions as our subject of study. Our methodology can be applied to  
303 evaluate OOD generalization for other domains, and we expect that our findings will hold true, as  
304 domain contamination is a general problem not tied to the specific domains we examined. However,  
305 we do anticipate challenges in accurately characterizing certain domain shifts, which could impede  
306 training the domain classifier. Nonetheless, if a small labeled dataset can be created to differentiate  
307 between these shifts, the subsequent processes should proceed smoothly. Given the manual effort  
308 required and the potential redundancy in findings, we defer this task to future work.

## 309 7 CONCLUSION

310 With the emergence of models trained on web-scale datasets containing abundant samples from  
311 seemingly all possible domains, the study of domain generalization mostly came to a halt. Hence, the  
312 question of how dataset scale actually affects the ability of models to generalize between domains  
313 remains unanswered. Here, we try to answer this question thoroughly by fully controlling the  
314 domains used for model training. By creating clean subsets of LAION containing either natural  
315 images or renditions, and by training models on various mixtures and dataset sizes, we show that  
316 the generalization performance of CLIP trained on only one domain drops to levels similar to what  
317 has been observed for ImageNet-trained models. Hence, we conclude that the domain generalization  
318 problem remains unsolved even for very large-scale datasets. We release all training set splits as  
319 well as pretrained models and encourage the field to re-consider domain generalization as a central  
320 benchmark for future progress on model architectures, inductive biases, and learning objectives.  
321  
322  
323

324 REPRODUCIBILITY STATEMENT

325  
326 We describe the methodology to create all of the datasets we use in Sec. 3.2, A.2, A.3. We also  
327 sketch the training details of all our models in Sec. A.4,4, C. This should be sufficient to reproduce  
328 all our datasets and experimental results. We aim to host our datasets and models shortly. The  
329 code to train the domain classifiers is available at [https://anonymous.4open.science/r/  
330 clip-dg-68D1/](https://anonymous.4open.science/r/clip-dg-68D1/).

331  
332 REFERENCES

- 333  
334 Amro Abbas, Kushal Tirumala, Dániel Simig, Surya Ganguli, and Ari S Morcos. Semdedup:  
335 Data-efficient learning at web-scale through semantic deduplication. *arXiv preprint*, 2023. URL  
336 <https://arxiv.org/abs/2303.09540>.
- 337  
338 Reza Abbasi, Mohammad Hossein Rohban, and Mahdieh Soleymani Baghshah. Deciphering the  
339 role of representation disentanglement: Investigating compositional generalization in clip models,  
340 2024. URL <https://arxiv.org/abs/2407.05897>.
- 341  
342 Martin Arjovsky. Out of Distribution Generalization in Machine Learning. *arXiv preprint*, 2021.  
343 URL <https://arxiv.org/abs/2103.02667>.
- 344  
345 Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant Risk Minimization.  
346 *arXiv preprint*, 2019. URL <https://arxiv.org/abs/1907.02893>.
- 347  
348 Zechen Bai, Yuta Nakashima, and Noa Garcia. Explain me the painting: Multi-topic knowledgeable  
349 art description generation. In *CVPR*, 2021.
- 350  
351 Andrei Barbu, David Mayo, Julian Alverio, William Luo, Christopher Wang, Dan Gutfreund, Josh  
352 Tenenbaum, and Boris Katz. Objectnet: A large-scale bias-controlled dataset for pushing the limits  
353 of object recognition models. *Advances in neural information processing systems*, 2019.
- 354  
355 Wei-Ta Chu and Yi-Ling Wu. Image style classification based on learnt deep correlation features.  
356 *IEEE Transactions on Multimedia*, 2018.
- 357  
358 Benjamin Cohen-Wang, Joshua Vendrow, and Aleksander Madry. Ask your distribution shift if  
359 pre-training is right for you. *arXiv preprint*, 2024a. URL [https://arxiv.org/abs/2403.  
360 00194](https://arxiv.org/abs/2403.00194).
- 361  
362 Benjamin Cohen-Wang, Joshua Vendrow, and Aleksander Madry. Ask your distribution shift if  
363 pre-training is right for you. *arXiv preprint*, 2024b. URL [https://arxiv.org/abs/2403.  
364 00194](https://arxiv.org/abs/2403.00194).
- 365  
366 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas  
367 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image  
368 is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint*, 2020. URL  
369 <https://arxiv.org/abs/2010.11929>.
- 370  
371 Alex Fang, Gabriel Ilharco, Mitchell Wortsman, Yuhao Wan, Vaishaal Shankar, Achal Dave, and  
372 Ludwig Schmidt. Data determines distributional robustness in contrastive language image pre-  
373 training (clip). In *ICML*, 2022.
- 374  
375 Noa Garcia and George Vogiatzis. How to read paintings: semantic art understanding with multi-  
376 modal retrieval. In *ECCV*, 2018.
- 377  
378 Leon A. Gatys, Alexander S. Ecker, and Matthias Bethge. A neural algorithm of artistic style. *ArXiv*,  
379 2015. URL <https://api.semanticscholar.org/CorpusID:13914930>.
- 380  
381 Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A. Wichmann, and  
382 Wieland Brendel. Imagenet-trained cnns are biased towards texture; increasing shape bias improves  
383 accuracy and robustness. In *ICLR*, 2019.
- 384  
385 Ishaan Gulrajani and David Lopez-Paz. In search of lost domain generalization. In *ICLR*, 2021.

- 378 Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul  
379 Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, Dawn Song, Jacob Steinhardt, and Justin Gilmer.  
380 The many faces of robustness: A critical analysis of out-of-distribution generalization. In *ICCV*,  
381 2021a.
- 382 Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial  
383 examples. In *CVPR*, 2021b.
- 384 Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori,  
385 Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali  
386 Farhadi, and Ludwig Schmidt. Openclip, July 2021. URL [https://doi.org/10.5281/  
387 zenodo.5143773](https://doi.org/10.5281/zenodo.5143773). If you use this software, please cite it as below.
- 388 Akshay Joshi, Ankit Agrawal, and Sushmita Nair. Art style classification with self-trained ensemble  
389 of autoencoding transformations. *arXiv preprint*, 2020. URL [https://arxiv.org/abs/  
390 2012.03377](https://arxiv.org/abs/2012.03377).
- 391 Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Balsub-  
392 ramani, Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Irena Gao, Tony Lee, Etienne  
393 David, Ian Stavness, Wei Guo, Berton A. Earnshaw, Imran S. Haque, Sara Beery, Jure Leskovec, An-  
394 shul Kundaje, Emma Pierson, Sergey Levine, Chelsea Finn, and Percy Liang. Wilds: A benchmark  
395 of in-the-wild distribution shifts, 2021. URL <https://arxiv.org/abs/2012.07421>.
- 396 Xuanlin Li, Yunhao Fang, Minghua Liu, Zhan Ling, Zhuowen Tu, and Hao Su. Distilling large  
397 vision-language model with out-of-distribution generalizability, 2023. URL [https://arxiv.  
398 org/abs/2307.03135](https://arxiv.org/abs/2307.03135).
- 399 Jiashuo Liu, Zheyang Shen, Yue He, Xingxuan Zhang, Renzhe Xu, Han Yu, and Peng Cui. Towards  
400 out-of-distribution generalization: A survey, 2023. URL [https://arxiv.org/abs/2108.  
401 13624](https://arxiv.org/abs/2108.13624).
- 402 Zhuang Liu and Kaiming He. A decade’s battle on dataset bias: Are we there yet? *arXiv preprint*,  
403 2024. URL <https://arxiv.org/abs/2403.08632>.
- 404 Spandan Madan, Timothy Henry, Jamell Dozier, Helen Ho, Nishchal Bhandari, Tomotake Sasaki,  
405 Frédo Durand, Hanspeter Pfister, and Xavier Boix. When and how cnns generalize to out-of-  
406 distribution category-viewpoint combinations, 2021. URL [https://arxiv.org/abs/2007.  
407 08032](https://arxiv.org/abs/2007.08032).
- 408 Spandan Madan, Li You, Mengmi Zhang, Hanspeter Pfister, and Gabriel Kreiman. What makes  
409 domain generalization hard? *arXiv preprint arXiv:2206.07802*, 2022.
- 410 Prasanna Mayilvahanan, Thaddäus Wiedemer, Evgenia Rusak, Matthias Bethge, and Wieland Brendel.  
411 Does clip’s generalization performance mainly stem from high train-test similarity? *arXiv preprint*,  
412 2023. URL <https://arxiv.org/abs/2310.09562>.
- 413 Orfeas Menis-Mastromichalakis, Natasa Sofou, and Giorgos Stamou. Deep ensemble art style  
414 recognition. In *2020 International Joint Conference on Neural Networks (IJCNN)*, 2020.
- 415 John P Miller, Rohan Taori, Aditi Raghunathan, Shiori Sagawa, Pang Wei Koh, Vaishaal Shankar,  
416 Percy Liang, Yair Carmon, and Ludwig Schmidt. Accuracy on the line: on the strong correlation  
417 between out-of-distribution and in-distribution generalization. In *ICML*, 2021.
- 418 Bac Nguyen, Stefan Uhlich, Fabien Cardinaux, Lukas Mauch, Marzieh Edraki, and Aaron Courville.  
419 Saft: Towards out-of-distribution generalization in fine-tuning, 2024. URL [https://arxiv.  
420 org/abs/2407.03036](https://arxiv.org/abs/2407.03036).
- 421 Thao Nguyen, Gabriel Ilharco, Mitchell Wortsman, Sewoong Oh, and Ludwig Schmidt. Quality not  
422 quantity: On the interaction between dataset design and robustness of clip. *NeurIPS*, 2022.
- 423 Xingchao Peng, Qinxun Bai, Xide Xia, Zijun Huang, Kate Saenko, and Bo Wang. Moment matching  
424 for multi-source domain adaptation. In *ICCV*, 2019.



- 432 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
433 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
434 models from natural language supervision. In *ICML*, 2021.
- 435 Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do imagenet classifiers  
436 generalize to imagenet? In *ICML*, 2019.
- 437 Evgenia Rusak, Steffen Schneider, Peter Vincent Gehler, Oliver Bringmann, Wieland Brendel, and  
438 Matthias Bethge. Imagenet-d: A new challenging robustness dataset inspired by domain adaptation.  
439 In *ICML 2022 Shift Happens Workshop*, 2022. URL [https://openreview.net/forum?](https://openreview.net/forum?id=LiC2vmzbpMO)  
440 [id=LiC2vmzbpMO](https://openreview.net/forum?id=LiC2vmzbpMO).
- 441 Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang,  
442 Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition  
443 challenge. *International journal of computer vision*, 2015.
- 444 Catherine Sandoval, Elena Pirogova, and Margaret Lech. Two-stage deep learning approach to the  
445 classification of fine-art paintings. *IEEE Access*, 2019.
- 446 Catherine Sandoval Rodriguez, Margaret Lech, and Elena Pirogova. Classification of style in fine-  
447 art paintings using transfer learning and weighted image patches. In *2018 12th International*  
448 *Conference on Signal Processing and Communication Systems (ICSPCS)*, 2018.
- 449 Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis,  
450 Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of  
451 clip-filtered 400 million image-text pairs. *arXiv preprint*, 2021. URL [https://arxiv.org/](https://arxiv.org/abs/2111.02114)  
452 [abs/2111.02114](https://arxiv.org/abs/2111.02114).
- 453 Yang Shu, Xingzhuo Guo, Jialong Wu, Ximei Wang, Jianmin Wang, and Mingsheng Long. Clipood:  
454 Generalizing clip to out-of-distributions, 2023. URL [https://arxiv.org/abs/2302.](https://arxiv.org/abs/2302.00864)  
455 [00864](https://arxiv.org/abs/2302.00864).
- 456 Gowthami Somepalli, Anubhav Gupta, Kamal Gupta, Shramay Palta, Micah Goldblum, Jonas  
457 Geiping, Abhinav Shrivastava, and Tom Goldstein. Measuring style similarity in diffusion models.  
458 *arXiv preprint*, 2024. URL <https://arxiv.org/abs/2404.01292>.
- 459 Rohan Taori, Achal Dave, Vaishaal Shankar, Nicholas Carlini, Benjamin Recht, and Ludwig Schmidt.  
460 Measuring robustness to natural distribution shifts in image classification. *NeurIPS*, 2020.
- 461 Haohan Wang, Songwei Ge, Zachary C. Lipton, and Eric P. Xing. Learning robust global representa-  
462 tions by penalizing local predictive power. In *NeurIPS*, 2019.
- 463 Sheng-Yu Wang, Alexei A Efros, Jun-Yan Zhu, and Richard Zhang. Evaluating data attribution for  
464 text-to-image models. In *CVPR*, 2023.
- 465 Ross Wightman, Hugo Touvron, and Hervé Jégou. Resnet strikes back: An improved training  
466 procedure in timm. *arXiv e-prints*, 2021.
- 467 Yihao Xue, Siddharth Joshi, Dang Nguyen, and Baharan Mirzasoleiman. Understanding the robust-  
468 ness of multi-modal contrastive learning to distribution shift. In *ICLR*, 2024.
- 469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485

## 486 A DETAILS ON THE DOMAIN CLASSIFIER

487  
488 We briefly talk about some related work in the domain classification literature in Sec. A.1. We  
489 then describe our labeling procedure based on this demarcation in Sec. A.2 and explore different  
490 ways to train a domain classifier on the resulting dataset in Sec. A.4. In Sec. A.6, we employ the  
491 best-performing classifier to analyze the composition of different training and test sets and finally use  
492 it to subsample LAION-Natural and LAION-Rendition in Sec. 3.2. For the remainder of this work,  
493 we substitute LAION-400M by LAION-200M, which we obtain by de-duplicating LAION-400M  
494 based on perceptual similarity as introduced by Abbas et al. (2023). They demonstrate that CLIP  
495 trained on LAION-200M obtains comparable downstream performance while greatly increasing data  
496 efficiency.

### 497 A.1 RELATED WORK

498  
499 **Domain Classification** The primary goal of our work necessitates creating web-scale datasets  
500 of different domains. This entails building a robust domain classifier that can reliably distinguish  
501 *natural images* from *renditions*. This task can be regarded as classifying the style of an image, which  
502 Gatys et al. (2015) proposed to measure using Gram matrices and which has been widely explored  
503 since then (Sandoval et al., 2019; Menis-Mastromichalakis et al., 2020; Sandoval Rodriguez et al.,  
504 2018; Joshi et al., 2020; Garcia and Vogiatzis, 2018; Chu and Wu, 2018; Bai et al., 2021). More  
505 recently, Cohen-Wang et al. (2024a) use a fine-tuned CLIP model from OpenCLIP (Ilharco et al.,  
506 2021) to distinguish between ImageNet and test sets with a domain shift, such as ImageNet-Sketch,  
507 ImageNet-R, and ImageNet-V2 (Recht et al., 2019). Wang et al. (2023) and Somepalli et al. (2024)  
508 develop a dataset classifier using a backbone trained by self-supervised learning and classification  
509 through retrieval via a database. Liu and He (2024) report high performance when training image  
510 classifiers to distinguish between different large-scale and diverse datasets.

