

WHAT MAKES IMAGENET LOOK UNLIKE LAION

Anonymous authors

Paper under double-blind review

ABSTRACT

ImageNet was famously created from Flickr image search results. What if we recreated ImageNet instead by searching the massive LAION dataset based on image captions alone? In this work, we carry out this counterfactual investigation. We find that the resulting ImageNet recreation, which we call LAIONet, looks distinctly unlike the original. Specifically, the intra-class similarity of images in the original ImageNet is dramatically higher than it is for LAIONet. Consequently, models trained on ImageNet perform significantly worse on LAIONet. We propose a rigorous explanation for the discrepancy in terms of a subtle, yet important, difference in two plausible causal data-generating processes for the respective datasets, that we support with systematic experimentation. In a nutshell, searching based on an image caption alone creates an information bottleneck that mitigates the selection bias otherwise present in image-based filtering. Our explanation formalizes a long-held intuition in the community that ImageNet images are stereotypical, unnatural, and overly simple representations of the class category. At the same time, it provides a simple and actionable takeaway for future dataset creation efforts.

1 INTRODUCTION

For nearly a decade, ImageNet (Deng et al., 2009) was the focal benchmark for much of computer vision and deep learning. Created from image web search results and human filtering, ImageNet contributed curated images suitable for supervised learning at the time. In recent years, however, the community has seen a new generation of models trained on massive amounts of noisy image-text data gathered from the web with minimal curation. Available to the academic public is the massive scale LAION dataset, in two versions, featuring 400 million (Schuhmann et al., 2021) and 5 billion (Schuhmann et al., 2022) crawled image-text pairs, filtered by the OpenAI CLIP model (Radford et al., 2021) for sufficient image-text relevance rather than human annotators.

At the outset, LAION works much like text-based web image search. We can specify a query and retrieve images with high similarity between the query and the text surrounding the image on the website that it was crawled from. We can therefore search LAION for each of the 1000 categories in the ImageNet ILSVRC-2012 dataset¹ and retrieve images corresponding to each of the classes. This process is much like the first step of creating ImageNet from Flickr search results, except that LAION replaces Flickr, but either way, both are based on web crawls. Where the creators of ImageNet hired human annotators to filter images, we analyze image captions to ensure that the resulting images have high fidelity to the class category.

We might expect that for a suitably chosen textual similarity threshold, the resulting dataset would bear resemblance to the original ImageNet. However, we demonstrate that this is anything but the case. The dataset, so created from LAION, very much looks unlike ImageNet. And we explain *why*, supported by independent evidence from other well-curated datasets. This explanation, although subtle, reveals a fundamental fact about the difference between ImageNet and LAION that has consequences for understanding dataset creation at large.

¹Unless otherwise stated, by ImageNet we mean the ImageNet ILSVRC-2012 dataset.

1.1 OUR CONTRIBUTIONS

We introduce a new research artifact, called *LAIONet*, that aims at a recreation of ImageNet on the basis of LAION. We start from LAION-400M, a collection of 400M image-text pairs extracted from web pages in Common Crawl (commoncrawl.org) between 2014 and 2021. The relevance of images and their corresponding texts was quality-controlled with OpenAI CLIP model, excluding instances with a cosine similarity of image and text embeddings less than 0.3.

Creation of LAIONet. We create LAIONet solely on the basis of text-based selection. We require the exact “lemmas” (terms) in a so-called “synset” of an ImageNet category to appear in the text corresponding to an image. Moreover, we require a high similarity between the text and the synset name and definition. We use the cosine similarity of CLIP text embeddings to calculate this similarity, however, we make consistent observations using MPNet (Song et al., 2020) as the text encoder. LAIONet selection criteria are conservative in that they tend toward images that are easy to classify; at least from the CLIP point of view, there is no evidence that LAIONet images are harder to classify than ImageNet.

Contrasting LAIONet and ImageNet. To begin to understand the differences between LAIONet and ImageNet, we evaluate a slew of Imagenet models on LAIONet. As we show, the accuracy of models trained on ImageNet drops by 5 to 12 percentage points when evaluated on LAIONet (Figure 1). In calculating accuracy, we weight classes uniformly as is done in ImageNet. When classes are weighted based on the frequency of each class in LAIONet, accuracy drops by another 5 to 10 percentage points.

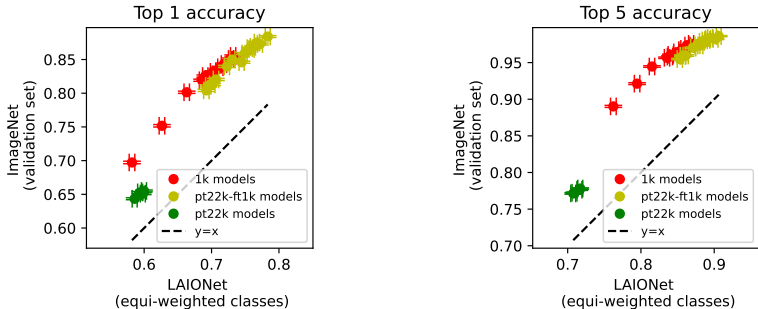


Figure 1: Accuracy of ImageNet-trained models when evaluated on ImageNet validation set versus LAIONet. Three types of models are distinguished based on whether they are pre-trained on ImageNet-22k and whether they are fine-tuned on ImageNet-1k. Accuracy is defined as the average of the recalls calculated for each class that is present in LAIONet.

Drops in accuracy, such as these, are a well-documented phenomenon in machine learning at this point. In this work, we go a step further by providing a substantive explanation for the difference between LAIONet and ImageNet.

Diagnosing the difference. In a first step, we observe that the intra-class similarity, measured as the pairwise similarity of the images within a class, is lower for LAIONet than for ImageNet. In other words, LAIONet images are more diverse in each class. The recall of the models is also lower in the classes with lower intra-class similarity. Hence, lower intra-class similarity gives a concrete reason for why the accuracy of ImageNet models drops on LAIONet. But why does LAIONet have lower intra-class similarity in the first place?

We answer this question in terms of two plausible causal graphs for the respective data-generating processes (Figure 2). Both graphs are based on the standard anti-causal representation of classification problems (Schölkopf et al., 2012), whereby for each category Y there is a mechanism to generate data (here, image X and text T) given Y . But, the graphs differ in one important aspect.

In the case of LAIONet (Figure 2a), selection is based on text alone. The causal graph has the important property that the distribution of the image is independent of the selection decision conditional on the text. In other words the text serves as an information bottleneck between the selection

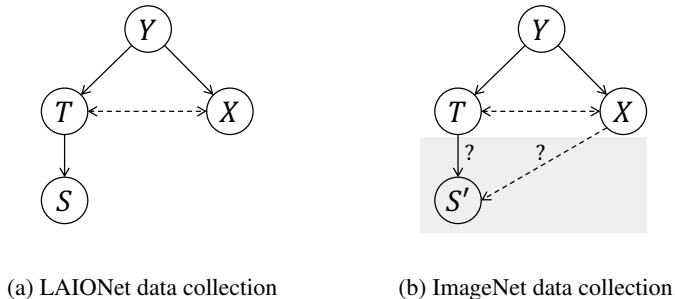


Figure 2: The suggested underlying mechanism of data generation and selection in LAIONet and ImageNet. Class Y , text description T , image X , selection S or S' .

mechanism and the image. Choosing an image reveals nothing more about the image than what can be learned from its textual representation. This powerful conditional independence property limits how much selection can bias the distribution of the image. In contrast, in the case of ImageNet (Figure 2b), there is a link from the image to the selection decision. For example, this link exists when human annotators see the full image and decide to select or discard an image. The existence of this link is what can strongly bias the distribution of the image conditional on selection. It is this selection bias that is visible in the higher intra-class similarity.

Our case hinges on the existence and strength of the image-to-selection link in the causal graph for ImageNet. We then go beyond LAIONet and provide three complementary arguments as evidence:

- We can weaken the image-to-selection link by considering ImageNet instances of different *selection frequencies*. The selection frequency describes the rate at which Amazon MTurk workers selected a candidate image into the dataset within a target class. This allows us to modulate the strength of the image-to-selection link. Looking at three versions of ImageNetV2 (Recht et al., 2019), we find that for a lower selection frequency, the resulting images come closer to LAIONet.
- We show that text alone cannot explain why an image was selected into ImageNet. The ImageNet-Captions dataset (Fang et al., 2022) has restored the captions for one-third of the original ImageNet images. If the text was the only factor in determining the relevance to a synset, it should explain why the images in ImageNet-Captions are selected. Looking at the similarity between texts and their synsets, a majority of text-synset pairs exhibit high similarity, but the distribution has a heavy tail and there are instances with low similarity. For pairs with low similarity, there are often many other synsets more similar to the text. This makes these instances unlikely to have been selected solely based on their text.
- We search LAION for the texts most similar to the texts from the ImageNet-Captions dataset. The resulting images show significantly higher variability (in other words, lower intra-class similarity) than ImageNet. This suggests that another mechanism must have been at play.

In conclusion, we argue that the image-to-selection mechanism was significantly at play in the creation of ImageNet. It is this mechanism that makes ImageNet look unlike LAION. This insight has direct prescriptive value for dataset creation efforts in general. When creating a dataset and diversity is desired, we should select candidates on the basis of an information bottleneck. A succinct text caption, for example, generally carries much less information than the entire image. Selecting on the basis of the text caption, therefore, retains much of the entropy present in the image distribution.

