

Towards Fairer Datasets: Filtering and Balancing the Distribution of the People Subtree in the ImageNet Hierarchy

Kaiyu Yang
Princeton University
Princeton, NJ
kaiyuy@cs.princeton.edu

Klint Qinami
Princeton University
Princeton, NJ
kqinami@cs.princeton.edu

Li Fei-Fei
Stanford University
Stanford, CA
feifeili@cs.stanford.edu

Jia Deng
Princeton University
Princeton, NJ
jiadeng@cs.princeton.edu

Olga Russakovsky
Princeton University
Princeton, NJ
olgarus@cs.princeton.edu

ABSTRACT

Computer vision technology is being used by many but remains representative of only a few. People have reported misbehavior of computer vision models, including offensive prediction results and lower performance for underrepresented groups. Current computer vision models are typically developed using datasets consisting of manually annotated images or videos; the data and label distributions in these datasets are critical to the models' behavior. In this paper, we examine ImageNet, a large-scale ontology of images that has spurred the development of many modern computer vision methods. We consider three key factors within the person subtree of ImageNet that may lead to problematic behavior in downstream computer vision technology: (1) the stagnant concept vocabulary of WordNet, (2) the attempt at exhaustive illustration of all categories with images, and (3) the inequality of representation in the images within concepts. We seek to illuminate the root causes of these concerns and take the first steps to mitigate them constructively.

1 INTRODUCTION

As computer vision technology becomes widespread in people's Internet experience and daily lives, it is increasingly important for computer vision models to produce results that are appropriate and fair. However, there are notorious and persistent issues. For example, face recognition systems have been demonstrated to have disproportionate error rates across race groups, in part attributed to the underrepresentation of some skin tones in face recognition datasets [10]. Models for recognizing human activities perpetuate gender biases after seeing the strong correlations between gender and activity in the data [36, 83]. The downstream effects range from perpetuating harmful stereotypes [55] to increasing the likelihood of being unfairly suspected of a crime (e.g., when face recognition models are used in surveillance cameras).

Many of these concerns can be traced back to the datasets used to train the computer vision models. Thus, questions of fairness and representation in datasets have come to the forefront. In this

work, we focus on one dataset, ImageNet [18], which has arguably been the most influential dataset of the modern era of deep learning in computer vision. ImageNet is a large-scale image ontology collected to enable the development of robust visual recognition models. The dataset spearheaded multiple breakthroughs in object recognition [35, 45, 70]. In addition, the feature representation learned on ImageNet images has been used as a backbone for a variety of computer vision tasks such as object detection [33, 61], human activity understanding [69], image captioning [75], and recovering depth from a single RGB image [49], to name a few. Works such as Huh et al. [37] have analyzed the factors that contributed to ImageNet's wide adoption. Despite remaining a free education dataset released for non-commercial use only,¹ the dataset has had profound impact on both academic and industrial research.

With ImageNet's large scale and diverse use cases, we examine the potential social concerns or biases that may be reflected or amplified in its data. It is important to note here that references to "ImageNet" typically imply a subset of 1,000 categories selected for the image classification task in the ImageNet Large Scale Visual Recognition Challenge (ILSVRC) of 2012-2017 [64], and much of the research has focused on this subset of the data. So far, Dulhanty and Wong [22] studied the demographics of people in ILSVRC data by using computer vision models to predict the gender and age of depicted people, and demonstrated that, e.g., males aged 15 to 29 make up the largest subgroup.² Stock and Cisse [71] did not explicitly analyze the dataset but demonstrate that models trained on ILSVRC exhibit misclassifications consistent with racial stereotypes. Shankar et al. [67] and DeVries et al. [19] showed that most images come from Europe and the United States, and the resulting models have difficulty generalizing to images from other places. Overall, these studies identify a small handful of protected attributes and analyze their distribution and/or impact within the ILSVRC dataset, with the goal of *illuminating* the existing bias.

Goals and contributions. There are two key distinctions of our work. First, we look beyond ILSVRC [64] to the broader ImageNet dataset [18]. As model accuracy on the challenge benchmark is now near-perfect, it is time to examine the larger setting of ImageNet. The 1,000 categories selected for the challenge contain only 3 people

FAT* '20, January 27–30, 2020, Barcelona, Spain

© 2020 Copyright held by the owner/author(s). Publication rights licensed to ACM. This is the author's version of the work. It is posted here for your personal use. Not for redistribution. The definitive Version of Record was published in *Conference on Fairness, Accountability, and Transparency (FAT* '20)*, January 27–30, 2020, Barcelona, Spain, <https://doi.org/10.1145/3351095.3375709>.

¹Please refer to ImageNet terms and conditions at image-net.org/download-faq

²As noted by the authors themselves, this approach poses a chicken-or-egg problem as trained models are likely to exhibit gender or age biases, thus limiting their ability to accurately benchmark dataset bias.

categories (scuba diver, bridegroom, and baseball player) while the full ImageNet contains 2,832 people categories under the person subtree (accounting for roughly 8.3% of the total images). Their use can be problematic and raises important questions about fairness and representation. In this work, we focus on the person subtree of the full ImageNet hierarchy.