### 511 A.2 LABELING

512  
513 LAION-200M contains diverse images from a multitude of sources. The images vary from naturally  
514 occurring to synthetically generated. We encourage the reader to glance at Fig. 19 to get a sense of  
515 the dataset and the difficulty of determining the domain of each image. As explained above, we aim  
516 to classify images belonging to the *natural* image or *rendition* domain. We also add an *ambiguous*  
517 class for images with elements of both domains and edge-cases.

518  
519 We provide the human annotator with a comprehensive set of guidelines derived from analyzing the  
520 existing OOD test sets, which we outline in App. A.3. In general, we adopt a *texture-centric* approach  
521 to distinguish renditions of a scene or object from their natural depictions. That is, depictions where  
522 *fine-grained texture information* is preserved are generally considered *natural*, while depictions with  
523 *simplified or flat textures* are considered *renditions*. Fig. 4 illustrates this demarcation on samples  
524 from LAION-200M, ImageNet test sets and DomainNet test sets.

525  
526 Overall, we label 19 000 random images from LAION-200M and 1000 images from each of the  
527 ImageNet and DomainNet distribution shifts (12 000 in total). Notably, almost all ImageNet and  
528 DomainNet test sets that are usually assumed to contain only images of a single domain exhibit some  
529 domain contamination. We discuss this in detail in Sec. A.6. Tab. 3 contains a detailed breakdown of  
530 labels for each data set. We show more samples grouped by domain for each data set in Figs. 22- 33.

### 531 A.3 LABELING

532  
533 As mentioned in Sec. A.2, we take a *texture-centric* approach in domain labeling. We resolve further  
534 ambiguities with respect to labeling in the following way:

- 535 • Natural objects with watermark or text, infographs with natural objects, signs with human  
536 symbol (eg. walking signal), objects with common logos (eg. Nike), naturalistic books  
537 or movie covers, images that are retro / low resolution / blurry / grainy / or with fake  
538 background but with texture information preserved, graphically altered natural images with  
539 significant texture information, and real objects with fake backgrounds **are all classified as  
natural.**

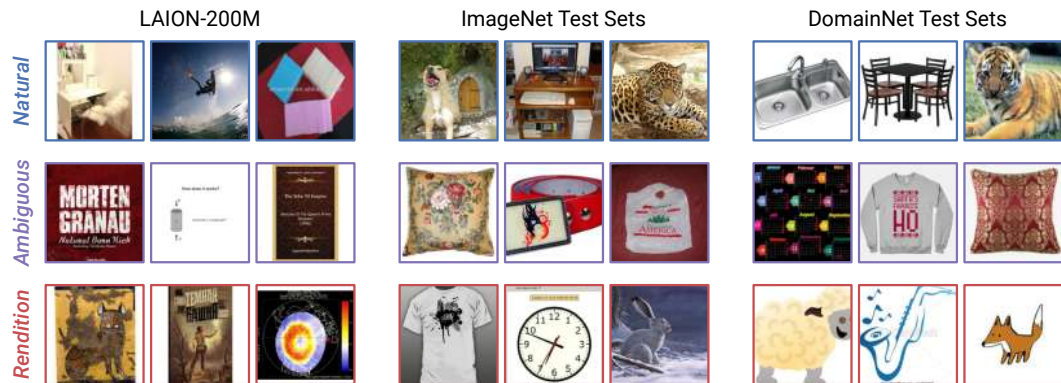


Figure 4: Labeled *Natural*, *ambiguous*, and *rendition* samples from different data sets. *Natural* images are photos or high-quality renders with minor filters that preserve *fine-grained textures*, while *renditions* are typically sketches, paintings, or graphics with *flat or simplified textures*. Images with elements of both, such as collages or natural images with large stylized elements, and images that mainly contain text are labelled as *ambiguous*.

- Stylistic: Infographs with stylized objects, stylized books or movie covers, retro / low resolution / blurry / grainy /graphically altered images with significant loss in texture information, stylized objects on plain or common natural background (eg. wall, bedsheet etc.) **are all classified as stylistic**.
- Ambiguous: Tattoos where hand / back is very visible, sculpture with real objects around, real images with distinct drawing of logos with objects, images that are retro / low resolution / blurry / grainy / or with fake background but with little texture information preserved **are all classified as ambiguous**.

To further ease the labeling procedure, we first build a rough binary classifier by fine-tuning CLIP ViT-L/14 with a linear readout to differentiate between some of the *natural* ImageNet and DomainNet test sets (namely, ImageNet-Val, ObjectNet (Barbu et al., 2019), ImageNet-V2 (Recht et al., 2019), ImageNet-A (Hendrycks et al., 2021b), and DomainNet-Real) and *stylistic* test sets (namely, ImageNet-Sketch, ImageNet-R, DomainNet-Painting, DomainNet-Sketch, and DomainNet-Clipart). We use this classifier to roughly pre-label samples and provide the annotator with 25 images from the same group at a time. This setup is shown in Fig. 5. The labeling was done by one labeler who labeled about 750-1000 images per hour. The labeler also did a checking of these labels by regrouping and going over them again. Below we visualize our labeling setup:

Final labeled images breakdown:

#### A.4 TRAINING AND CHOOSING THE DOMAIN CLASSIFIER

With the domain-labeled dataset, we can train a domain classifier to partition all of LAION-200M into *natural* images, *renditions*, or *ambiguous* images. Since we aim to obtain datasets that contain only images from a single domain we need a domain classifier that is as precise as possible. To this end, we train classifiers on 13 000 labelled LAION-200M images, retaining 3000 samples each for a validation and test set. From the domain classification literature discussed in Sec. ??, we evaluate four methods with publicly available code that we outline below. All methods build on CLIP ViT-L/14 pretrained on LAION-2B, which we choose for its balance between accuracy and inference speed.

**Contrastive Style Descriptors (CSD)** Somepalli et al. (2024) fine-tune pre-trained backbones via multi-label supervised contrastive learning and self-supervised learning with only style-preserving augmentations (random flips, resize, rotation). The resulting final-layer embeddings serve as style descriptors: During inference, they find the  $k$  stylistically nearest neighbors in a database of labelled images (e.g., the training set) by computing pairwise embedding-similarities to the test images. An image is classified as belonging to a style if at least one of the  $k$  neighbors has that style. We can

594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647



Figure 5: **Labeling setup.** By clicking on the image, the border changes to red, green, or blue, each representing natural, ambiguous, or rendition. By pressing the right or the left button the previous or next set of 25 images are rendered and the labels of the previous images are updated in a json file.

directly set up their method using the 13 000 labelled LAION-200M images as both the training set and the database for inference. From that, we obtain two binary classifiers, CSD-N (classifying natural vs. non-natural) and CSD-R (classifying renditions vs. non-renditions) that, together, can be used for our ternary classification.

**Density Ratios** Cohen-Wang et al. (2024b) aim to estimate the probability that a given sample is drawn from a reference distribution  $p_{\text{ref}}$ . Since high dimensional density estimation is challenging, they build a classifier to distinguish between a reference and a shifted distribution and compute the density ratio  $\frac{p_{\text{ref}}}{p_{\text{shifted}}}$  which they threshold at 0.2 to classify a given sample. We deploy their method unchanged to our task. We again obtain two binary classifiers, DR-N (classifying natural vs. non-natural) and DR-R (classifying renditions vs. non-renditions).

**Centroid Embeddings** Inspired by the baselines in Somepalli et al. (2024), we implement a simple model (embedding model plus linear readout) where we take the pretrained CLIP ViT-L/14 as the embedding model and create a linear readout by comparing to the centroid embeddings for each domain. We use this as a ternary untrained nearest-neighbor classifier, dubbed CE.

Table 3: **Number of labeled data points from several datasets and their domain-wise breakdown.** For training our domain classifier, we use the LAION-200M (Train), and LAION-200M (Val) for validation, and everything else to evaluate the final test performance.

Dataset	Natural	Stylistic	Ambiguous	Total
LAION-200M (Train)	7268	2978	2754	13000
LAION-200M (Val)	1000	1000	1000	3000
LAION-200M (Test)	1000	1000	1000	3000
ImageNet-A	974	7	19	1000
ObjectNet	917	2	81	1000
ImageNet-R	22	859	119	1000
ImageNet-Sketch	49	937	14	1000
ImageNet-V2	945	5	50	1000
ImageNet-Val	934	16	50	1000
DomainNet-Clipart	48	933	19	1000
DomainNet-Infograph	134	720	146	1000
DomainNet-Painting	101	795	104	1000
DomainNet-Quickdraw	0	1000	0	1000
DomainNet-Real	836	111	53	1000
DomainNet-Sketch	24	942	34	1000

**Fine-Tuning** We fine-tune the pretrained CLIP ViT-L/14 with a linear readout on the training dataset to obtain a ternary classifier, dubbed FT.

For the baselines (Cohen-Wang et al., 2024b; Somepalli et al., 2024), we simply use the training code detailed in their works and their public code. For the FT (Finetuning) model, as mentioned in Sec. A.4, we finetune a CLIP ViT-L/14 pretrained on LAION-2B with a linear readout. We finetune all models on 4 A100 GPUs, using a batch size of 256, weight decay of  $5e - 4$ , using an SGD optimizer, with step scheduler (0.1 every 20 epochs), at a learning rate of 0.1, for 50 epochs. All models converge. Each model took about 2 A100 GPU hours to train, therefore all the models took around 30 A100 GPU hours. The storage requirement for these datasets were less than 100 GB memory.

We use the validation set to determine the two best classifiers, one for natural images and one for renditions. Since the domain classifier should maximize precision above all else, we set the confidence threshold for each model such that it achieves 98% per-class precision. For CSD, we instead choose  $k$  to reach this precision. We then pick the classifier with the highest per-class recall to minimize the number of datapoints that are discarded when subsampling LAION-200M to build LAION-Natural and LAION-Rendition. We end up with FT, the fine-tuned ternary classifier, as our classifier for natural images, and DR-R, the binary classifier using density ratios as our rendition classifier. We use these classifiers for all subsequent experiments. Tab. 4 reports each model’s precision and recall on the *natural* and *rendition* class across ImageNet and DomainNet test sets. For raw accuracy numbers of all models, which in general are high for most, please refer to Tabs. 5 and 6 in App. A.5.

#### A.5 DOMAIN CLASSIFIER PERFORMANCE WITHOUT PRECISION THRESHOLDING

In Sec.A.4 we only compute the precision and recall obtained from the threshold at which we get 98% precision on LAION-200M Val domain dataset. We here report the accuracy of these classifiers on these test sets at their own standard precision of these models. We also train additional classifiers binary and ternary classifiers and by balancing the dataset sizes. To compare with the models from Cohen-Wang et al. (2024b), we train binary classifiers where we club natural with ambiguous and differentiate it from rendition (we name this FT-R), or we club rendition with ambiguous and differentiate it from natural (we name this FT-N). Further, we create several subsets for each of the ternary and the binary classification problem by balancing the number of datapoints in each class. We add the prefix ‘(balanced)’ to these models.

Table 4: We chose the **best natural classifier** and the **best rendition classifier** amongst binary classifiers based on Contrastive Style Descriptors (CSD) (Somepalli et al., 2024) and Density Ratios (DR) (Cohen-Wang et al., 2024b) as well as ternary classifiers using a linear readout based on either each domain’s centroid embedding (CE) or a fine-tuned CLIP (FT). All models use CLIP ViT-L/14 pretrained on LAION-2B. We report precision and recall on for the *natural* class (top) and *rendition* class (bottom) on ImageNet (IN) and DomainNet (DN) test sets and average performance across all test sets. Model hyperparameters are chosen for a validation precision of 98 % if possible. For each class, we select the classifier with the highest recall on the validation.

<i>cls=natural</i>		Val		Test		IN-Val		IN-v2		IN-A		ON		DN-R		Average	
Model		P	R	P	R	P	R	P	R	P	R	P	R	P	R	P	R
CSD-N	k=1	0.61	0.85	0.58	0.85	0.96	0.93	0.97	0.92	0.98	0.91	0.93	0.94	0.92	0.88	0.85	0.90
CSD-R	k=23	0.98	0.26	0.99	0.29	1.00	0.22	1.00	0.27	1.00	0.27	1.00	0.59	0.99	0.32	0.99	0.32
DR-N		0.98	0.00	0.50	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.21	0.00
DR-R		0.98	0.08	0.72	0.08	1.00	0.00	1.00	0.00	1.00	0.00	0.95	0.20	1.00	0.00	0.95	0.05
CE		0.98	0.35	0.89	0.33	0.95	0.02	1.00	0.04	1.00	0.02	0.99	0.16	0.99	0.11	0.97	0.15
sbbluedeep!50 FT		0.98	0.41	0.95	0.44	1.00	0.36	0.99	0.40	1.00	0.46	0.99	0.53	1.00	0.42	0.99	0.43

<i>cls=rendition</i>		Val		Test		IN-R		IN-S		DN-S		DN-Q		DN-P		DN-C		DN-I		Average	
Model		P	R	P	R	P	R	P	R	P	R	P	R	P	R	P	R	P	R	P	R
CSD-N	k=6	0.98	0.26	0.99	0.24	1.00	0.20	1.00	0.18	1.00	0.25	0.00	0.00	1.00	0.24	1.00	0.22	0.98	0.34	0.88	0.21
CSD-R	k=1	0.64	0.56	0.68	0.60	0.93	0.62	0.98	0.63	0.98	0.62	0.00	0.00	0.92	0.59	0.98	0.63	0.82	0.46	0.77	0.52
DR-N		0.98	0.20	0.98	0.23	1.00	0.29	1.00	0.20	1.00	0.27	1.00	0.01	1.00	0.28	1.00	0.28	0.98	0.11	0.99	0.21
sbreddeep!50 DR-R		0.98	0.35	0.98	0.41	1.00	0.60	1.00	0.71	1.00	0.74	1.00	0.33	0.99	0.60	1.00	0.65	0.98	0.39	0.99	0.53
white CE		0.98	0.11	0.99	0.12	0.99	0.43	1.00	0.39	1.00	0.30	1.00	0.09	0.98	0.47	1.00	0.38	1.00	0.01	0.99	0.26
FT		0.98	0.27	0.95	0.26	1.00	0.38	1.00	0.57	1.00	0.61	1.00	0.68	1.00	0.21	1.00	0.50	1.00	0.30	0.99	0.42

Table 5: Accuracy on each of the natural test sets on class natural without thresholding. Some classifiers give the illusion of being good but have very low precision or recall(see Sec. A.4).