1.2 RELATED WORK

Recreating an ImageNet test set, called ImageNetV2, although with a different motivation, was the subject of the seminal paper by Recht, Roelofs, Schmidt, and Shankar (2019). Engstrom et al. (2020) argue that there is a subtlety in thresholding empirical estimates of the true underlying selection frequency of an image in ImageNetV2. Our argument, however, does not rely on any specific threshold of the selection frequency. We only need to observe what happens as we vary it from small to large. In contrast to ImageNetV2, our goal is not to recreate ImageNet as closely as possible. Rather it is the differences between ImageNet and LAION that are the focus of our investigation.

Many other works have modified ImageNet for a variety of reasons. Geirhos et al. (2019) created a stylized version of ImageNet to reduce the reliance of the trained model on texture. Xiao et al. (2021) disentangled the foreground and background of ImageNet images to show the tendency of the models to rely on the background. Li et al. (2023b) proposed ImageNet-W test set by inserting a transparent watermark into the images of ImageNet validation set, revealing the reliance of the models on watermarks. ImageNet undergoes ongoing augmentation over time. For example, the ImageNet-Captions (Fang et al., 2022) project has restored the captions of about one-third of original ImageNet images from Flickr. ImageNet-X (Idrissi et al., 2023) provides a set of human annotations pinpointing 16 failure types for ImageNet such as pose, background, or lighting. The peculiarities of ImageNet have been the subject of multiple studies. For example, Huh et al. (2016) found the large size and many classes, including very similar classes, do not affect the successful transfer performance of ImageNet-trained features.

On the side of LAION, researchers are keenly interested in understanding the strong zero-shot accuracy of contrastive language image models using LAION (Vogel et al., 2022). Fang et al. (2022) found none of the large training set size, language supervision, and contrastive loss function determines this robustness and a more diverse training distribution should be the main cause. Our work demystifies this distributional advantage by contrasting ImageNet and LAION. Nguyen et al. (2022) compared various large image-text datasets differing in the creation process and found the robustness induced by each varies widely in different aspects, suggesting further studies of the role of dataset design. Our work highlights an important mechanism at play in dataset design that can move the dataset further away from a natural distribution.

2 LAIONET: AN IMAGENET OUT OF LAION

Our starting point is to create an ImageNet-like dataset from LAION. This dataset is a research artifact intended to highlight the differences between LAION and ImageNet. Our goal is not to provide a new benchmark or a new training set. However, LAIONet might be of interest to obtain diverse samples, or variants of LAIONet may be created to improve our understanding of benchmarks.

To start, recall that every ImageNet class corresponds to a WordNet (Miller, 1998) *synset* which consists of so-called *lemmas*. Synsets also come with a short definition known as gloss. We label a LAION instance with a WordNet synset if 1) at least one lemma from the synset exists in the text of the instance, and 2) this text is sufficiently similar to the name and definition of the synset. Out of LAION 400M samples, 21M of them passed the first condition. The second condition ensures the lemma as found in the LAION sample has the intended meaning. To quantify the similarity of the LAION text and a synset, we first create a textual representation for the synset by concatenating its name and definition (to be called the synset text). We then calculate the embedding vectors for both the synset text and LAION text using CLIP and compute their cosine similarity. Alternatively, one may use any sufficiently powerful text encoder for this purpose. For instance, we repeat this process using MPNet (Song et al., 2020) in Appendix A.

Figure 3a illustrates the distribution of LAION text to synset text similarities. In general, a high value for textual similarity ensures the LAION text is describing the same object as the synset. But as Figure 3b shows, we cannot set a very high similarity threshold since the extracted dataset will lose its coverage over the ImageNet’s 1k classes. We found the threshold of 0.82 the highest reasonable choice as it allows for covering most classes while going beyond it sharply reduces the number of covered classes (Figure 3b) with no significant reduction in the dataset size (Figure 3c). To further support this choice, in Section 4 (Figure 8b), we demonstrate that using the restored captions of ImageNet, a textual similarity of above 0.7 is sufficient to ensure that a sample belongs uniquely to the synset. Refer to Appendix C for an example of when the second step of filtering is necessary and why the chosen threshold is conservative.

We take a few additional measures to guarantee the safety and quality of the chosen instances. First, we drop samples with more than one label to simplify evaluation on the dataset. Second, we drop images tagged as not-safe-for-work in LAION. Finally, we exclude images that contain text matching the name of their synset. This will ensure the captions are describing an object in the image and not just reflecting on another text. To achieve this, we employ EAST for text detection (Zhou et al., 2017) and TrOCR for text recognition (Li et al., 2023a). This step eliminates 1.1% of the samples. The final dataset, which we call *LAIONet*, consists of 822k samples from 915 ImageNet

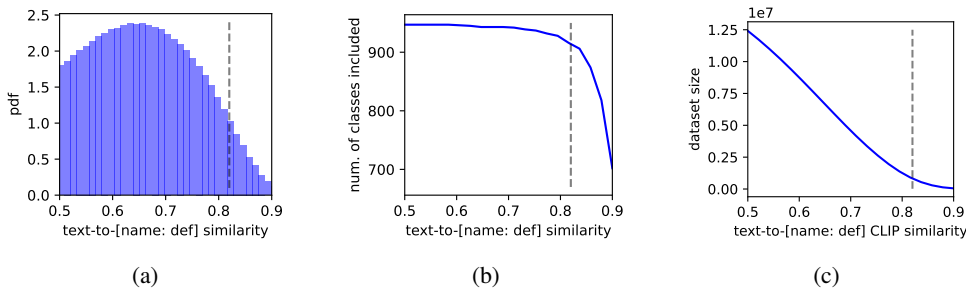


Figure 3: Filtering LAION samples based on their textual similarity to the candidate synsets. The dashed line shows the chosen threshold. (a) The overall distribution of the similarities prior to the second step of filtering. (b and c) The number of ImageNet classes covered by the dataset and the size of the dataset for different levels of similarity threshold.

classes, sufficiently large for fine-grained evaluation purposes at statistical significance. Unlike ImageNet which provides about the same number of images per class, the large variation in the relative frequency of the classes in LAIONet reflects the natural distribution of each class (Figure 4). We will use these frequencies to compare the performance of models in frequent and infrequent classes. We can create a more conservative version of LAIONet mimicking ImageNet validation by retaining only the top 50 most similar instances for each class. This version of LAIONet exhibits the same properties (Appendix B). Find sample images of LAIONet in Appendix G.

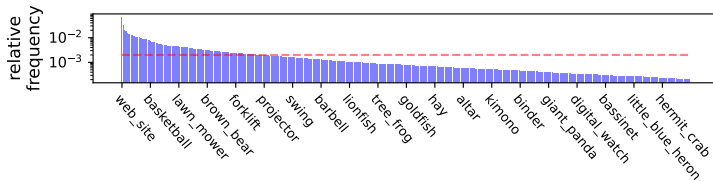


Figure 4: Relative frequencies of different classes in LAIONet sorted in descending order for the 500 most frequent classes. Some class names are shown. The red line shows uniform weight.

Are LAIONet images harder to classify? To find out, we compare CLIP zero-shot accuracy on LAIONet and ImageNet. For every image, we predict the label of the image based on what synset has the highest cosine similarity between the image embedding and the synset text embedding. To make accuracy estimates on LAIONet comparable with ImageNet, we calculate accuracy as the average recall across the classes present in LAIONet. This uniform weighting is consistent with the setup of ImageNet validation with 50 images per class. We found CLIP zero-shot top 1 accuracy to only differ by 2% across datasets. Hence, at least from the CLIP view, LAIONet images are not harder to classify. We note the limitation that CLIP text embeddings are jointly trained with image embeddings, possibly giving CLIP an advantage on LAIONet. Appendix D offers a more direct assessment of the level of difficulty involved in identifying the intended object in LAIONet. This is achieved by directly computing the cross-modality similarity between an image and its associated synset. Overall, LAIONet images do not exhibit significant difficulty compared to ImageNet.

3 LAIONET VERSUS IMAGENET

We begin to understand the differences between the two datasets by looking at the accuracy of various ImageNet classifiers on LAIONet. After observing a significant accuracy drop, we consider the disparity in intra-class similarity as a possible explanation.

3.1 COMPARING ACCURACY

We consider four model families: ResNet (He et al., 2016), Vision Transformers (ViT) (Dosovitskiy et al., 2021), modernized ConvNet (ConvNeXt) (Liu et al., 2022), and Bidirectional Encoder repre-

resentation from Image Transformers (BEiT) (Bao et al., 2022). All models are trained on ImageNet without extra training data. We use various versions of each model in terms of the size (small, base, large, etc.), image resolution (224x224 or 384x384), patch resolution (16x16 or 32x32), and whether models are pre-trained on the complete ImageNet with 22k classes or not. All models come from HuggingFace (huggingface.co) checkpoints.

We first compare the (equally weighted) accuracy defined by the average of recalls across the classes covered by LAIONet. Figure 1 compares the top 1 and top 5 accuracy on ImageNet and LAIONet. In most of the highly accurate models, accuracy drops by at least 10 percentage points when estimated on LAIONet with models pre-trained on ImageNet-22k showing slightly more robustness.