Second, in contrast to prior work, our goal is to do a deeper analysis of the root causes of bias and misrepresentation, and to propose concrete steps towards mitigating them. We identify three key factors that may lead to problematic behavior in downstream technology: (1) the stagnant concept vocabulary from WordNet [53], (2) the attempt at exhaustive illustration of all categories with images, and (3) the inequality of demographic representation in the images. For each factor, we seek to illuminate the root causes of the concern, take the first steps to mitigate them through carefully designed annotation procedures, and identify future paths towards comprehensively addressing the underlying issues.

Concretely, we thoroughly analyze the person subtree of ImageNet and plan to modify it along several dimensions. First, in Sec. 4, we examine the 2,832 people categories that are annotated within the subtree, and determine that 1,593 of them are potentially offensive labels that should not be used in the context of an image recognition dataset. We plan to remove all of these from ImageNet. Second, in Sec. 5, out of the remaining 1,239 categories we find that only 158 of them are visual, with the remaining categories simply demonstrating annotators' bias. We recommend further filtering the person subtree down to only these 158 categories when training visual recognition models. Finally, in Sec. 6 we run a large-scale crowdsourcing study to manually annotate the gender, skin color, and age of the people depicted in ImageNet images corresponding to these remaining categories. While the individual annotations may be imperfect despite our best efforts (e.g., the annotated gender expression may not correspond to the depicted person's gender identity), we can nevertheless compute the approximate demographic breakdown. We believe that releasing these sensitive attribute annotations directly is not the right step for ethical reasons, and instead plan to release a Web interface that allows an interested user to *filter* the images within a category to achieve a new target demographic distribution.

We are working on incorporating these suggestions into the dataset. We will additionally release our annotation interfaces³ to allow for similar cleanup of other computer vision benchmarks.

2 RELATED WORK ON FAIRNESS IN MACHINE LEARNING

We begin with a more general look at the literature that identified or attempted to mitigate bias in modern artificial intelligence systems. In short, datasets often have biased distributions of demographics (gender, race, age, etc.); machine learning models are trained to exploit whatever correlations exist in the data, leading to discriminatory behavior against underrepresented groups [6, 7]. A great overview of the history of machine learning fairness can be found in Hutchinson and Mitchell [38]. The approaches to address fairness concerns fall loosely into two categories: (1) identifying

and correcting issues in datasets or (2) studying and encouraging responsible algorithmic development and deployment.

Identifying and addressing bias in datasets. There are several issues raised in the conversation around dataset bias. The first common issue is the lack of transparency around dataset design and collection procedures. Datasheets for Datasets [31] (and, relatedly, for models [54]) have been proposed as solutions, encouraging dataset creators to release detailed and standardized information on the collection protocol which can be used by downstream users to assess the suitability of the dataset. Throughout this work we dive deep into understanding the data collection pipelines of ImageNet and related datasets, and consider their implications.

The second issue is the presence of ethically questionable concepts or annotations within datasets. Examples range from quantifying beauty [47] to predicting sexual orientation [76] to (arguably) annotating gender [25, 42]. In Sec. 4 (and to a lesser extent in Sec. 5), we consider the underlying cause for such concepts to appear in large-scale datasets and propose the first steps of a solution.

A related consideration is the ethics and privacy of the subjects depicted in these computer vision datasets. Here we refer the reader to, e.g., Whittaker et al. [78] for a recent detailed discussion as this is outside the scope of our work.

The final and perhaps best-known source of dataset bias is the imbalance of representation, e.g., the underrepresentation of demographic groups within the dataset as a whole or within individual classes. In the context of computer vision, this issue has been brought up at least in face recognition [10], activity recognition [83], facial emotion recognition [62], face attribute detection [65] and image captioning [11] – as well as more generally in pointing out the imaging bias of datasets [73]. This is not surprising as many of the images used in computer vision datasets come from Internet image search engines, which have been shown to exhibit similar biased behavior [14, 41, 55]. In some rare cases, a dataset has been collected explicitly to avoid such influences, e.g., the Pilot Parliaments Benchmark (PPB) by Buolamwini and Gebru [10]. In Sec. 6 we propose a strategy for balancing ImageNet across several protected attributes while considering the implications of such design (namely the concerns with annotating pictures of individual people according to these protected attributes).

Responsible algorithmic development. Beyond efforts around better dataset construction, there is a large body of work focusing on the development of fair and responsible algorithms that aim to counteract the issues which may be present in the datasets. Researchers have proposed multiple fairness metrics including statistical parity [12, 13, 24, 39], disparate impact [27, 80], equalized odds [34] and individual fairness [23], and analyzed the relationship between them [43, 60]. Algorithmic solutions have ranged from removing undesired bias by preprocessing the data [39, 59], striking a tradeoff between performance and fairness by posing additional regularization during training or inference [34, 40, 52, 80, 81, 83], or designing application-specific interventions (such as of Burns et al. [11] for reducing gender bias in image captioning models).