Model	(Val)	(Test)	IN-Val	IN-V2	IN-A	ON	DN-R	DN-I
FT	0.90	0.89	0.93	0.94	0.96	0.95	0.94	0.72
CE	0.75	0.78	0.80	0.84	0.86	0.95	0.81	0.19
FT-N	0.89	0.90	0.94	0.95	0.97	0.97	0.93	0.49
DR-N (balanced)	0.89	0.91	0.94	0.94	0.95	0.98	0.92	0.50
DR-R	0.98	0.97	0.99	0.99	1.00	1.00	0.97	0.90
FT (balanced)	0.78	0.82	0.84	0.86	0.86	0.88	0.83	0.46
FT-R	0.96	0.95	0.93	0.95	0.97	0.98	0.96	0.90
FT-N (balanced)	0.85	0.85	0.92	0.95	0.96	0.95	0.91	0.43
DR-R (balanced)	0.93	0.92	0.93	0.94	0.95	0.99	0.90	0.75
FT-R (balanced)	0.86	0.86	0.88	0.88	0.90	0.89	0.88	0.84
DR-N	0.93	0.92	0.94	0.95	0.94	0.99	0.92	0.76

## A.6 ANALYZING THE DOMAIN MAKE-UP OF DIFFERENT DATA SETS

Both ImageNet and DomainNet are web-scraped datasets that were refined through extensive human annotation. In contrast, LAION-400M is obtained purely through web scraping without subsequent human domain filtering. Since human annotators can make mistakes, and LAION-400M’s domain composition is inherently unknown, we use our domain classifiers to understand it.

To this end, we deploy the chosen classifiers from Sec. 3 and label a sample *ambiguous* if the *natural* and *rendition* classifier disagree. We apply the classifiers both with their strict thresholds at 98 % validation precision which yields a strong lower bound for the number of samples in each domain, as well as with their default thresholds which yields a more rounded estimate. From Tab. 9, it is clear that the LAION-200M contains a considerable portion of strictly stylistic images (with a lower bound of 7.90 % corresponding to 16 million images), and potentially many more images with some rendition elements are contained in the ambiguous group. In contrast, for ImageNet, we find a much smaller fraction of renditions (at least 0.4 % of samples). We additionally observe that many evaluation datasets are considerably domain-contaminated (at least 5 % of samples stem from

Table 6: **Accuracy on each of the rendition test sets on class natural without thresholding.** Some classifiers give the illusion of being good but have very low precision or recall(see Sec. A.4).

Model	(Val)	(Test)	IN-R	IN-S	DN-S	DN-Q	DN-P	DN-C	DN-I
DR-R	0.77	0.80	0.93	0.98	0.98	0.96	0.92	0.93	0.88
FT (balanced)	0.78	0.88	0.82	0.94	0.94	0.91	0.80	0.85	0.77
FT	0.76	0.75	0.75	0.91	0.90	0.95	0.73	0.80	0.74
DR-N	0.89	0.92	0.99	0.99	0.99	0.98	0.97	0.97	0.94
FT-R	0.69	0.68	0.69	0.81	0.80	0.79	0.65	0.72	0.67
DR-N (balanced)	0.93	0.94	0.97	0.99	0.99	1.00	0.95	0.94	0.99
FT-R (balanced)	0.86	0.84	0.80	0.92	0.91	0.90	0.75	0.83	0.88
CE	0.61	0.62	0.95	0.90	0.89	0.96	0.95	0.93	0.32
DR-R (balanced)	0.90	0.93	0.99	0.99	0.99	0.99	0.98	0.97	0.96
FT-N	0.84	0.83	0.72	0.83	0.82	0.48	0.63	0.77	0.97
FT-N (balanced)	0.87	0.86	0.75	0.93	0.91	0.96	0.64	0.88	0.98

Table 7: **Domain composition of training sets.** We apply our *natural* and *rendition* domain classifiers with their strict thresholds at 98 % validation precision to get a lower bound of samples from each domain and with their default thresholds to obtain a more balanced estimate. ImageNet-Train has a much smaller fraction of *rendition* samples than LAION-200M. We also note that ‘combined-pruned’, the training set from Mayilvahanan et al. (2023) that corrected for test set contamination still contains a large fraction of renditions.

Dataset	# Samples	Classifier Precision				
		<i>Natural</i>	<i>Rendition</i>	<i>Natural</i>	<i>Ambiguous</i>	<i>Rendition</i>
LAION-200M	199 663 250	0.79	0.77	60.74 %	25.41 %	13.86 %
		0.98	0.98	28.40 %	63.70 %	7.90 %
ImageNet-Train	1 281 167	0.79	0.77	89.20 %	9.62 %	1.18 %
		0.98	0.98	36.00 %	63.60 %	0.40 %
combined-pruned	187 471 515	0.79	0.77	62.98 %	25.18 %	11.83 %
		0.98	0.98	29.58 %	64.02 %	6.40 %

the opposite domain), especially ImageNet-R, DomainNet-Real, DomainNet-Clipart, DomainNet-Painting, and DomainNet-Infograph (refer to Tab. 8, App. A.7).

LAION-Natural ~57 million samples



LAION-Rendition ~16 million samples



Figure 6: **Random samples from LAION-Natural and LAION-Rendition.**

We also analyze the domain composition of datasets from Mayilvahanan et al. (2023), who created several subsets of LAION-200M that do not contain samples that are perceptually *highly similar* to ImageNet OOD test sets. These removed images are expected to be (near-) duplicates of test images in terms of both content and style. Their dataset ‘combined-pruned’ is a subset of LAION-200M where highly similar images to ImageNet-Sketch, ImageNet-R, ImageNet-Val2, ImageNet-Val, ImageNet-A, and ObjectNet were pruned. In their work, it remained unclear whether pruning also

effectively removed all images of the rendition domain, which we can now answer. Tab. 9 reveals that a considerable number of renditions remains in the pruned dataset (at least 6.4% corresponding to around 11 million images). These remaining renditions might have played a significant role in the generalization performance of their CLIP models, especially on ImageNet-Sketch and ImageNet-R. As a result, CLIP’s domain generalization performance is yet to be evaluated fairly.

#### A.7 DOMAIN COMPOSITION AT DIFFERENT PRECISION

We provide a detailed overview over the domain composition of datasets at standard precision in Table 8, and over the domain composition of datasets at 98% precision in Table 9.

Table 8: **Domain composition of datasets at standard precision (without thresholding).** The first three columns show the fraction of samples in the original dataset classified as natural, stylistic, or ambiguous, respectively, while the latter column shows the dataset’s total number of samples.

Dataset	Natural [%]	Stylistic [%]	Ambiguous [%]	Total
LAION-200M	60.74	13.86	25.41	199 663 250
ImageNet (Train)	89.2	1.18	9.62	1 281 167
ImageNet (Val)	89.1	1.18	9.72	50 000
ObjectNet	90.22	0.1	9.68	18 574
ImageNet-V2	88.49	1.38	10.13	10 000
ImageNet-A	93.79	0.52	5.69	7 500
ImageNet-R	9.75	64.42	25.83	30 000
ImageNet-Sketch	3.69	85.34	10.97	50 889
DomainNet-Real	80.07	7.59	12.34	175 327
DomainNet-Quickdraw	1.35	93.27	5.38	172 500
DomainNet-Clipart	8.28	75.89	15.83	48 833
DomainNet-Painting	13.97	56.33	29.7	75 759
DomainNet-Sketch	3.1	84.18	12.71	70 386
DomainNet-Infograph	11.17	53.41	35.41	53 201

Table 9: **Domain composition of datasets at 98% precision.** The first three columns show the fraction of samples in the original dataset classified as natural, stylistic, or ambiguous, respectively, while the latter column shows the dataset’s total number of samples.

Dataset	Natural [%]	Stylistic [%]	Ambiguous [%]	Total
LAION-200M	28.4	7.9	63.7	199 663 250
ImageNet (Train)	36.0	0.4	63.6	1 281 167
ImageNet (Val)	35.73	0.37	63.9	50 000
ObjectNet	50.32	0.0	49.68	18 574
ImageNet-V2	36.04	0.29	63.67	10 000
ImageNet-A	43.25	0.16	56.59	7 500
ImageNet-R	3.56	52.82	43.61	30 000
ImageNet-Sketch	1.21	67.92	30.87	50 889
DomainNet-Real	34.31	3.98	61.71	175 327
DomainNet-Quickdraw	0.09	34.41	65.5	172 500
DomainNet-Clipart	3.46	62.53	34.01	48 833
DomainNet-Painting	5.3	47.55	47.15	75 759
DomainNet-Sketch	1.38	69.58	29.04	70 386
DomainNet-Infograph	1.59	28.11	70.3	53 201



## A.8 ON THE DOMAIN COMPOSITION OF MAYILVAHANAN ET AL. (2023)

Please find in Tab. 10 the exact number of rendition examples calculated by deploying our domain classifier on each the 3 datasets (pruned using rendition test sets) from Mayilvahanan et al. (2023). We see that at least 11-13M images are not pruned away from the datasets, therefore explaining the insignificant drop in performance.

Table 10: **Number datapoints within the dataset vs number of datapoints pruned away in Mayilvahanan et al. (2023).**

Dataset	Size	Within	Pruned
sketch-pruned	191 481 491	24 016 047	3 654 180
r-pruned	194 088 525	24 304 991	3 365 236
combined-pruned	187 471 515	22 173 006	5 497 221
sketch-pruned (98% precision)	19 1481 491	13 266 999	2 482 751
r-pruned (98% precision)	194 088 525	13 338 759	2 410 991
combined-pruned (98% precision)	187 471 515	11 999 276	3 750 474

## A.9 PREPARING CLEAN DATASETS

In Sec. 3.2, we created several train and test sets from LAION-200M and ImageNet / DomainNet shifts respectively, by deploying our classifier at 98% precision. The exact number of samples and the number of (remaining) classes are in Tab. 11.

Table 11: **Clean datasets composition.** Obtained by deploying the domain classifiers from Sec.A.4 at 98% precision.

Dataset	Classes	Size
LAION-Natural	-	56 685 759
LAION-Stylistic	-	15 749 750
ImageNet-Val	985	17 864
ImageNet-V2	926	3 604
ImageNet-Sketch	991	34 564
ImageNet-R	200	15 847
ImageNet-A	197	3 244
ObjectNet	113	9 347
DomainNet-Real	339	60 148
DomainNet-Quickdraw	344	59 353
DomainNet-Infograph	345	14 957
DomainNet-Clipart	345	30 536
DomainNet-Sketch	344	48 974
DomainNet-Painting	345	36 020

## B NOTES ON THE CLIP MODELS

### B.1 RESOURCES SPENT

We train about 28 CLIP ViT-B/32 models on several subsets of LAION-200M. These models took about 8000 A100 GPU hours. We also needed about 18 TB of memory to store these datasets.

## B.2 RAW ACCURACY NUMBERS OF CLIP TRAINED ON LAION-N VS LAION

In Sec. 4, in Fig. 2, we only reported the relative numbers. Here, in Fig. 7, 9, 8, 10, we report the actual numbers as a function of dataset size.

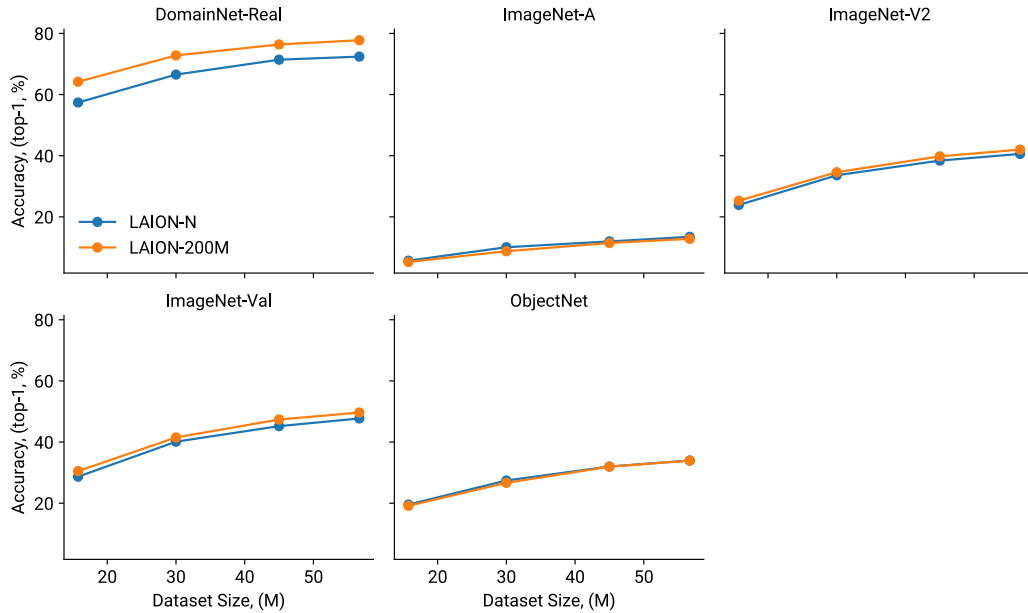


Figure 7: CLIP trained on LAION v LAION-N performance on standard natural test sets.

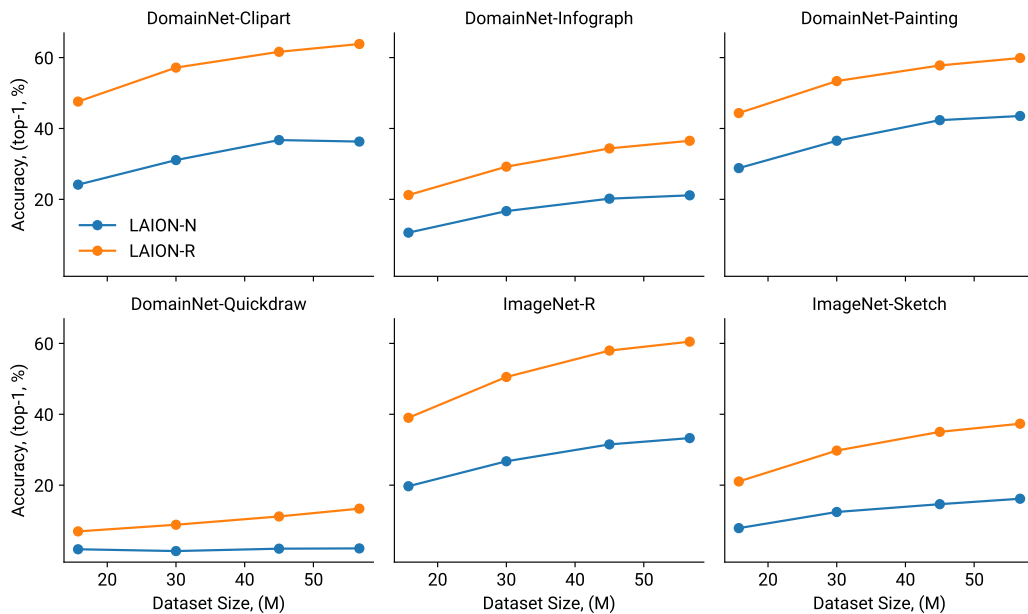


Figure 8: CLIP trained on LAION v LAION-N performance on standard rendition test sets.