Next, we use the relative frequency of each class in LAIONet to weight its recall and obtain a LAION-weighted accuracy. Figure 5 compares LAION-weighted and equally-weighted accuracy on LAIONet. The LAION-weighted accuracy is consistently lower by 5 to 10 percentage points (similar observations made in Appendix H when evaluated on ImageNet). This can partially be explained by the observation that ImageNet-trained models are performing worse when the class is describing a more common object (Appendix F.1).

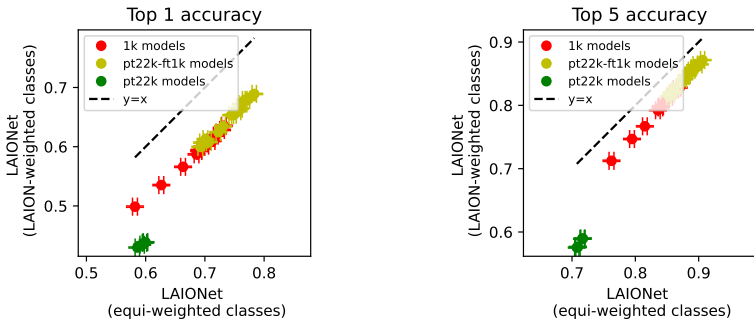


Figure 5: A LAION-weighted accuracy is calculated according to the relative frequency of the classes in LAIONet and compared to the accuracy with equally weighted classes.

3.2 COMPARING INTRA-CLASS SIMILARITY

While LAIONet images are in a precise sense not more difficult than ImageNet, there is another factor that can explain the accuracy drop: the intra-class similarity of images. We define this similarity as the pairwise similarity of the images from the same class, measured by the cosine similarity of their CLIP image embeddings. The lower these similarity values, the more diverse the images from that class. Figure 6a shows the distribution of intra-class similarities aggregated over all the classes. To make the distributions comparable, we sampled (with replacement) the similarities from LAIONet to match ImageNet. The left tail of the LAIONet intra-class similarity distribution makes it clear that LAIONet overall provides a more diverse set of images. To observe the effect in greater detail, for each class, Figure 6b shows the average intra-class similarity of LAIONet images subtracted by the average intra-class similarity of ImageNet images from the same class. In almost two-thirds of the classes, LAIONet has significantly lower intra-class similarity. This provides further evidence that LAIONet images exhibit greater variability within each class.

In Appendix F.2, we show that models struggle more with classes where LAIONet and ImageNet have significantly different intra-class similarity. This, combined with our observation of LAIONet having lower intra-class similarity, supports our argument that intra-class similarity plays a crucial role in reducing accuracy.

4 DIAGNOSING IMAGENET

As is standard modeling practice, we think of a data-generating process that for a given class Y generates a pair of image X and text T . Ideally, when we search for images of a particular class y , we would like to draw samples from distribution $p(X|Y = y)$. Unless we have access to the

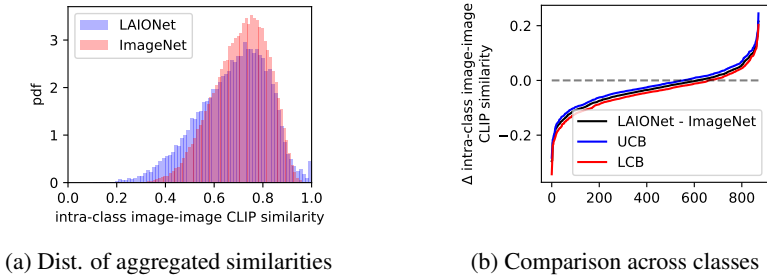


Figure 6: Comparing the intra-class similarity of LAIONet and ImageNet. (a) In each class, pairwise similarities of LAIONet images are sampled to match ImageNet in number. All the classes combined, the distribution of intra-class similarity is depicted. (b) For each class, the average intra-class similarity of ImageNet images is subtracted from the same value in LAIONet. The blue and red curves show upper and lower 95% confidence intervals. All values are sorted ascendingly.

generative process or we have a completely random set of images all correctly labeled, drawing samples directly from $p(X|Y = y)$ will not be possible. In particular, none of these options are available when researchers collect a new dataset. Instead, researchers have to define a selection mechanism S for choosing images. What we observe is the conditional distribution of X given S .

In creating LAIONet, we relied on texts to select the samples (Figure 2a). LAIONet images follow $p(X|S = 1)$, where $S = 1$ if T is sufficiently similar to Y . With our conservative selection criteria, we can assume every T passed our similarity threshold is generated from the intended $Y = y$. Therefore, $p(X|S = 1) = p(X|S = 1, Y = y)$. Generally, an image carries much more information than the text. So, for the images of a certain class, conditioning on the text alone should not alter the distribution significantly. Intuitively speaking, $p(X|Y = y, T = t) \approx p(X|Y = y)$. In our setting, a weaker independence is sufficient to show LAIONet images follow the desired distribution. Even if information from X beyond Y is present in T , since we deliberately refrained from searching for visual descriptions in the text, we expect S to be independent from X for a given $Y = y$. Therefore, we have reason to hope $p(X|S = 1) \approx p(X|S = 1, Y = y) \approx p(X|Y = y)$.

In general, a selection S' can rely on both text and image directly (Figure 2b). In this case, the distribution of observed images $p(X|S' = 1)$ can be far from the desired distribution $p(X|Y = y)$. We believe this has happened in the collection of ImageNet, primarily through human annotators examining and acting on images. Incorporation of visual features at the side of the search engine provider is another plausible mechanism. While we may not be able to pinpoint the exact mechanism at play, we will now move beyond LAIONet and demonstrate, through three independent experiments, a strong link between the image X and the selection criterion S' in the creation of ImageNet.

4.1 A WEAKER IMAGE-TO-SELECTION LINK MAKES IMAGENET MORE LIKE LAIONET

Image annotation is one clear mechanism by which the image X influences selection S' . Changing the strictness of annotation allows us to modulate the strength of this mechanism and measure its effect. This experiment is possible due to the availability of ImageNetV2 (Recht et al., 2019) that comes with three different versions. The three versions of ImageNetV2, called a, b, and c, differ in the level of agreement among annotators. More precisely, each image comes with a *MTurk selection frequency* which is what fraction of MTurk workers selected the image to be from the target class. ImageNetV2 versions a, b, and c have an average MTurk selection frequency of 0.85, 0.73, and 0.93, respectively. Note that version b has the lowest and version c has the highest selection frequency.

We first observe that allowing for more disagreement among annotators results in the inclusion of more diverse images. Figure 7a shows the distribution of intra-class similarity for ImageNetV2 versions b and c. One can see that in version b with the lowest average MTurk selection frequency, the intra-class similarity is shifted toward lower values. We further show as the average MTurk selection frequency increases, ImageNetV2 becomes more similar to ImageNet and less similar to LAIONet. In this regard, to compare two datasets, we count the number of classes in which the first dataset has significantly lower intra-class similarity than the second dataset, and vice versa. Figure 7b compares LAIONet and three versions of ImageNetV2. As the figure suggests, LAIONet

and ImageNetV2 are quite similar when the average MTurk selection frequency is low (ImageNetV2 version b) but as the MTurk selection frequency increases, ImageNetV2 shows higher intra-class similarity than LAIONet. At the same time, Figure 7c shows ImageNetV2 becomes more similar to ImageNet as we increase the MTurk selection frequency. These observations show the impact the image has on the selection, particularly during annotation, is significant and can partially explain the divergence between LAIONet and ImageNet. Further, the extra intra-class diversity of LAIONet is achievable from less stringent human annotation and can explain the consistent accuracy drop on LAIONet and ImageNetV2.

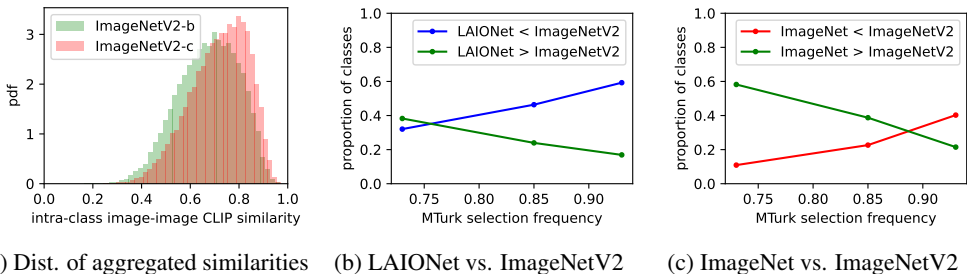


Figure 7: The effect of MTurk selection frequency on intra-class similarity. (a) The distribution of intra-class similarity aggregated over all classes for ImageNetV2 versions b and c. (b) LAIONet versus three versions of ImageNetV2. Vertical axis shows the proportion of classes in which one dataset exhibits significantly lower intra-class similarity than the other dataset (significance determined using 95% confidence intervals). Blue curve: LAIONet has lower intra-class similarity. Green curve: ImageNetV2 has lower intra-class similarity. (c) ImageNet versus ImageNetV2. Red curve: ImageNet has lower intra-class similarity. Green curve: ImageNetV2 has lower intra-class similarity.