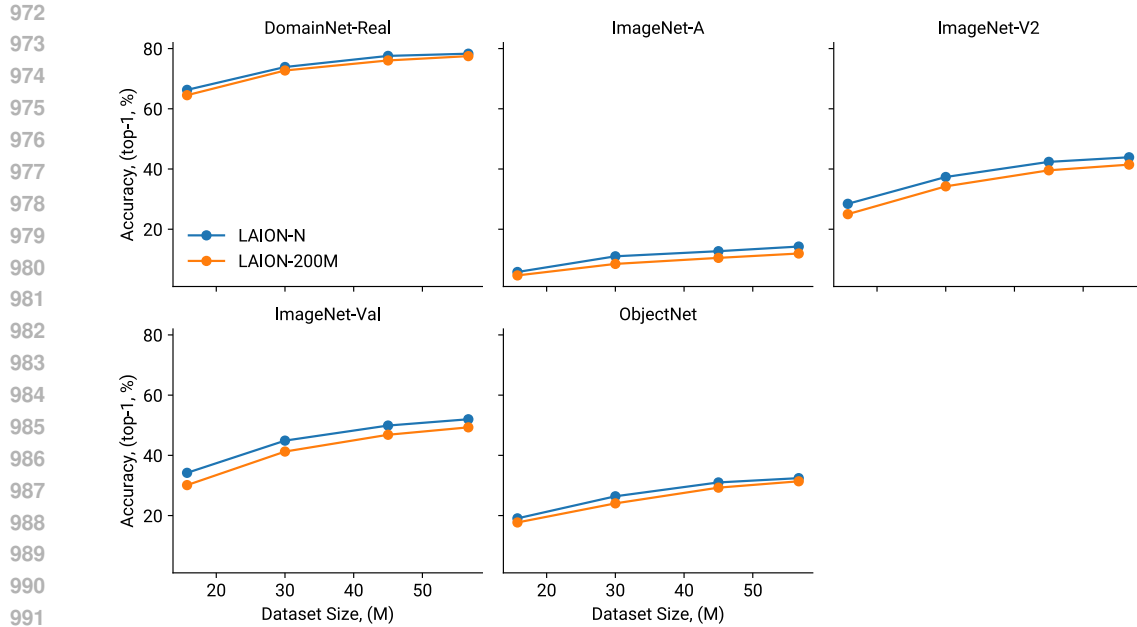


Figure 9: CLIP trained on LAION v LAION-N performance on clean natural test sets.

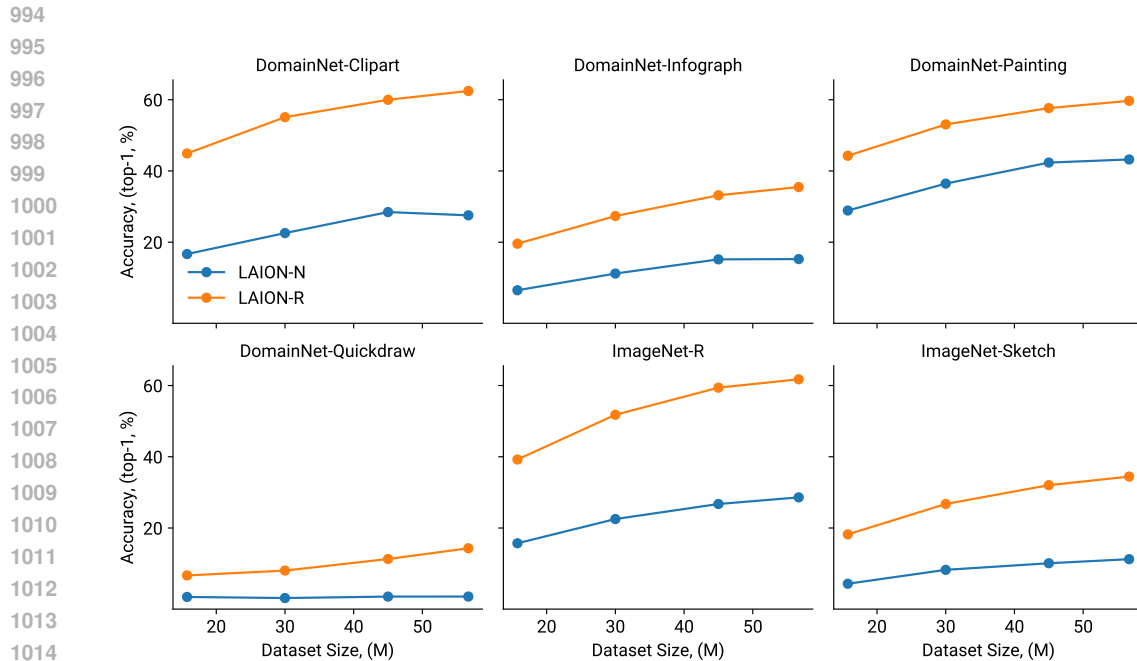


Figure 10: CLIP trained on LAION v LAION-N performance on clean rendition test sets.

## C TRAINING RESNETS ON IMAGENET

We deploy our natural domain classifier from Sec/3 at 90% precision (threshold obtain from LAION 13K Val set) on ImageNet-Train to obtain about 1M datapoints belonging to the natural domain (dubbed ImageNet-N). We create several datasets of smaller sizes subsampling from ImageNet-N. We also create randomly sampled datasets of similar sizes from the original ImageNet. We train ResNet-50 models on all of these datasets. We follow the training recipe A3 of Wightman et al. (2021) and train the models for 200 epochs. We then evaluate these models on standard test sets and

1026 clean test sets from Sec.3.2. The accuracies of ResNets trained on subsets of original ImageNet is  
 1027 used for the effective robustness plots in Sec. 4, D. Further, the comparison of accuracies between  
 1028 the models trained on subsets from ImageNet-N and ImageNet is in Fig. 11, 13, 12, 14. As such  
 1029 there is no significant performance difference anywhere, thus indicating that ImageNet does not have  
 1030 substantial domain leakage.

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

1052 **Figure 11: Resnets trained on ImageNet v ImageNet-N performance on standard natural test**  
 1053 **sets.**

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1072

1073

1074

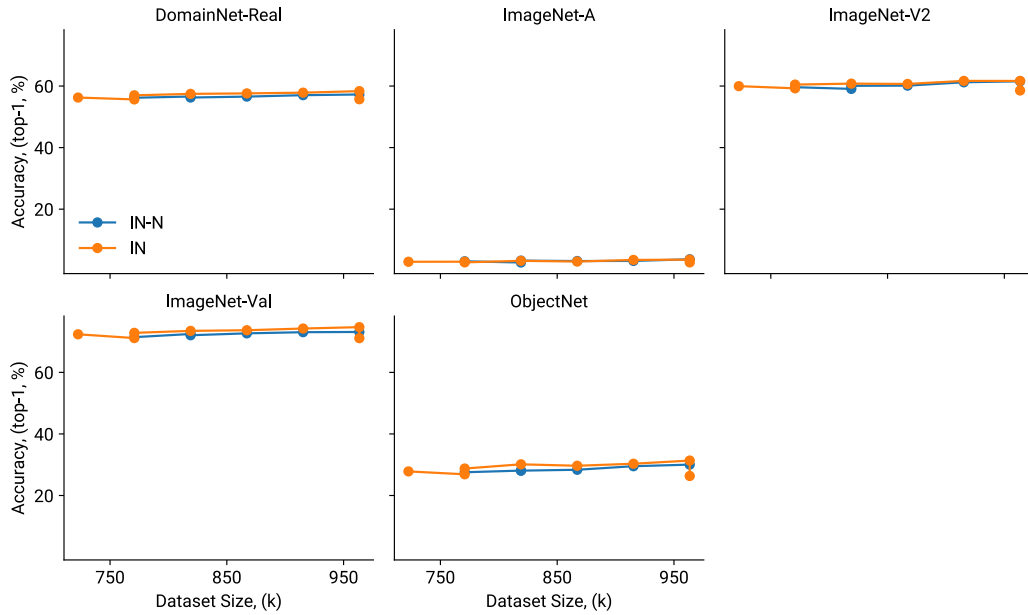
1075

1076

1077

1078

1079



1052 **Figure 11: Resnets trained on ImageNet v ImageNet-N performance on standard natural test**  
 1053 **sets.**

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1072

1073

1074

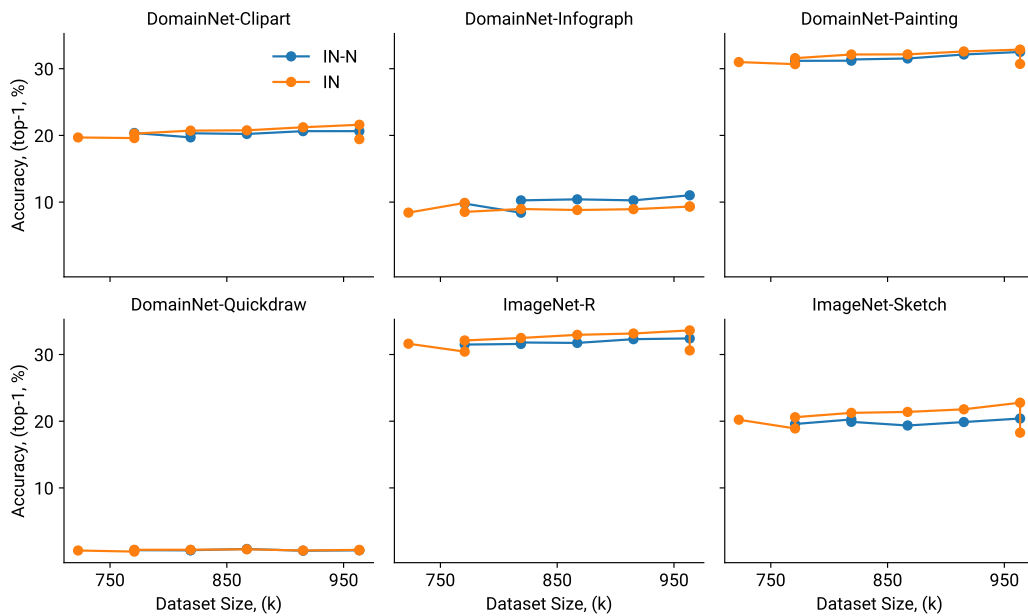
1075

1076

1077

1078

1079



1077 **Figure 12: Resnets trained on ImageNet v ImageNet-N performance on standard rendition test**  
 1078 **sets.**

1080  
 1081  
 1082  
 1083  
 1084  
 1085  
 1086  
 1087  
 1088  
 1089  
 1090  
 1091  
 1092  
 1093  
 1094  
 1095  
 1096  
 1097  
 1098  
 1099  
 1100  
 1101  
 1102  
 1103  
 1104  
 1105  
 1106  
 1107  
 1108  
 1109  
 1110  
 1111  
 1112  
 1113  
 1114  
 1115  
 1116  
 1117  
 1118  
 1119  
 1120  
 1121  
 1122  
 1123  
 1124  
 1125  
 1126  
 1127  
 1128  
 1129  
 1130  
 1131  
 1132  
 1133

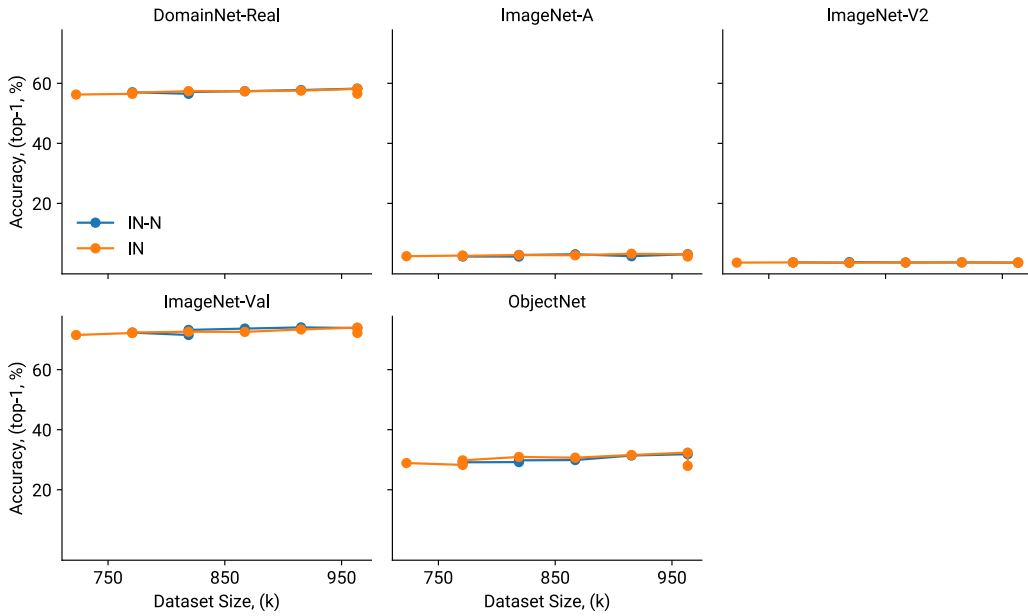


Figure 13: Resnets trained on ImageNet v ImageNet-N performance on clean natural test sets.

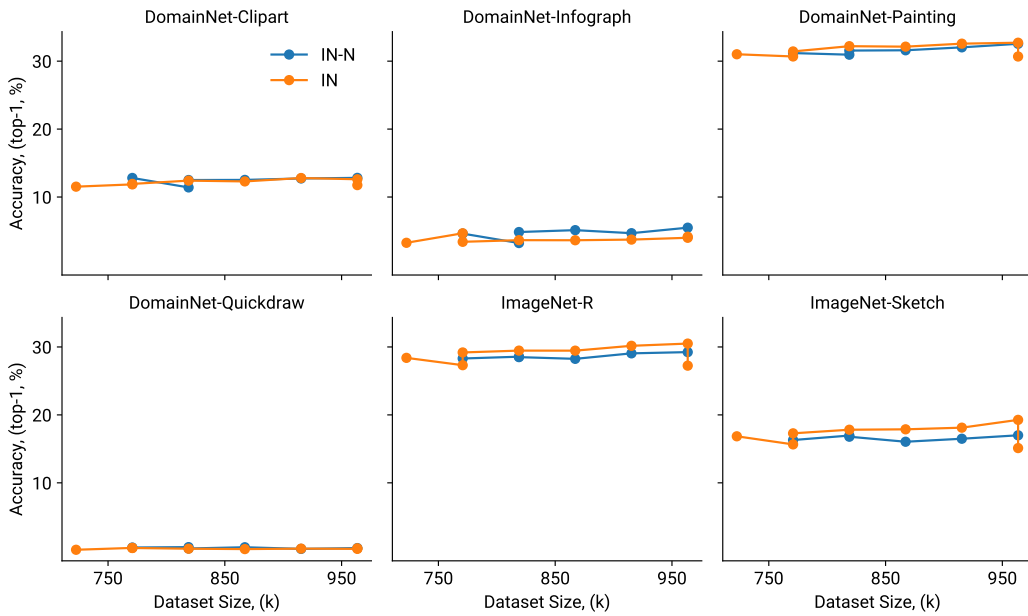


Figure 14: Resnets trained on ImageNet v ImageNet-N performance on clean rendition test sets.

D DETAILED EFFECTIVE ROBUSTNESS PLOTS ON INDIVIDUAL SHIFTS

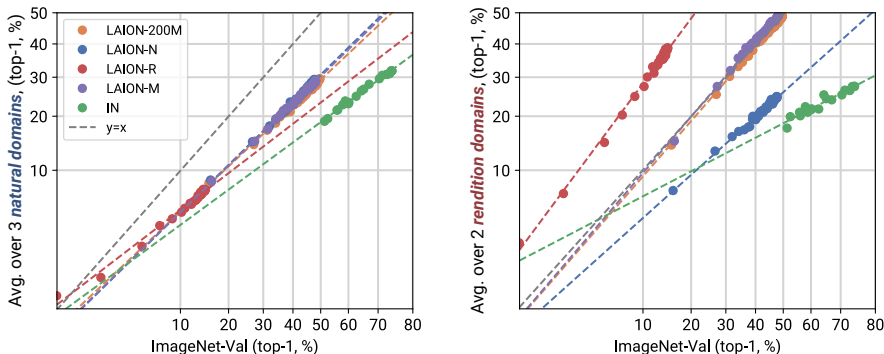


Figure 15: **CLIP’s effective robustness to renditions is driven by domain-contamination.** We evaluate effective robustness (Fang et al., 2022; Taori et al., 2020) for models trained on different LAION-200M subsets. Most notably, CLIP trained on LAION-Natural matches the effective robustness of a LAION-200M-trained CLIP on the *natural* domain (left), but has significantly lower effective robustness on the *rendition* domain, indicating that CLIP requires rendition samples in its training distribution to perform well on this domain.

In Fig. 15, we report aggregated results where we average over natural and stylistic ImageNet distribution shifts. We display the results on the individual distribution shifts in Fig. 16. On ImageNet-R and ImageNet-Sketch (bottom row), we observe that the effective robustness of the CLIP models can be modulated by training it on the different dataset splits, i.e. LAION-Natural, LAION-Rendition, LAION-Mix. The model trained on LAION-Natural is much closer to the ImageNet trained model in terms of effective robustness compared to the model trained on LAION-Rendition. In contrast, effective robustness is barely affected on the natural splits (top row). This can be explained by the final data distributions of the different training splits: Our filtering procedure does not affect natural images which are most responsible for the performance on natural datasets which explains the consistency in performance.

We also investigate effective robustness on the DomainNet shifts in Fig. 17. We note that the ImageNet model’s accuracy numbers on DomainNet are not comparable to the CLIP models because the ImageNet model has been evaluated on a subset of DomainNet (ImageNet-D, Rusak et al., 2022) which is compatible with ImageNet classes. DomainNet has many classes which are not present in ImageNet, such as for example “The Great Wall of China” or “paper clip” which have been removed in ImageNet-D to enable evaluating ImageNet trained models without the need for training an additional readout layer. In contrast, we evaluate the CLIP trained models on the full DomainNet splits following standard zero-shot evaluation procedure. We will add a Figure where we control for the missing classes and evaluate the CLIP models on ImageNet-D in the next version of the manuscript.

On DomainNet, we similarly observe strong changes in effective robustness of the CLIP trained models when evaluating on the stylistic domains (all domains except for DomainNet-Real), and barely any changes when evaluating on the DomainNet-Real domain.

E VISUALIZATION OF ERRORS MADE BY THE DOMAIN CLASSIFIER

We show images which have been misclassified by our domain classifier Fig. 18. We observe that the errors are interpretable. For example, the “natural” images which have been classified as “ambiguous” are indeed ambiguous: We see a sculpture in one image, a large woodwork of an ant in another and a pencil drawing of an airplane with a partly visible human hand drawing it in a third image.

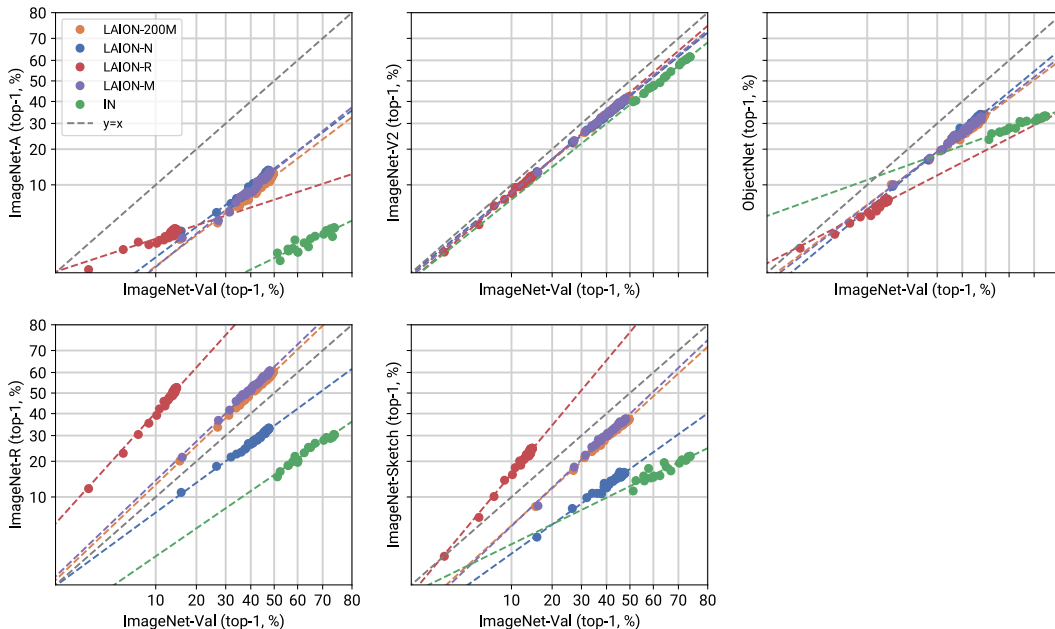


Figure 16: **Effective Robustness of different models on different ImageNet distribution shifts.** On ImageNet-R and ImageNet-Sketch (bottom row), we observe that the effective robustness of the CLIP models can be modulated by training it on the different dataset splits, i.e. LAION-Natural, LAION-Rendition, LAION-Mix. The model trained on LAION-Natural is much closer to the ImageNet trained model in terms of effective robustness compared to the LAION-Rendition model.

## F VISUALIZATION OF SAMPLES FROM THE LAION DATASET

We visualize random examples from the “Natural”, “Rendition” and “Ambiguous” domains from LAION in Figs. 19-21.

## G VISUALIZATIONS OF IMAGENET DISTRIBUTION SHIFTS

We visualize random examples from the “Natural”, “Rendition” and “Ambiguous” domains from the considered ImageNet shifts datasets in Figs. 22-27. We show 20 images per split; occasionally, there are fewer than 20 images in some of these splits, such as e.g. there are very few renditions in ImageNet-A. In that case, we plot all images from that split and leave the remaining subplots blank.

## H VISUALIZATIONS OF DOMAINNET DISTRIBUTION SHIFTS

We visualize random examples from the “Natural”, “Rendition” and “Ambiguous” domains from different DomainNet datasets in Figs. 28-33. We show 20 images per split; occasionally, there are fewer than 20 images in some of these splits, such as e.g. no natural images in the Quickdraw domain. In that case, we plot all images from that split and leave the remaining subplots blank.

1242  
 1243  
 1244  
 1245  
 1246  
 1247  
 1248  
 1249  
 1250  
 1251  
 1252  
 1253  
 1254  
 1255  
 1256  
 1257  
 1258  
 1259  
 1260  
 1261  
 1262  
 1263  
 1264  
 1265  
 1266  
 1267  
 1268  
 1269  
 1270  
 1271  
 1272  
 1273  
 1274  
 1275  
 1276  
 1277  
 1278  
 1279  
 1280  
 1281  
 1282  
 1283  
 1284  
 1285  
 1286  
 1287  
 1288  
 1289  
 1290  
 1291  
 1292  
 1293  
 1294  
 1295

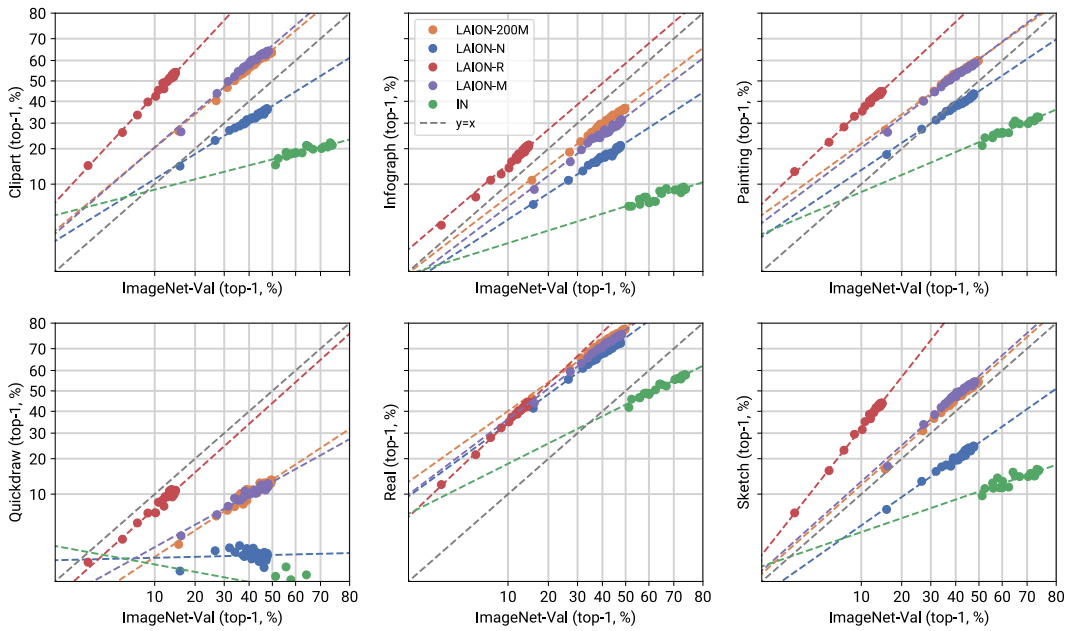


Figure 17: **Effective Robustness of different models on different DomainNet distribution shifts.** On the stylistic domains, we observe that the effective robustness of the CLIP models can be modulated by training it on the different dataset splits, i.e. LAION-Natural, LAION-Rendition, LAION-Mix. Effective robustness barely changes when evaluating different CLIP models on DomainNet-Real.



1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349

		Predicted					
		Natural		Rendition		Ambiguous	
True	Natural						
	Rendition						
	Ambiguous						

Figure 18: Confusion matrix of example images which have been misclassified by our domain classifier.

1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403

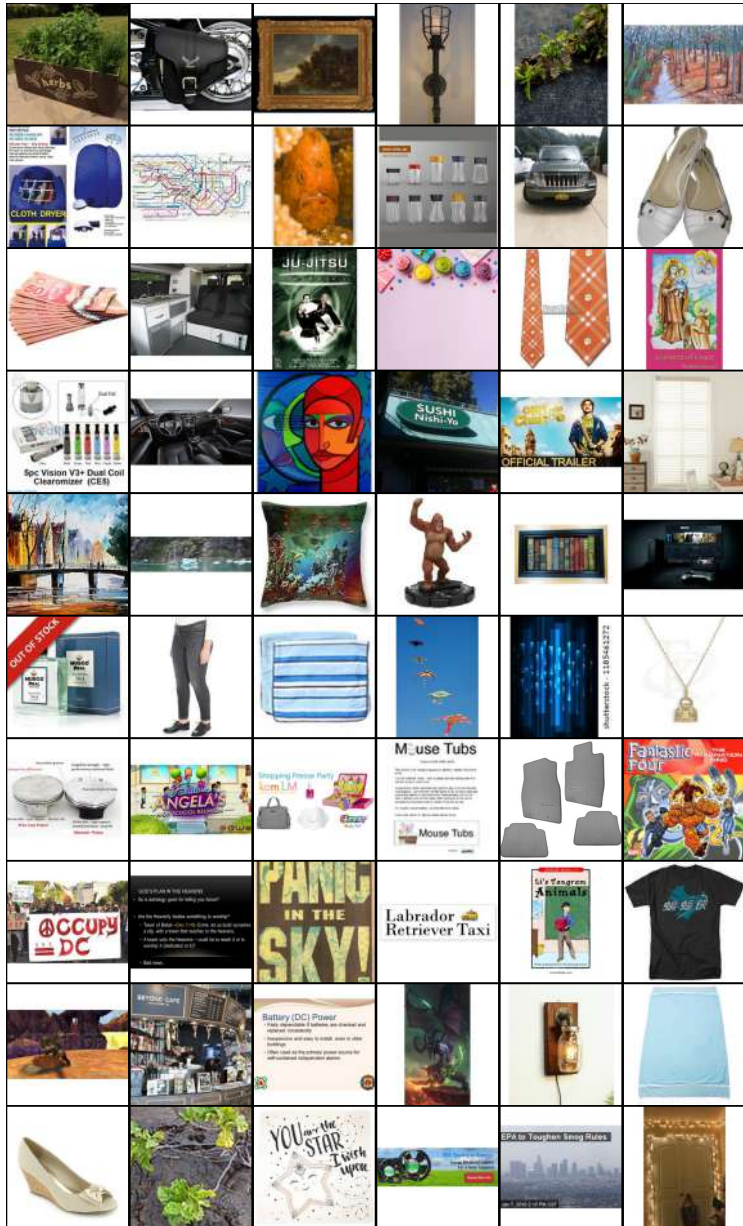


Figure 19: Random samples from LAION-200M. We omit NSFW images and images of humans.