4.2 TEXT ALONE CANNOT EXPLAIN WHY AN IMAGE IS SELECTED INTO IMAGE NET

ImageNet-Captions (Fang et al., 2022) is a subset of ImageNet-1k training data with restored title, description, and tags from Flickr. We assume the samples in ImageNet-Captions are a random subset of the original ImageNet and the captions are accurately restored. If there was no link $X \rightarrow S'$, the accompanying caption of an image in ImageNet-Captions should be able to explain why this image is selected. We follow Fang et al. (2022) and define the text as the title, description, and tags concatenated. Figure 8a illustrates the similarity between the texts and their respective synsets using CLIP text embeddings. Although most of the texts have a high similarity of 0.6 or above to their synsets, the distribution has a heavy left tail. The fact that a text has low similarity to the intended synset does not necessarily mean it could not be chosen by the search engine. However, we show many of the texts that have low similarity to the intended synsets actually have high similarity to numerous other synsets, making them less likely to have appeared for the intended meaning. For every text, we find the similarity to all synsets, i.e. the similarity to their names and definitions, and count the proportion of unintended synsets (false classes) that are more similar to the text than the intended synset. A low value for this proportion shows the text well represents its intended synset whereas a significant non-zero value indicates that there are considerable other synsets that are more strongly present in the text. As Figure 8b demonstrates, for a text with low similarity to its synset there are on average 5% (equivalently, 200) or more other synsets more similar to the text. These observations show that at least based on the restored texts in ImageNet-Captions, the text alone cannot fully explain why an image is selected and another mechanism should have been at play.

4.3 IMAGE NET, HAD IT BEEN CREATED SOLELY SEARCHING TEXTS, DOES NOT RESEMBLE CURRENT IMAGE NET

If the link from X to S' did not exist, regardless of how the selection algorithms works, $p(X|T = t)$ would look similar in both graphs of Figure 2. To test this hypothesis, we extract a new dataset from LAION. For every image in ImageNet with corresponding text $T = t$ in ImageNet-Captions, we find the LAION sample with the most similar text to t . We only keep a LAION sample if the similarity is above 0.7. This choice ensures the two texts are sufficiently similar as we can consider them roughly the same while the dataset covers more than 95% of the ImageNet classes (Appendix E).

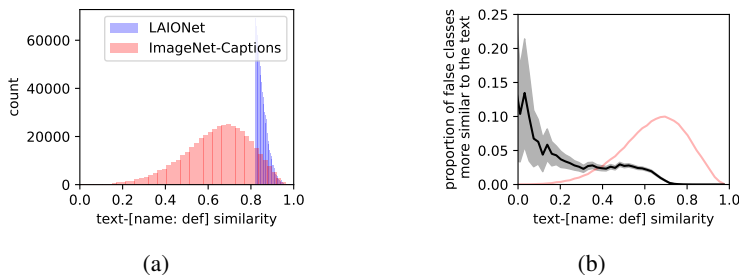


Figure 8: (a) The distribution of the text-to-synset similarity. (b) For every bin of text-to-synset similarity, the average proportion of unintended classes which are more similar to the text than the intended class is depicted in black.

As Figure 9a suggests, images in the new dataset have a significantly lower intra-class similarity. Looking at each class separately, Figure 9b shows in almost 70% of the classes, the images from the new dataset are significantly more diverse (have lower intra-class similarity). These observations reject the hypothesis that the graphs of Figure 2 have the same structure and show a potential leak from the image to the selection. We note the limitation that texts in the ImageNet-Captions dataset may not completely include the text available at the time of ImageNet creation. Second, for many cases, we were unable to find great matches for the ImageNet texts in LAION-400M and scaling our analysis to LAION-5B might help here.

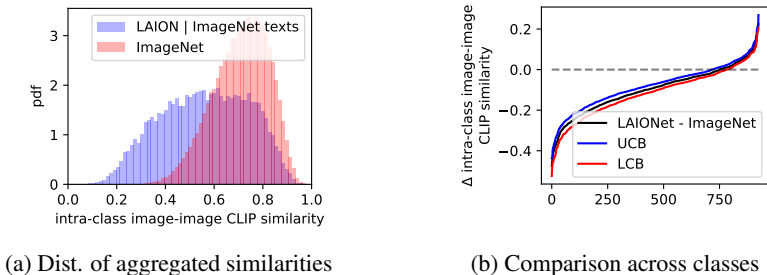


Figure 9: Comparing the intra-class similarity of the new dataset and ImageNet. The new dataset is obtained by selecting LAION examples with the most similar texts to the texts in ImageNet-Captions. (a) Distribution of intra-class similarity aggregated across all classes. In each class, pairwise similarities of the images in the new dataset are sampled to match ImageNet in number to make the distributions comparable. (b) For each class, the average of the intra-class similarity of the images in the new dataset minus the corresponding value in ImageNet is plotted in black. The upper and lower 95% confidence bounds are depicted in blue and red. All values are sorted ascendingly.

5 CONCLUSION

In conclusion, we argue that the image-to-selection mechanism played a significant role in the creation of ImageNet, distinguishing it from LAION. We demonstrated this through three experiments. First, we modulated the speculated link from image to selection, showing the significant contribution this mechanism has in reducing the diversity of the selected images. The next two experiments rejected the hypothesis that image plays no or negligible role in the selection by showing ImageNet captions cannot solely explain the selection.

This insight carries valuable implications for dataset creation efforts in general. When developing a new dataset and diversity is desired, we advise selecting candidate instances based on an information bottleneck, like a succinct textual description of the instance, rather than the full instance. This will mitigate the selection bias that may otherwise distort the distribution of data conditional on selection.

REFERENCES

- Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. BEit: BERT pre-training of image transformers. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=p-BhZSz59o4>.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=YicbFdNTTy>.
- Logan Engstrom, Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Jacob Steinhardt, and Aleksander Madry. Identifying statistical bias in dataset replication. In *International Conference on Machine Learning*, pp. 2922–2932. PMLR, 2020.
- Alex Fang, Gabriel Ilharco, Mitchell Wortsman, Yuhao Wan, Vaishaal Shankar, Achal Dave, and Ludwig Schmidt. Data determines distributional robustness in contrastive language image pre-training (clip). In *International Conference on Machine Learning*, pp. 6216–6234. PMLR, 2022.
- Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A. Wichmann, and Wieland Brendel. Imagenet-trained CNNs are biased towards texture; increasing shape bias improves accuracy and robustness. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=Bygh9j09KX>.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Minyoung Huh, Pulkit Agrawal, and Alexei A Efros. What makes imagenet good for transfer learning? *arXiv preprint arXiv:1608.08614*, 2016.
- Badr Youbi Idrissi, Diane Bouchacourt, Randall Balestriero, Ivan Evtimov, Caner Hazirbas, Nicolas Ballas, Pascal Vincent, Michal Drozdal, David Lopez-Paz, and Mark Ibrahim. Imagenet-x: Understanding model mistakes with factor of variation annotations. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=HXz7Vcm3VgM>.
- Minghao Li, Tengchao Lv, Jingye Chen, Lei Cui, Yijuan Lu, Dinei Florencio, Cha Zhang, Zhoujun Li, and Furu Wei. Trocr: Transformer-based optical character recognition with pre-trained models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 13094–13102, 2023a.
- Zhiheng Li, Ivan Evtimov, Albert Gordo, Caner Hazirbas, Tal Hassner, Cristian Canton Ferrer, Chenliang Xu, and Mark Ibrahim. A whac-a-mole dilemma: Shortcuts come in multiples where mitigating one amplifies others. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 20071–20082, 2023b.
- Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11976–11986, 2022.
- George A Miller. *WordNet: An electronic lexical database*. MIT press, 1998.
- Thao Nguyen, Gabriel Ilharco, Mitchell Wortsman, Sewoong Oh, and Ludwig Schmidt. Quality not quantity: On the interaction between dataset design and robustness of CLIP. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=LTCBavFWp5C>.

- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.
- Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do imagenet classifiers generalize to imagenet? In *International conference on machine learning*, pp. 5389–5400. PMLR, 2019.
- Bernhard Schölkopf, Dominik Janzing, Jonas Peters, Eleni Sgouritsa, Kun Zhang, and Joris Mooij. On causal and anticausal learning. In *International Conference on International Conference on Machine Learning*, pp. 459–466, 2012.
- Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114*, 2021.
- Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade W Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa R Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. LAION-5b: An open large-scale dataset for training next generation image-text models. In *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2022. URL <https://openreview.net/forum?id=M3Y74vmsMcY>.
- Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. Mpnet: Masked and permuted pre-training for language understanding. *Advances in Neural Information Processing Systems*, 33: 16857–16867, 2020.
- Felix Vogel, Nina Shvetsova, Leonid Karlinsky, and Hilde Kuehne. V1-taboo: An analysis of attribute-based zero-shot capabilities of vision-language models. *CoRR*, abs/2209.06103, 2022. URL <https://doi.org/10.48550/arXiv.2209.06103>.
- Kai Yuanqing Xiao, Logan Engstrom, Andrew Ilyas, and Aleksander Madry. Noise or signal: The role of image backgrounds in object recognition. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=g13D-xY7wLq>.
- Xinyu Zhou, Cong Yao, He Wen, Yuzhi Wang, Shuchang Zhou, Weiran He, and Jiajun Liang. East: an efficient and accurate scene text detector. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pp. 5551–5560, 2017.