1404  
1405  
1406  
1407  
1408  
1409  
1410  
1411  
1412  
1413  
1414  
1415  
1416  
1417  
1418  
1419  
1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427  
1428  
1429  
1430  
1431  
1432  
1433  
1434  
1435  
1436  
1437  
1438  
1439  
1440  
1441  
1442  
1443  
1444  
1445  
1446  
1447  
1448  
1449  
1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457

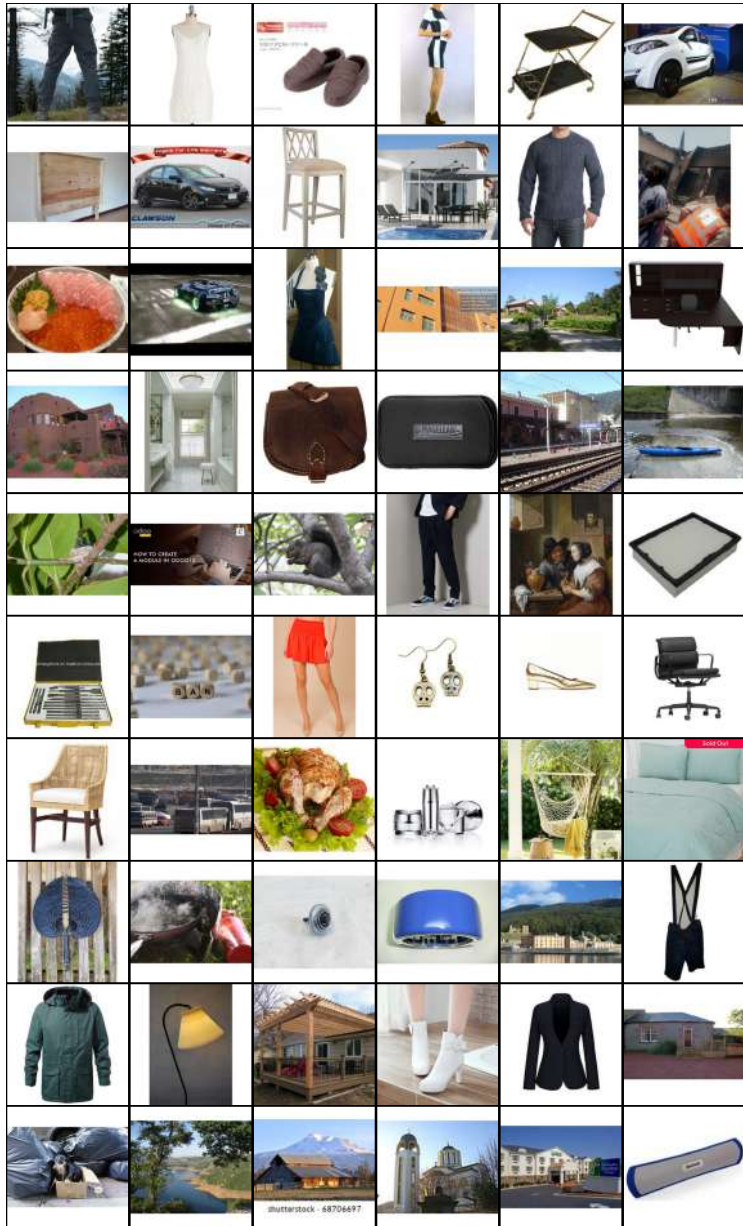


Figure 20: Random samples from LAION-Natural. We omit NSFW images and images of humans.

1458  
1459  
1460  
1461  
1462  
1463  
1464  
1465  
1466  
1467  
1468  
1469  
1470  
1471  
1472  
1473  
1474  
1475  
1476  
1477  
1478  
1479  
1480  
1481  
1482  
1483  
1484  
1485  
1486  
1487  
1488  
1489  
1490  
1491  
1492  
1493  
1494  
1495  
1496  
1497  
1498  
1499  
1500  
1501  
1502  
1503  
1504  
1505  
1506  
1507  
1508  
1509  
1510  
1511



Figure 21: Random samples from LAION-Rendition. We omit NSFW images and images of humans.

1512  
 1513  
 1514  
 1515  
 1516  
 1517  
 1518  
 1519  
 1520  
 1521  
 1522  
 1523  
 1524  
 1525  
 1526  
 1527  
 1528  
 1529  
 1530  
 1531  
 1532  
 1533  
 1534  
 1535  
 1536  
 1537  
 1538  
 1539  
 1540  
 1541  
 1542  
 1543  
 1544  
 1545  
 1546  
 1547  
 1548  
 1549  
 1550  
 1551  
 1552  
 1553  
 1554  
 1555  
 1556  
 1557  
 1558  
 1559  
 1560  
 1561  
 1562  
 1563  
 1564  
 1565

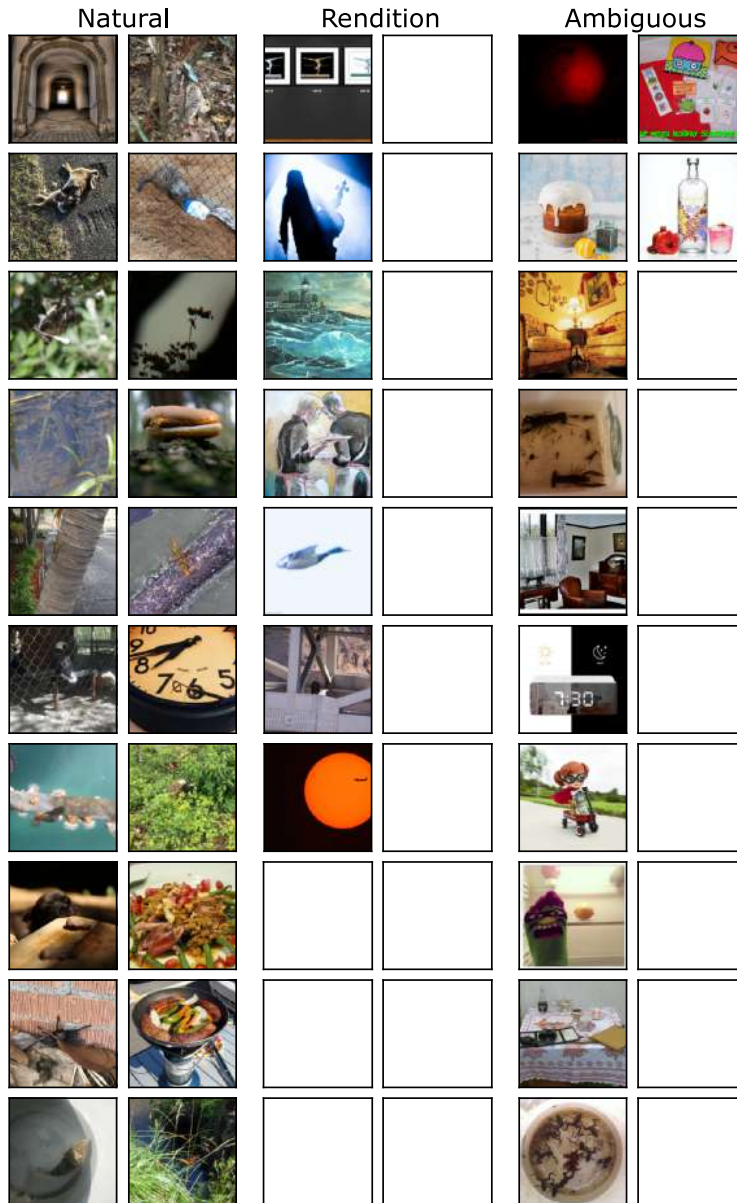


Figure 22: Random samples of ImageNet-A grouped by domain. We omit NSFW images and images of humans.

1566  
1567  
1568  
1569  
1570  
1571  
1572  
1573  
1574  
1575  
1576  
1577  
1578  
1579  
1580  
1581  
1582  
1583  
1584  
1585  
1586  
1587  
1588  
1589  
1590  
1591  
1592  
1593  
1594  
1595  
1596  
1597  
1598  
1599  
1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607  
1608  
1609  
1610  
1611  
1612  
1613  
1614  
1615  
1616  
1617  
1618  
1619

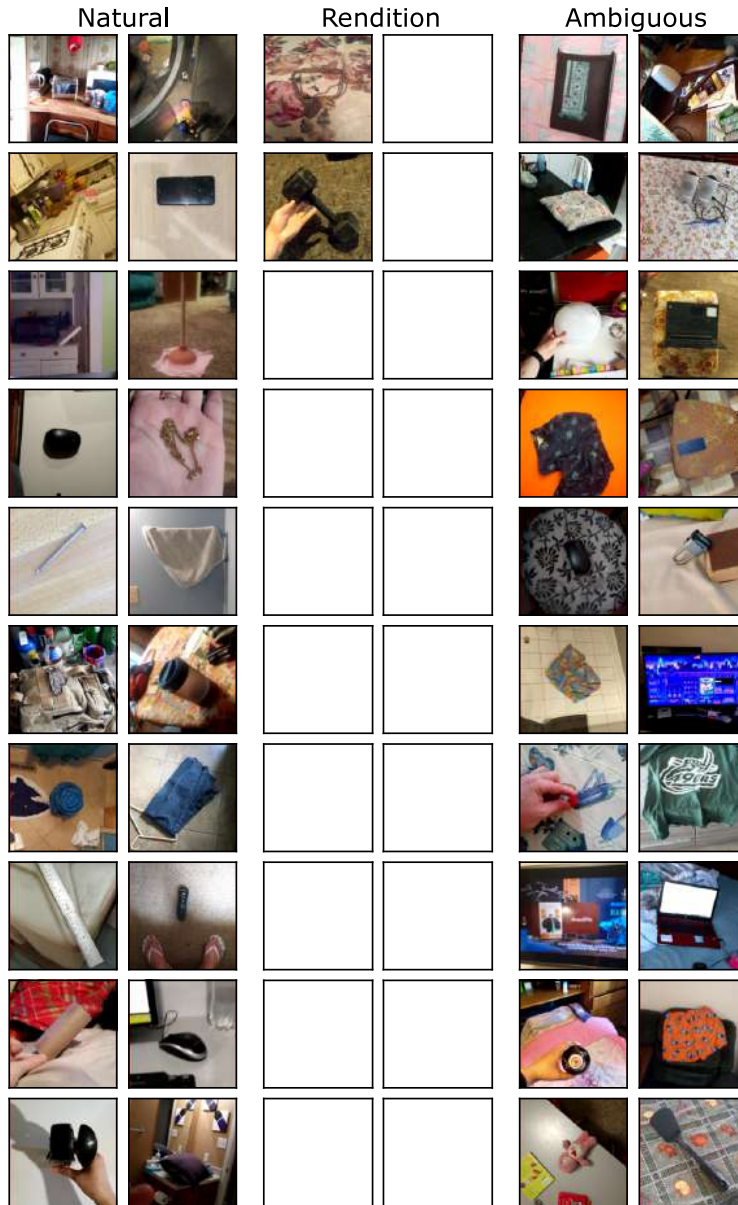


Figure 23: Random samples of ObjectNet grouped by domain. We omit NSFW images and images of humans.

1620  
 1621  
 1622  
 1623  
 1624  
 1625  
 1626  
 1627  
 1628  
 1629  
 1630  
 1631  
 1632  
 1633  
 1634  
 1635  
 1636  
 1637  
 1638  
 1639  
 1640  
 1641  
 1642  
 1643  
 1644  
 1645  
 1646  
 1647  
 1648  
 1649  
 1650  
 1651  
 1652  
 1653  
 1654  
 1655  
 1656  
 1657  
 1658  
 1659  
 1660  
 1661  
 1662  
 1663  
 1664  
 1665  
 1666  
 1667  
 1668  
 1669  
 1670  
 1671  
 1672  
 1673

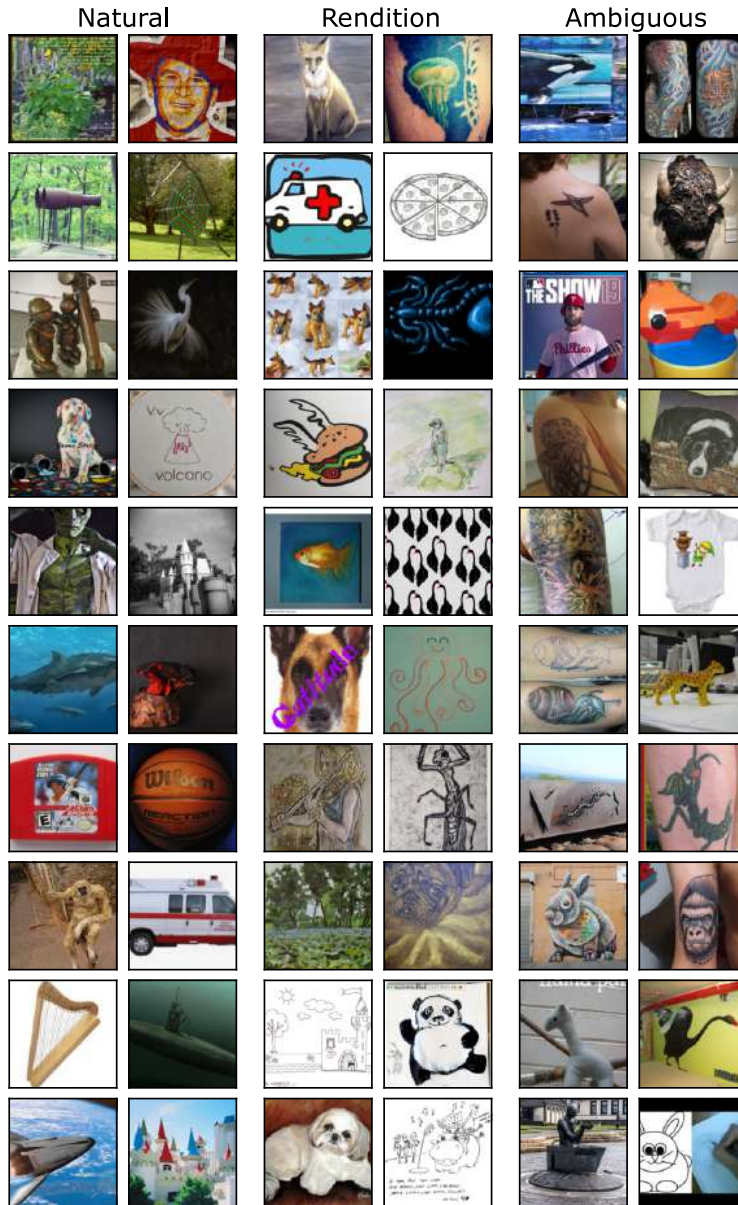


Figure 24: Random samples of ImageNet-R grouped by domain. We omit NSFW images and images of humans.

1674  
 1675  
 1676  
 1677  
 1678  
 1679  
 1680  
 1681  
 1682  
 1683  
 1684  
 1685  
 1686  
 1687  
 1688  
 1689  
 1690  
 1691  
 1692  
 1693  
 1694  
 1695  
 1696  
 1697  
 1698  
 1699  
 1700  
 1701  
 1702  
 1703  
 1704  
 1705  
 1706  
 1707  
 1708  
 1709  
 1710  
 1711  
 1712  
 1713  
 1714  
 1715  
 1716  
 1717  
 1718  
 1719  
 1720  
 1721  
 1722  
 1723  
 1724  
 1725  
 1726  
 1727

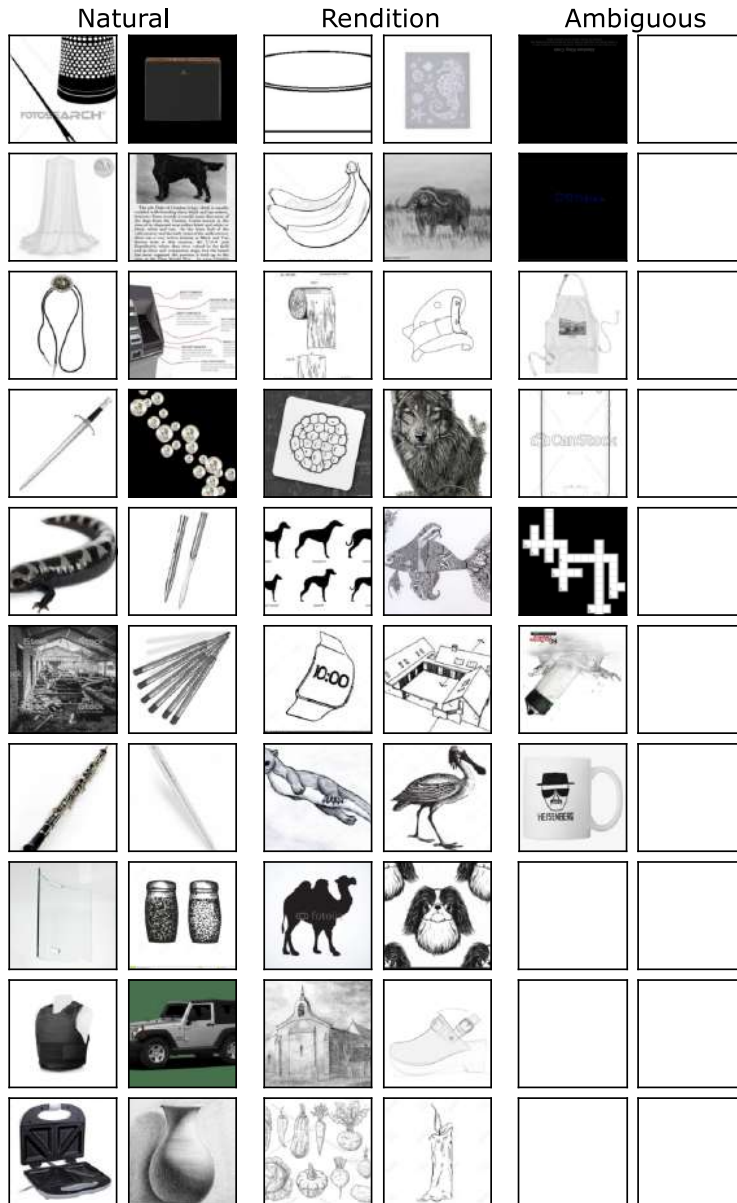


Figure 25: Random samples of ImageNet-Sketch grouped by domain. We omit NSFW images and images of humans.



1728  
1729  
1730  
1731  
1732  
1733  
1734  
1735  
1736  
1737  
1738  
1739  
1740  
1741  
1742  
1743  
1744  
1745  
1746  
1747  
1748  
1749  
1750  
1751  
1752  
1753  
1754  
1755  
1756  
1757  
1758  
1759  
1760  
1761  
1762  
1763  
1764  
1765  
1766  
1767  
1768  
1769  
1770  
1771  
1772  
1773  
1774  
1775  
1776  
1777  
1778  
1779  
1780  
1781

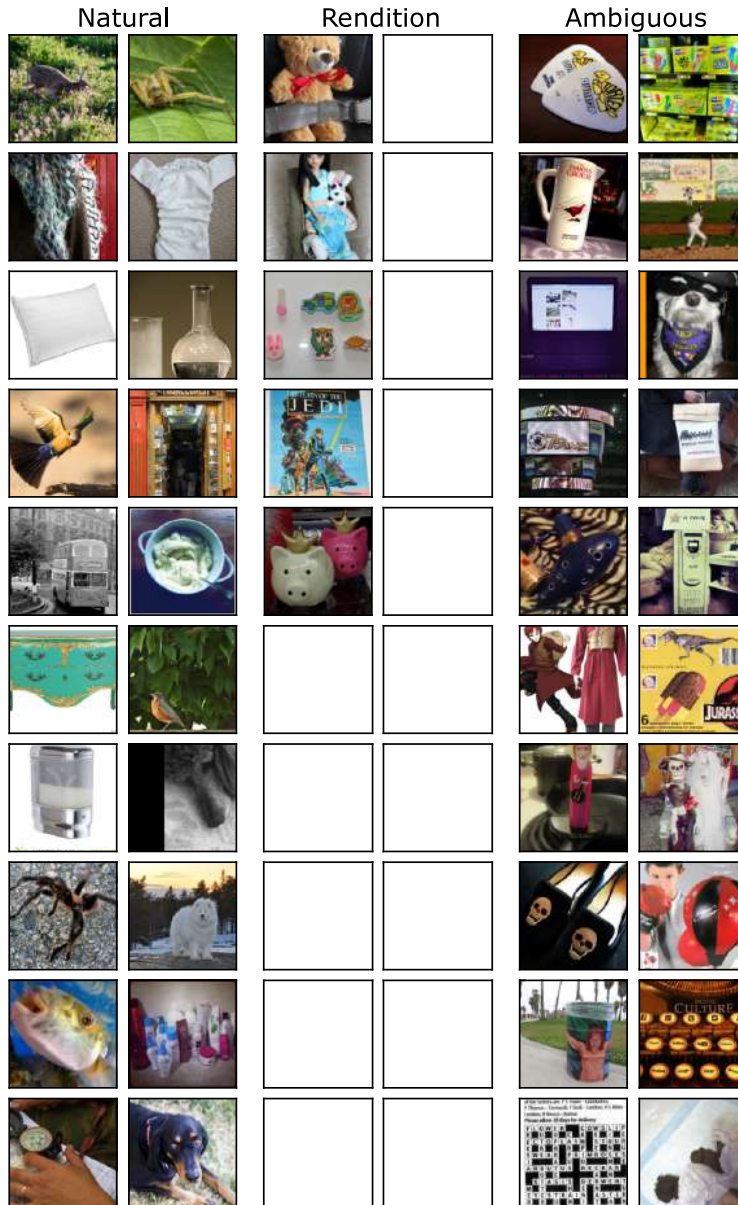


Figure 26: Random samples of ImageNet-V2 grouped by domain. We omit NSFW images and images of humans.



1836  
 1837  
 1838  
 1839  
 1840  
 1841  
 1842  
 1843  
 1844  
 1845  
 1846  
 1847  
 1848  
 1849  
 1850  
 1851  
 1852  
 1853  
 1854  
 1855  
 1856  
 1857  
 1858  
 1859  
 1860  
 1861  
 1862  
 1863  
 1864  
 1865  
 1866  
 1867  
 1868  
 1869  
 1870  
 1871  
 1872  
 1873  
 1874  
 1875  
 1876  
 1877  
 1878  
 1879  
 1880  
 1881  
 1882  
 1883  
 1884  
 1885  
 1886  
 1887  
 1888  
 1889

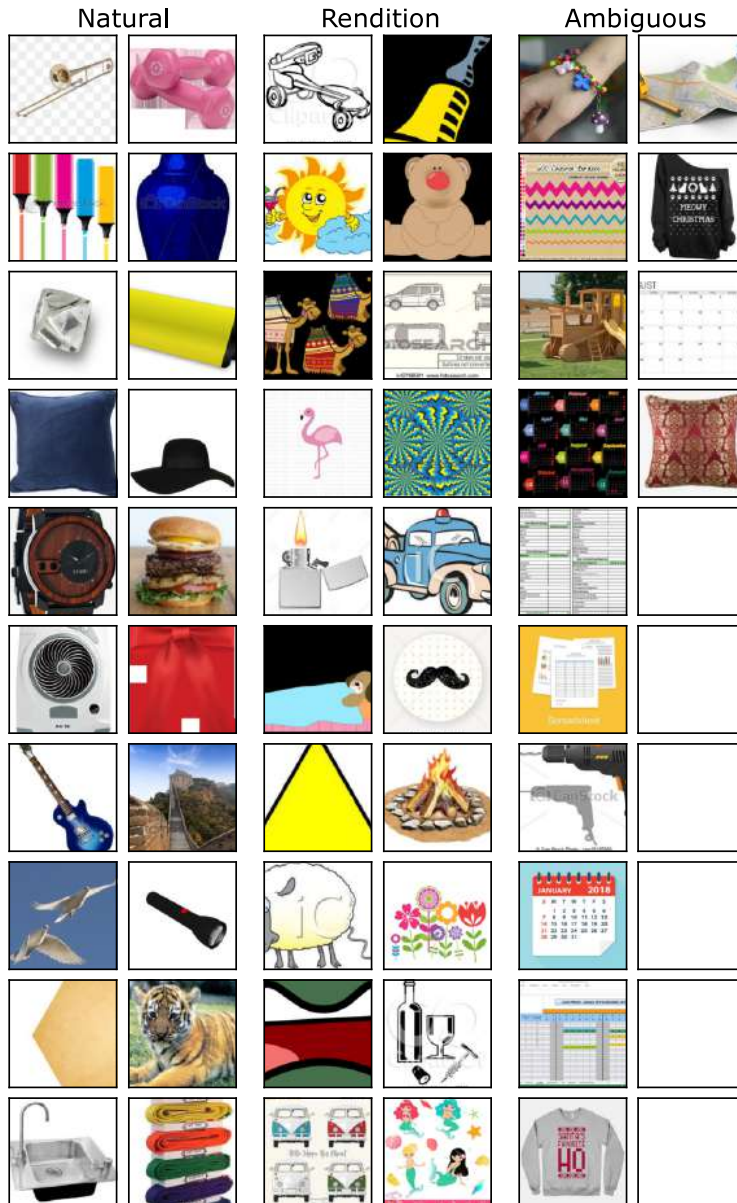


Figure 28: Random samples of DomainNet-Clipart grouped by domain. We omit NSFW images and images of humans.

1890  
 1891  
 1892  
 1893  
 1894  
 1895  
 1896  
 1897  
 1898  
 1899  
 1900  
 1901  
 1902  
 1903  
 1904  
 1905  
 1906  
 1907  
 1908  
 1909  
 1910  
 1911  
 1912  
 1913  
 1914  
 1915  
 1916  
 1917  
 1918  
 1919  
 1920  
 1921  
 1922  
 1923  
 1924  
 1925  
 1926  
 1927  
 1928  
 1929  
 1930  
 1931  
 1932  
 1933  
 1934  
 1935  
 1936  
 1937  
 1938  
 1939  
 1940  
 1941  
 1942  
 1943

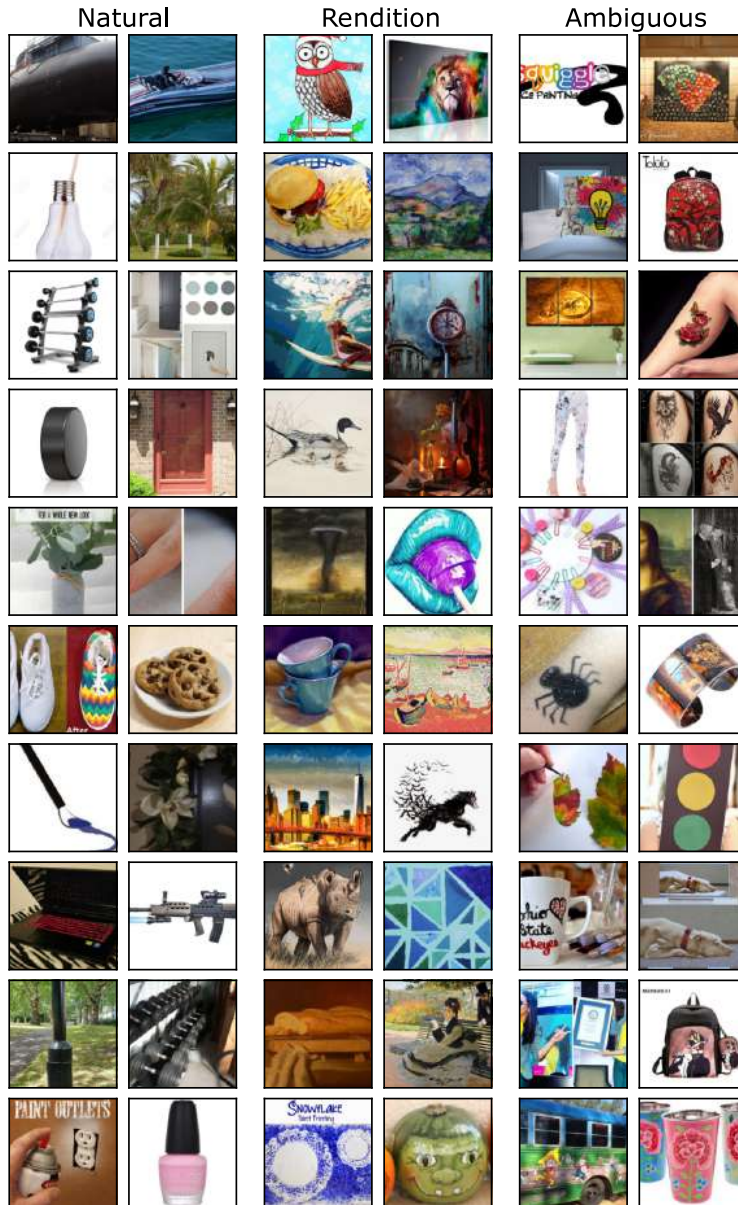


Figure 29: Random samples of DomainNet-Painting grouped by domain. We omit NSFW images and images of humans.