A AN MPNET-FILTERED LAIONET

The creation of LAIONet relies on textual similarity of the LAION text and synset text. In Section 2 we used the cosine similarity of CLIP text embeddings to calculate this similarity, however, any sufficiently strong text encoder can be used for this purpose. In particular, we use MPNet (Song et al., 2020) fine-tuned on 1B sentence pairs with a contrastive objective by HuggingFace.² We follow a similar procedure to Section 2 and choose the maximum similarity threshold so that the resulting dataset does not lose its coverage over classes. We select the similarity threshold of 0.58. As Figure 10 suggests, a threshold larger than 0.58 may exclude many classes without reducing the size of the resulting dataset. Refer to Appendix C for additional evidence that this threshold works.

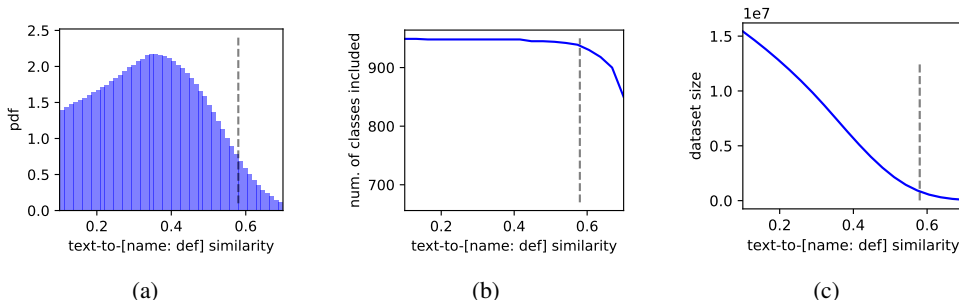


Figure 10: Filtering LAION samples based on their MPNet textual similarity to the candidate synsets. The dashed line shows the chosen threshold. (a) The overall distribution of the similarities prior to the second step of filtering. (b and c) The number of ImageNet classes covered by the dataset and the size of the dataset for different levels of similarity threshold.

Proceeding with the similarity threshold of 0.58, and after dropping samples labeled as not-safe-for-work, samples with multiple labels, and images containing text of their associated synsets, this version of LAIONet will have 831k samples covering 938 classes. As Figure 11 shows, consistent with our observation from CLIP-filtered LAIONet, models trained on ImageNet experience 10 to 15 percentage points of accuracy drop on MPNet-filtered LAIONet.

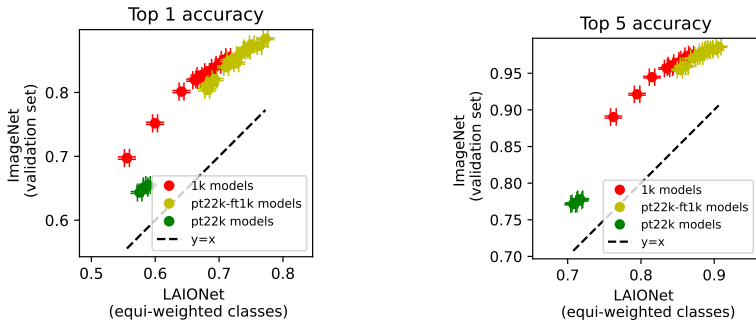


Figure 11: Accuracy of ImageNet-trained models when evaluated on ImageNet validation set versus MPNet-filtered LAIONet. Three types of models are distinguished based on whether they are pre-trained on ImageNet-22k and whether they are fine-tuned on ImageNet-1k. Accuracy is defined as the average of the recalls calculated for each class that is present in LAIONet.

Last but not least, Figure 12 suggests that MPNet-filtered LAIONet also exhibits lower intra-class similarity compared to ImageNet. In particular, in more than 70% of the classes, LAIONet has significantly lower intra-class similarity than ImageNet.

²<https://huggingface.co/sentence-transformers/all-mpnet-base-v2>

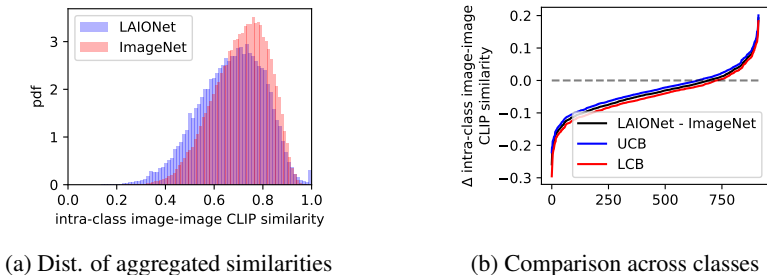


Figure 12: Comparing the intra-class similarity of LAIONet and ImageNet. (a) In each class, pairwise similarities of LAIONet images are sampled to match ImageNet in number. All the classes combined, the distribution of intra-class similarity is depicted. (b) For each class, the average intra-class similarity of ImageNet images is subtracted from the same value in LAIONet. The blue and red curves show upper and lower 95% confidence intervals. All values are sorted ascendingly.

B A LAIONET FROM MOST SIMILARS

We created LAIONet by ensuring the presence of at least one lemma from the associated synset in the LAION text and by ensuring sufficient similarity between the synset text and LAION text. The frequency of each class in LAIONet reflects the natural distribution of that class on the web and likely worldwide. However, we can create a more conservative version of LAIONet by retaining only the top 50 most similar instances for each class. This will make LAIONet more similar to the ImageNet validation set. Such a version of LAIONet will have 39k samples covering 915 classes if initially filtered by CLIP similarity threshold of 0.82, and 41k samples covering 938 classes if initially filtered by MPNet similarity of 0.58.

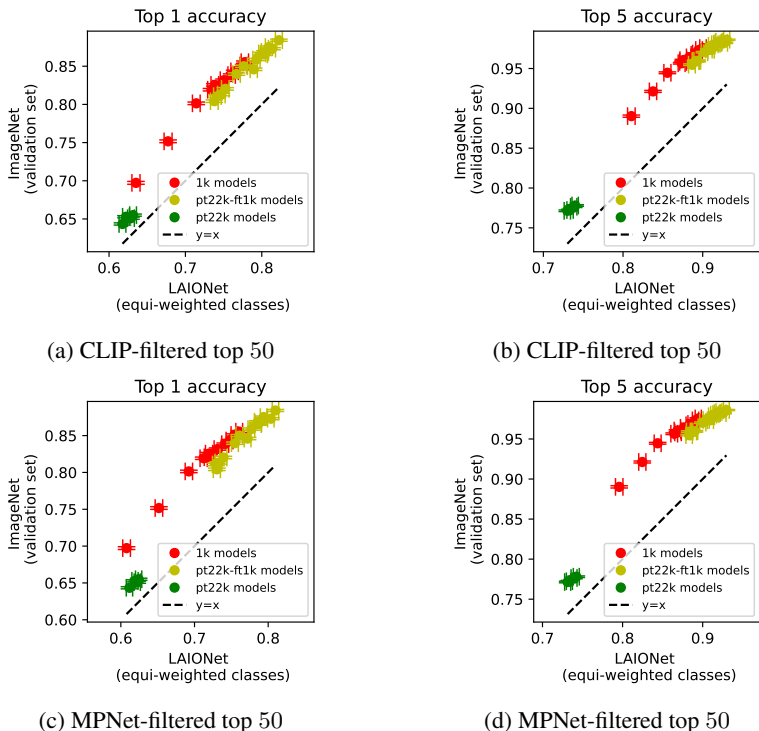


Figure 13: Accuracy of ImageNet-trained models when evaluated on ImageNet validation set versus LAIONet created by retaining top 50 most similar instances for each class.

Figure 13 illustrates that models performing well on ImageNet consistently experience a 7 to 10 reduction in accuracy on this version of LAIONet. Hence, the reduction in accuracy is consistent across all versions of LAIONets, including the most conservatively created ones. Figure 14 also confirms that this version of LAIONet still exhibits a longer tail of small intra-class similarity compared to ImageNet, potentially explaining the accuracy drop.

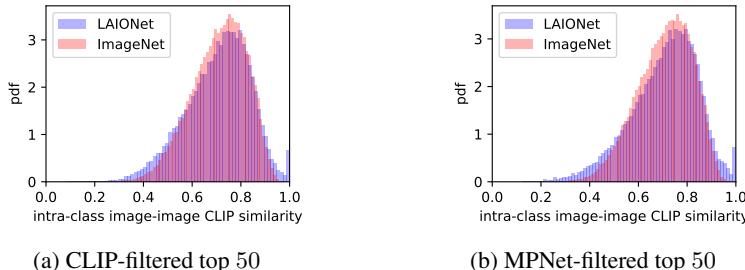


Figure 14: Comparing the intra-class similarity of LAIONet and ImageNet. In each class, pairwise similarities of LAIONet images are sampled to match ImageNet in number. All the classes combined, the distribution of intra-class similarity is depicted. LAIONet is created by retaining top 50 most similar instances to the synset text in each class after a textual similarity filtering with CLIP or MPNet.

C ON THE CHOICE OF LAION TEXT TO SYNSET TEXT SIMILARITY THRESHOLD

In Section 2, we described how LAIONet is generated through substring matching LAION texts with ImageNet synset lemmas, followed by filtering out the cases where the LAION text is not sufficiently similar to the synset name and definition. A critical choice in the second filtering step is the choice of the minimum required textual similarity. We conservatively chose this threshold to be the largest value such that the remaining examples cover a large number of ImageNet’s classes. To show this filtering is necessary and our threshold of 0.82 for CLIP-based filtering and threshold of 0.58 for MPNet-based filtering is conservative, we have provided an example in Figure 15. Here the synset “cougar” has lemma “puma”. From WordNet definition, “cougar” is a “large American feline resembling lion”. But the common usage of “puma” on the web is about a brand. As Figure 15 shows for small similarity to the synset, data most likely will represent the brand instead of the animal. As we increase the similarity threshold, the examples become more and more likely to be from the intended meaning. Our manual inspections show similar to this example, the chosen thresholds most likely result in high-quality matching to the intended meaning of the synset even if the web is dominated by other meanings.