1944  
 1945  
 1946  
 1947  
 1948  
 1949  
 1950  
 1951  
 1952  
 1953  
 1954  
 1955  
 1956  
 1957  
 1958  
 1959  
 1960  
 1961  
 1962  
 1963  
 1964  
 1965  
 1966  
 1967  
 1968  
 1969  
 1970  
 1971  
 1972  
 1973  
 1974  
 1975  
 1976  
 1977  
 1978  
 1979  
 1980  
 1981  
 1982  
 1983  
 1984  
 1985  
 1986  
 1987  
 1988  
 1989  
 1990  
 1991  
 1992  
 1993  
 1994  
 1995  
 1996  
 1997

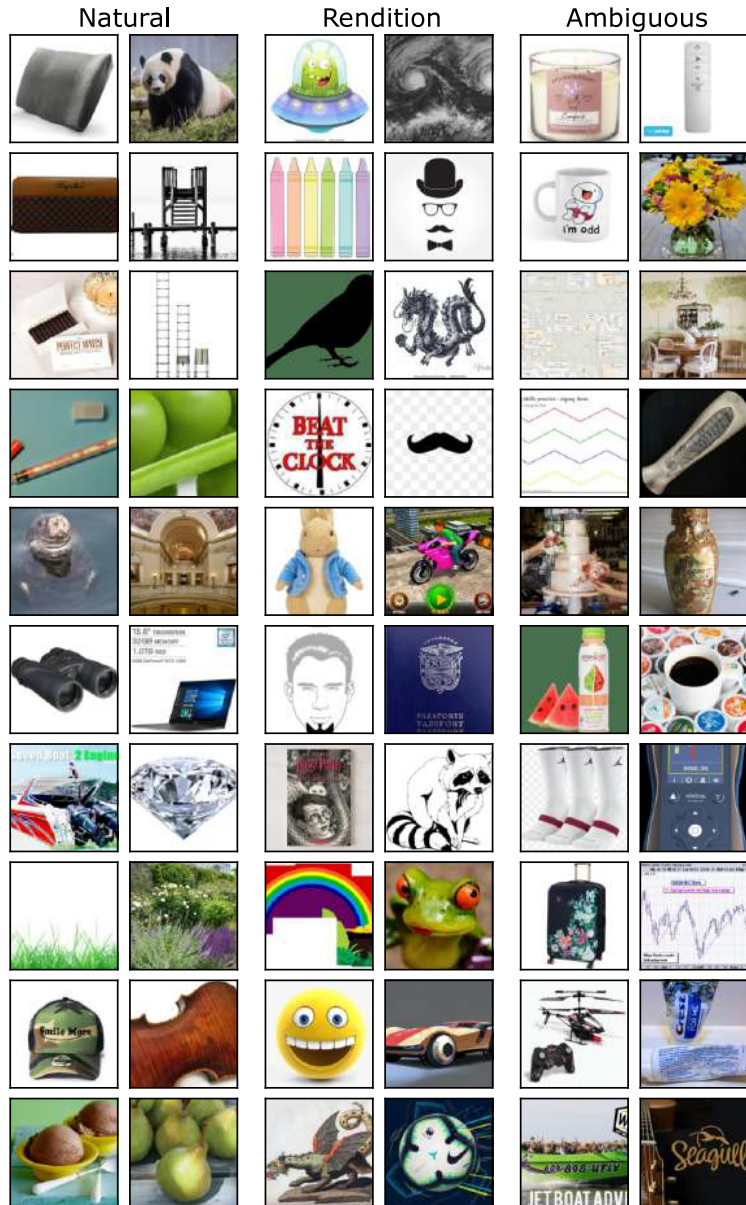


Figure 30: Random samples of DomainNet-Real grouped by domain. We omit NSFW images and images of humans.

1998  
 1999  
 2000  
 2001  
 2002  
 2003  
 2004  
 2005  
 2006  
 2007  
 2008  
 2009  
 2010  
 2011  
 2012  
 2013  
 2014  
 2015  
 2016  
 2017  
 2018  
 2019  
 2020  
 2021  
 2022  
 2023  
 2024  
 2025  
 2026  
 2027  
 2028  
 2029  
 2030  
 2031  
 2032  
 2033  
 2034  
 2035  
 2036  
 2037  
 2038  
 2039  
 2040  
 2041  
 2042  
 2043  
 2044  
 2045  
 2046  
 2047  
 2048  
 2049  
 2050  
 2051

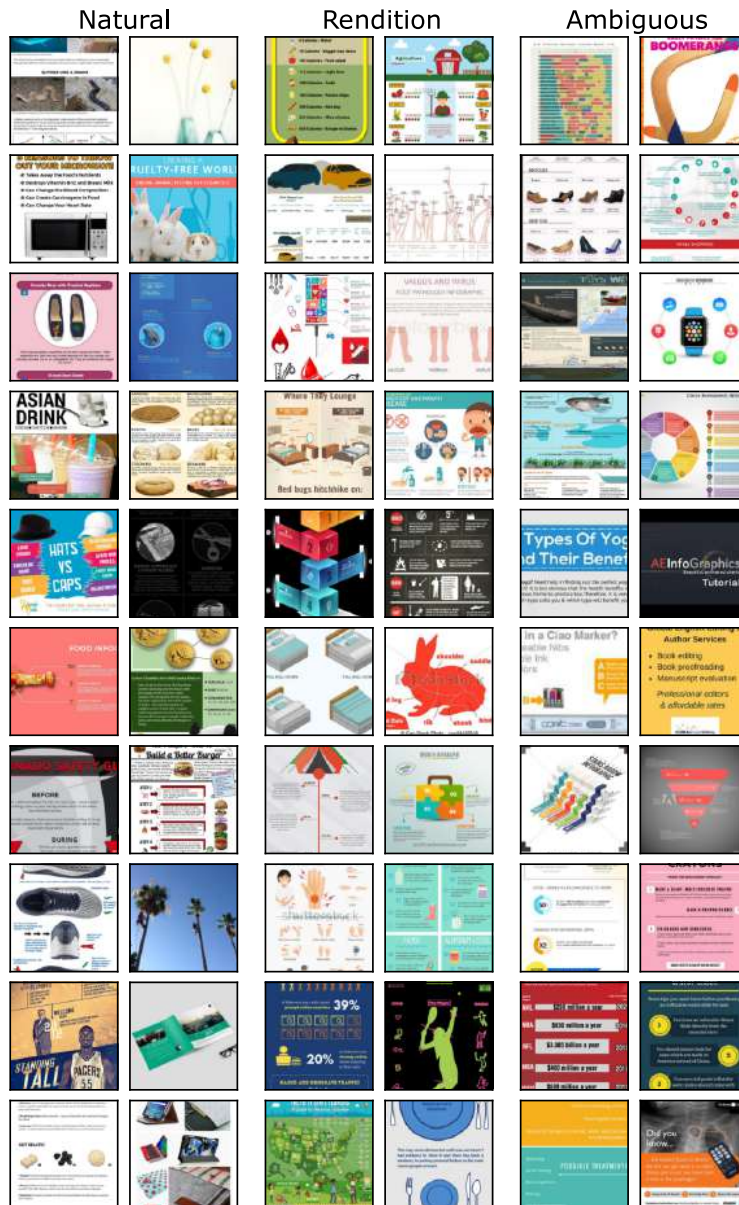


Figure 31: Random samples of DomainNet-Infograph grouped by domain. We omit NSFW images and images of humans.

2052  
 2053  
 2054  
 2055  
 2056  
 2057  
 2058  
 2059  
 2060  
 2061  
 2062  
 2063  
 2064  
 2065  
 2066  
 2067  
 2068  
 2069  
 2070  
 2071  
 2072  
 2073  
 2074  
 2075  
 2076  
 2077  
 2078  
 2079  
 2080  
 2081  
 2082  
 2083  
 2084  
 2085  
 2086  
 2087  
 2088  
 2089  
 2090  
 2091  
 2092  
 2093  
 2094  
 2095  
 2096  
 2097  
 2098  
 2099  
 2100  
 2101  
 2102  
 2103  
 2104  
 2105

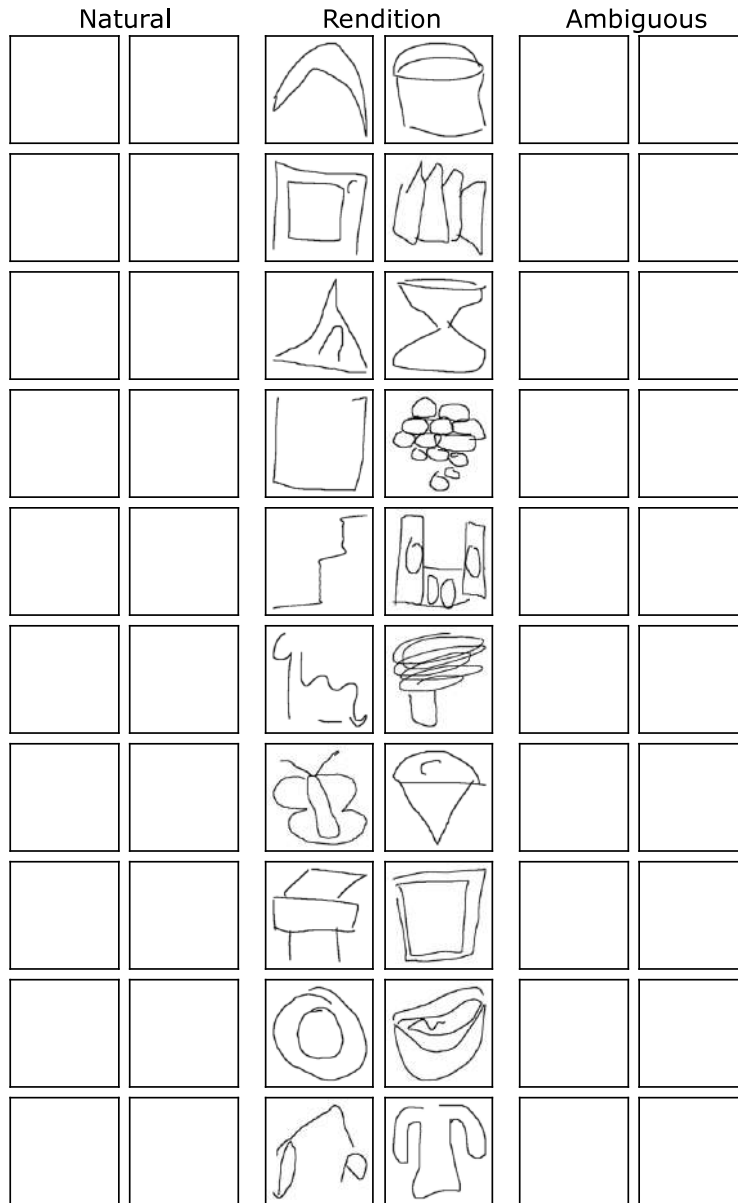


Figure 32: **Random samples of DomainNet-Quickdraw grouped by domain.** We omit NSFW images and images of humans.

2106  
 2107  
 2108  
 2109  
 2110  
 2111  
 2112  
 2113  
 2114  
 2115  
 2116  
 2117  
 2118  
 2119  
 2120  
 2121  
 2122  
 2123  
 2124  
 2125  
 2126  
 2127  
 2128  
 2129  
 2130  
 2131  
 2132  
 2133  
 2134  
 2135  
 2136  
 2137  
 2138  
 2139  
 2140  
 2141  
 2142  
 2143  
 2144  
 2145  
 2146  
 2147  
 2148  
 2149  
 2150  
 2151  
 2152  
 2153  
 2154  
 2155  
 2156  
 2157  
 2158  
 2159

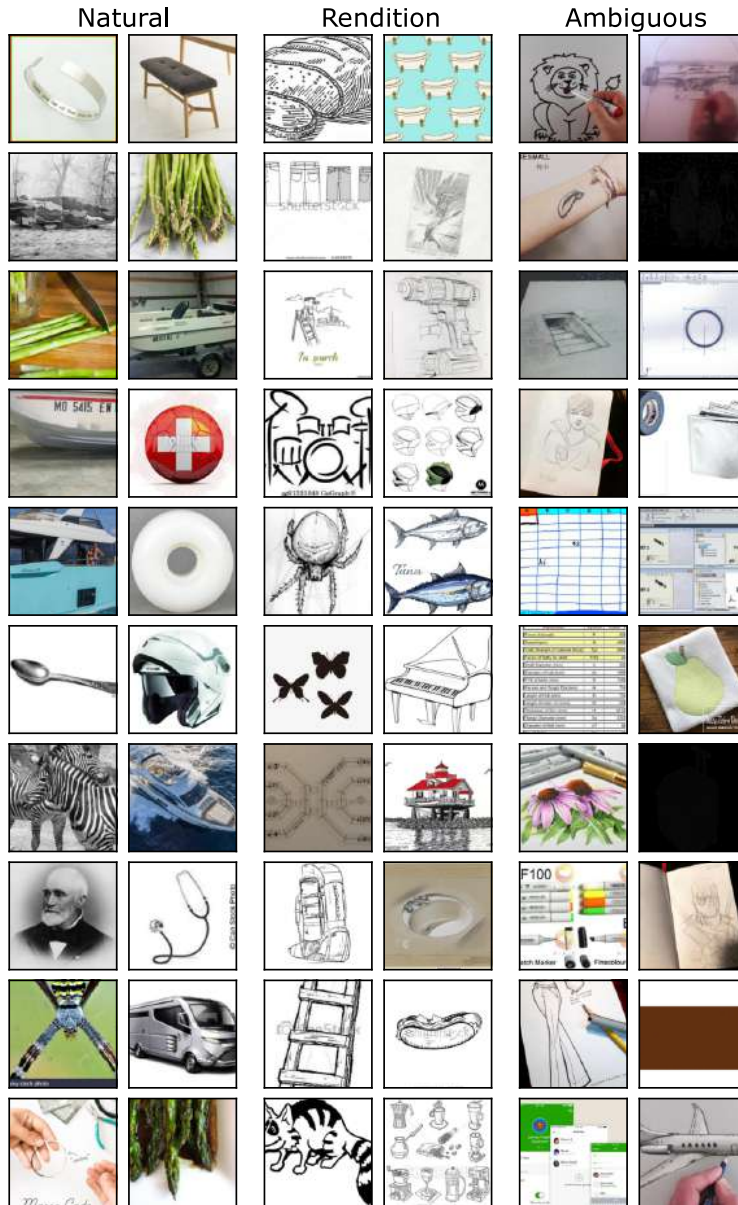


Figure 33: Random samples of DomainNet-Sketch grouped by domain. We omit NSFW images and images of humans.