D ON THE (NON)DIFFICULTY OF LAIONET IMAGE CLASSIFICATION

To obtain a better idea of how hard it is to recognize an object in LAIONet, we calculate the cross-modal similarity of the images to the texts of their associated synsets using CLIP embeddings. A high value of image-to-synset similarity indicates CLIP is able to identify an object from the synset in the image. On the other hand, a low value could indicate that the intended object is either absent from the image or difficult to recognize. We compare the image-to-synset similarities obtained from the ImageNet validation set and LAIONet.

Figure 16a illustrates the distribution of image-to-synset similarity for LAIONet and ImageNet. To ensure these distributions are comparable, we sampled LAIONet with replacement to match the number of images per class in the ImageNet validation set. As the figure suggests, the two datasets are not significantly different. In a more fine-grained test, we compared the image-to-synset similarity of the LAIONet and ImageNet for each class. Figure 16b shows the average similarity in each class for LAIONet subtracted by the average similarity in the same class for ImageNet along 95% upper and lower confidence bounds. Overall, there is no strong signal that LAIONet images are harder in particular.

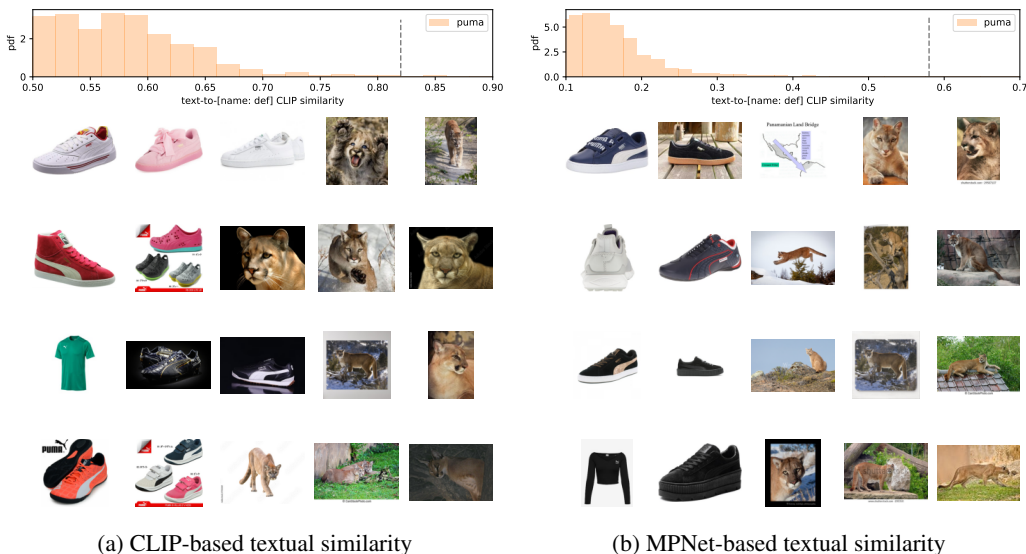


Figure 15: Sample images from five intervals of LAION text to synset text similarity.

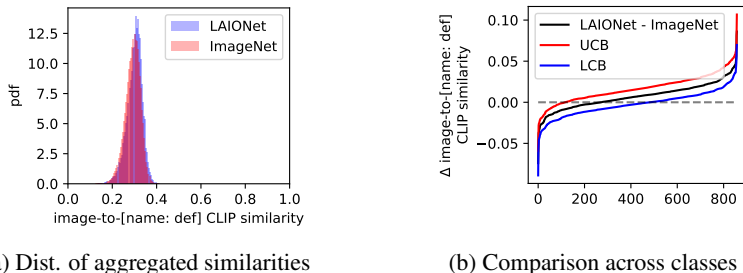


Figure 16: Comparing image-to-synset similarities of LAIONet and ImageNet. (a) For each class, LAIONet is sampled with replacement to have the same number of images as ImageNet, and all samples are aggregated to obtain the distribution. (b) For every class, the average similarity of the images to synset text is calculated for LAIONet and ImageNet and the difference is plotted. The upper and lower 0.95% confidence bound for this difference is plotted in red and blue. All values are sorted ascendingly.

E ON THE CHOICE OF TEXTUAL SIMILARITY THRESHOLD IN EXTRACTING MOST SIMILAR LAION INSTANCES TO IMAGENET-CAPTIONS

In Section 4.3, we selected a similarity threshold of 0.7 as the minimum requirement for similarity between LAION text and ImageNet text in order to include a sample from LAION. Ideally, we look for LAION examples with identical text as the ImageNet but due to the limited number of samples available in LAION, this is not possible. As Figure 3b shows, increasing the similarity threshold beyond the chosen level of 0.7 significantly decreases the number of covered classes. Meanwhile, for larger thresholds, the new dataset looks more like ImageNet but is still distinguishable. As Figure 17b shows, the proportion of classes with significantly lower intra-class similarity in ImageNet increases as the threshold increases, while the proportion of classes with significantly lower intra-class similarity in the new dataset decreases. The gap still persists but can potentially become smaller in the region our data cannot cover. In sum, the new dataset extracted based on ImageNet looks unlike ImageNet but to the extent it is possible to find similar texts in LAION.

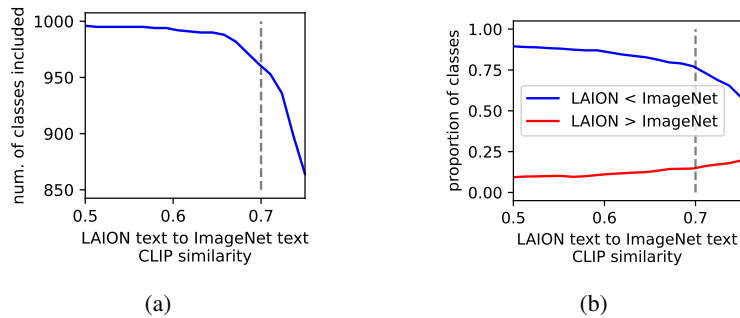


Figure 17: The effect of similarity threshold on the dataset extracted from LAION samples with most similar texts to the ImageNet texts. (a) Number of the classes covered in the new dataset versus the similarity threshold. (b) Proportion of classes with significantly lower intra-class similarity in the new dataset (blue) and proportion of classes with significantly lower intra-class similarity in ImageNet (red) versus the similarity threshold.

F THE RELATION OF RECALL, RELATIVE FREQUENCY, AND INTRA-CLASS SIMILARITY

F.1 RECALL VERSUS RELATIVE FREQUENCY

In Section 3.1 we observed accuracy drops when we weight different classes according to their frequency in LAIONet. This can be partially explained as models perform worse in more frequent classes. To directly observe this, Figure 18 shows the recall in each class versus the relative frequency of the class in LAIONet. Regardless of whether LAIONet is created by filtering based on CLIP textual similarity or MPNet similarity, there exists a weak but consistent trend that more frequent classes are more likely to be misclassified.

F.2 RECALL VERSUS INTRA-CLASS SIMILARITY

Section 3.2 introduced the hypothesis that higher intra-class similarity may account for the lower-than-expected performance of ImageNet models on LAIONet. To observe that intra-class similarity can be responsible for accuracy drop, Figure 19 demonstrates that models struggle on classes where LAIONet is more diverse than ImageNet, as shown by the recall rates plotted against the difference in average intra-class similarity. This is true regardless of what notion of accuracy and what version of LAIONet, CLIP-filtered or MPNet-filtered, is used.

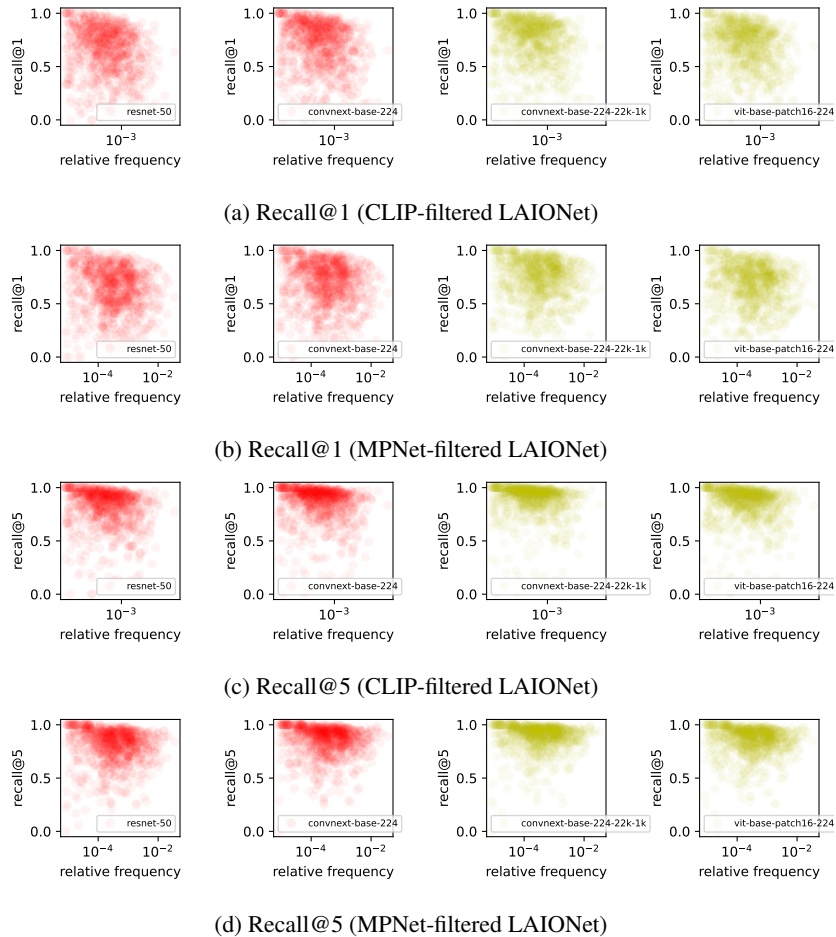


Figure 18: Recall per class evaluated on LAIONet versus how frequent the class is in LAIONet. Four different models are used, where two of them are pretrained on ImageNet-21k and two of them are not. Two versions of LAIONet, CLIP-filtered and MPNet-filtered are included. Trends are consistent.

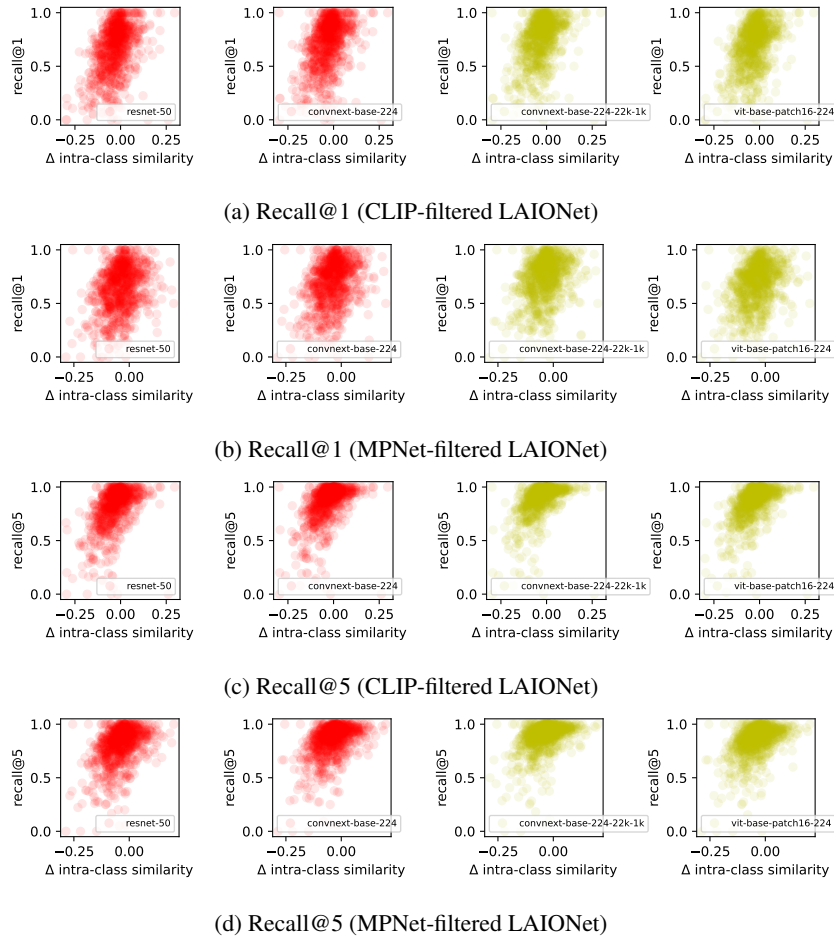


Figure 19: Recall on LAIONet for each class versus the disparity in intra-class similarity between LAIONet and ImageNet. This disparity (horizontal axis) is measured by subtracting the class-average intra-class similarity in ImageNet from that in LAIONet. Four exemplary models are shown, where two of them are pretrained on ImageNet-21k (yellow) and two of them are not (red). Two versions of LAIONet are considered. Trends are consistent.

G SAMPLE IMAGES FROM LAIONET

We provide randomly picked images from both CLIP-filtered and MPNet-filtered LAIONet (Appendix B) in this section. These images have been chosen based on various levels of difficulty. Figure 20 illustrates the distribution of the recall@5 difference for each common class between LAIONet and ImageNet. We choose recall@5 as a more reliable metric where the multiplicity of labels is less of a concern. One can see that there exist classes for which the recall on LAIONet is less than ImageNet for 0.5 or more. These are typically the classes for which LAIONet may have used a broader meaning for the synset or the images have appeared in a different context than ImageNet. It is worth noting that these classes make up a very small portion of all classes and have minimal impact on evaluations, whether or not including such images is desired.

For the classes labeled on the graphs of Figure 20, we have provided 10 random images from all datasets in the following. Each figure comes with a potential explanation for the failure of ImageNet models in the caption.

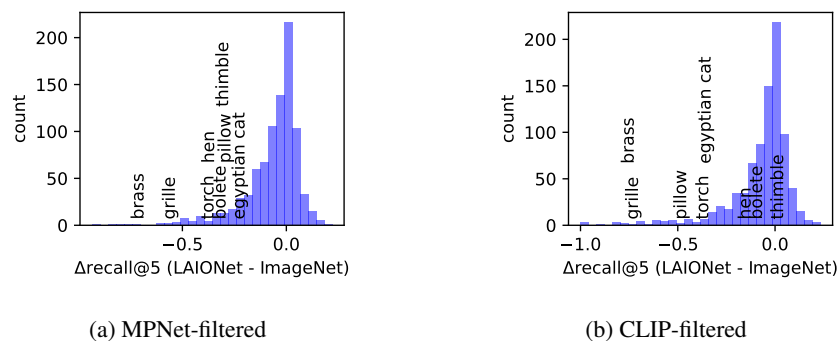


Figure 20: Distribution of recall@5 on LAIONet subtracted by recall@5 on ImageNet. Only common classes are considered. The texts show the chosen classes for which example images are provided. The position of each text on the horizontal axis is the difference in recalls for that class.

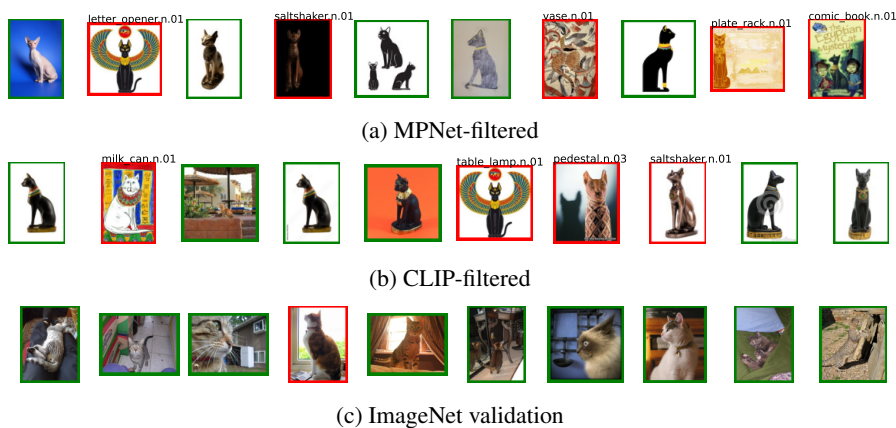


Figure 21: Egyptian cat. ImageNet models primarily struggle with Egyptian cat statues or painted graphics, which are not well-represented or are rare in the ImageNet dataset.

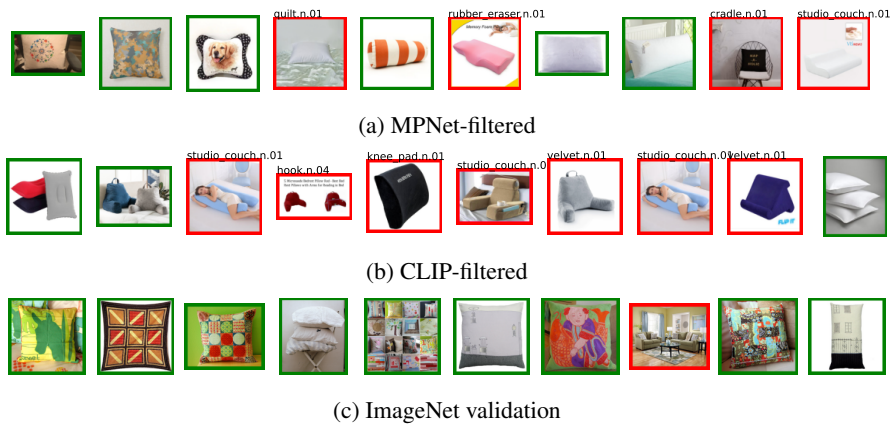


Figure 22: Pillow. ImageNet models struggle to identify pillows when they deviate from the predominantly rectangular shape that is common in ImageNet.

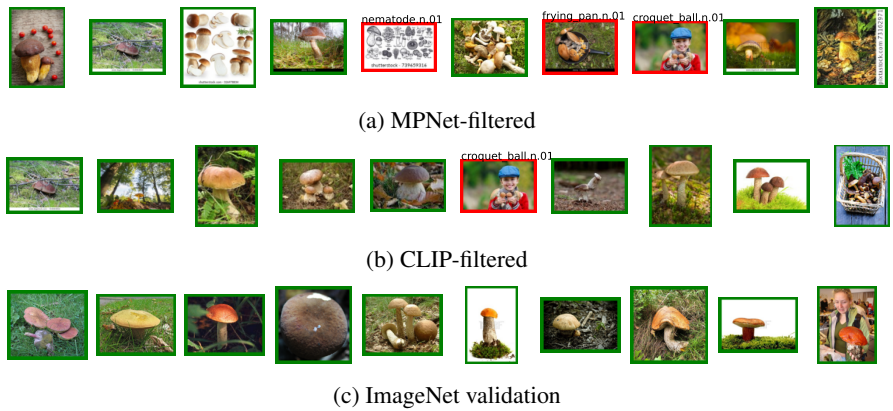


Figure 23: Bolete. ImageNet models are challenged when a bolete appears in contexts outside of nature, such as being picked by a girl or found in a pan.

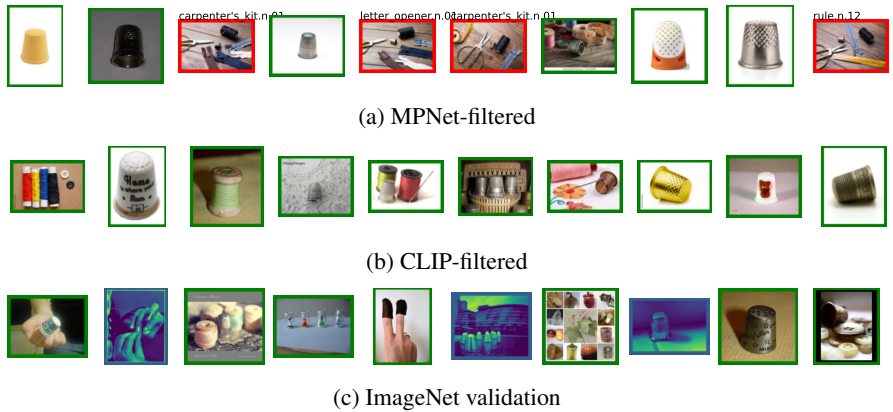


Figure 24: Thimble. ImageNet models are challenged when the thimble is among many other items.

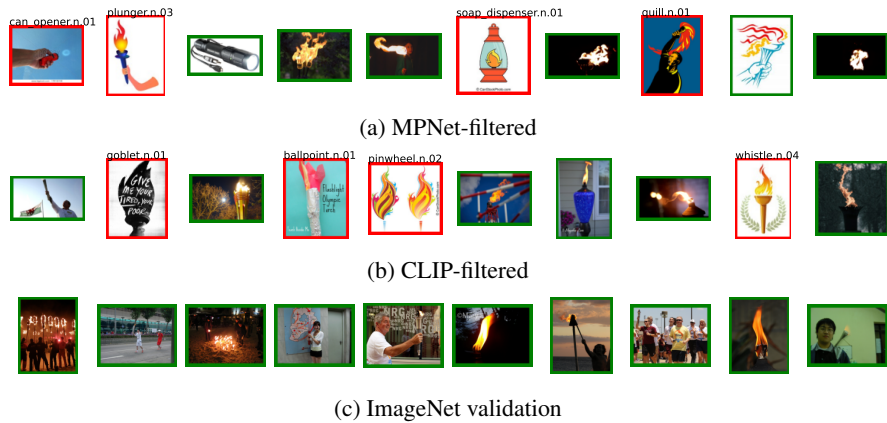


Figure 25: Torch. ImageNet models have difficulty with recognizing graphical depictions of torches and identifying variations in torch orientation.

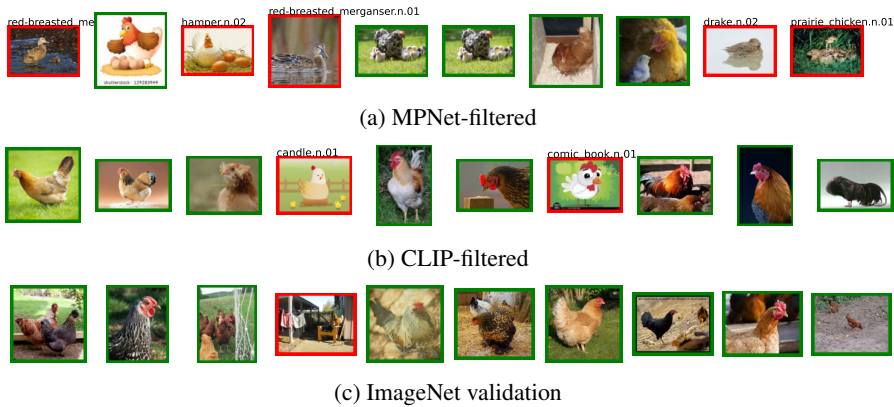


Figure 26: Hen. Graphical hens pose a challenge for ImageNet models. MPNet-filtered images also include blue and green-winged teal hens, which are not present in the ImageNet dataset.

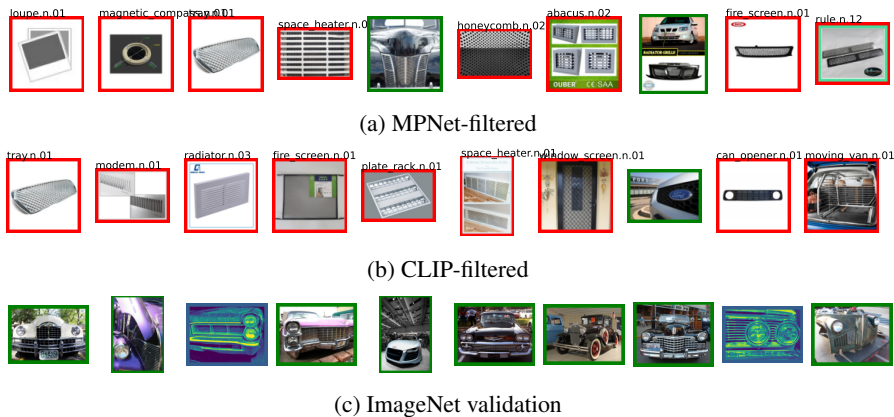


Figure 27: Grille. ImageNet models only recognize grille when installed on a car. LAIONet images also include various kinds of grille which are not meant by ImageNet class.

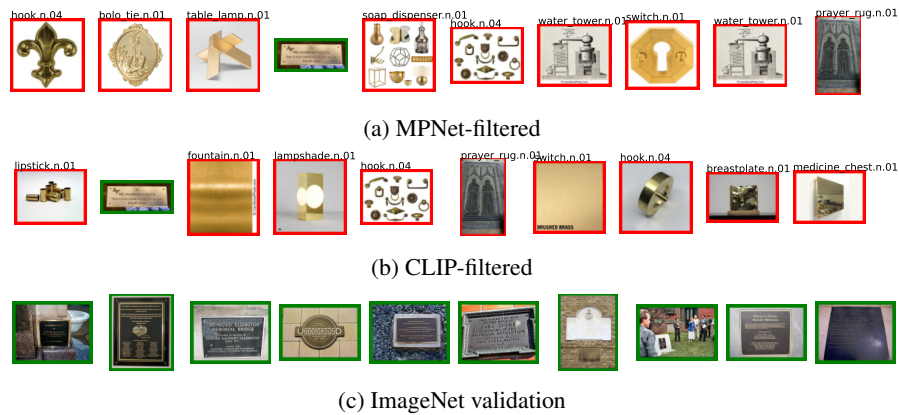


Figure 28: Brass. The intended concept of this class in ImageNet is a memorial made of brass. However, LAIONet images correspond to the broader meaning and the model is not expected to predict that.

H LAION-WEIGHTED ACCURACY EVALUATED ON IMAGENET

In Section 3.1 we introduced LAION-weighted accuracy where we use the relative frequency of each class in LAIONet to weight its recall. As we presented in Figure 5, the LAION-weighted accuracy is consistently lower than the equally-weighted accuracy when models are evaluated on LAIONet. This observation is not limited to evaluation on LAIONet. In fact, Figure 29 shows when we weight the classes according to their relative frequency on LAIONet, ImageNet accuracy also decreases. This can be attributed to the challenge of recognizing more frequent objects, given their potentially diverse types.

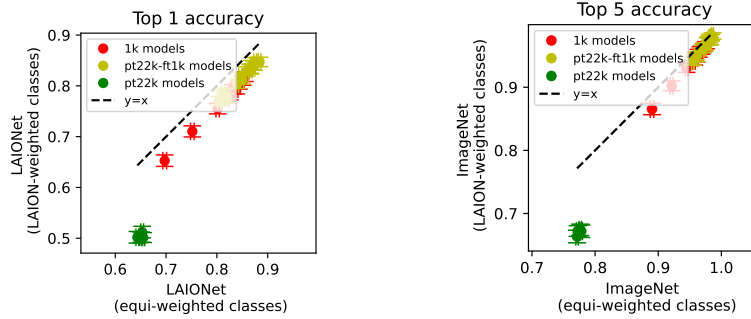


Figure 29: On ImageNet, a LAION-weighted accuracy is calculated according to the relative frequency of the classes in LAIONet and compared to the accuracy with equally weighted classes.