However, statistical machine learning models have three fundamental limitations that need to be considered. First, the accuracy of a machine learning model is strongly influenced by the number of training examples: underrepresented categories in datasets will

³image-net.org/filtering-and-balancing

be inherently more challenging for the model to learn [51]. Second, machine learning models are statistical systems that aim to make accurate predictions on the majority of examples; this focus on common-case reasoning encourages the models to ignore some of the diversity of the data and make simplifying assumptions that may *amplify* the bias present in the data [83]. Finally, learning with constraints is a difficult open problem, frequently resulting in satisfying fairness constraints at the expense of overall model accuracy [39, 40, 80, 81]. Thus, algorithmic interventions alone are unlikely to be the most effective path toward fair machine learning, and dataset interventions are necessary. Even more so, most algorithmic approaches are supervised and require the protected attributes to be explicitly annotated, again bringing us back to the need for intervention at the dataset level.

Datasets, algorithms and intention. Finally, we note that prior work in this space underscores a single important point: any technical fairness intervention will only be effective when done in the context of the broader awareness, intentionality and thoughtfulness in building applications. Poorly constructed datasets may introduce unnoticed bias into models. Poorly designed algorithms may exploit even well-constructed datasets. Accurate datasets and models may be used with malicious intent. The responsibility for downstream fair systems lies at all steps of the development pipeline.

3 BACKGROUND: THE IMAGENET DATA COLLECTION PIPELINE

To lay the groundwork for the rest of the paper, we begin by summarizing the data collection and annotation pipeline used in ImageNet as originally described in [18, 64]. This section can be safely skipped for readers closely familiar with ImageNet and related computer vision datasets, but we provide it here for completeness.

The goal of ImageNet is to illustrate English nouns with a large number of high-resolution carefully curated images as to “foster more sophisticated and robust models and algorithms to index, retrieve, organize and interact with images and multimedia data” [18]. We consider the entire ImageNet dataset consisting of 14,197,122 images illustrating 21,841 concepts rather than just the 1,431,167 images illustrating 1,000 concepts within the ImageNet challenge which are most commonly used [64]. There are three steps to the ImageNet data collection pipeline: (1) selecting the concept vocabulary to illustrate, (2) selecting the candidate images to consider for each concept, and (3) cleaning up the candidates to ensure that the images in fact correspond to the target concept. We describe each step and its similarities to the steps in other vision datasets.

(1) Concept vocabulary. When building a visual recognition benchmark, the first decision is settling on a concept vocabulary and decide which real-world concepts should be included. WordNet [53] emerges as a natural answer. It is a language ontology in which English nouns are grouped into sets of synonyms (synsets) that represent distinct semantic concepts.⁴ The synsets are then organized into a hierarchy according to the “is a” relation, such as “coffee table is a table”. WordNet serves as the concept vocabulary for

ImageNet, which provides images for grounding the synsets visually. Similarly, subsets of the WordNet backbone have been used in datasets like Places [84], Visual Genome [44] and ShapeNet [15].

(2) Candidate images. The natural and easiest-to-access source of visual data is the Internet. For every concept in WordNet, the ImageNet creators queried image search engines and aimed to increase the variety of retrieved images through using multiple search engines, employing query expansion, and translating the search terms into multiple languages [18]. Similarly, the vast majority of vision datasets rely on images collected from the Internet [26, 44, 46, 48, 74, 84], with many collected by first defining the set of target concepts and then obtaining the associated images using query expansion: e.g., PASCAL [26], COCO [48], Places [84].

(3) Manual cleanup. As noted in Torralba et al. [74], image search engines were only about 10% accurate, and thus a manual cleanup of the candidate images is needed. The cleanup phase of ImageNet consists of a set of manual annotation tasks deployed on the Amazon Mechanical Turk (AMT) marketplace. The workers are provided with a single target concept (e.g., Burmese cat), its definition from WordNet, a link to Wikipedia, and a collection of candidate images. They are instructed to click on all images that contain the target concept, irrespective of any occlusion, scene clutter, or the presence of other concepts. Images that reach a desired level of positive consensus among workers are added to ImageNet. Similarly, most computer vision datasets rely on manual annotation although the details change: e.g., PASCAL was annotated in-house rather than using crowdsourcing [26], Places relies on both positive and negative verification [84], COCO favors very detailed annotation per image [48], and Open Images [46] and Places [84] both use a computer-assisted annotation approach.

Outline. In the following sections, we consider the fairness issues that arise as a result of this pipeline, and propose the first steps to mitigate these concerns.

4 PROBLEM 1: STAGNANT CONCEPT VOCABULARY

The backbone of WordNet [53] provides a list of synsets for ImageNet to annotate with images. However, born in the past century, some synsets in WordNet are no longer appropriate in the modern context. People have reported abusive synsets in the WordNet hierarchy, including racial and sexual slurs (e.g., synsets like n10585077 and n10138472).⁵ This is especially problematic within the person subtree of the concept hierarchy (i.e., synset n00007846 and its descendants). During the construction of ImageNet in 2009, the research team removed any synset explicitly denoted as “offensive”, “derogatory”, “pejorative”, or “slur” in its gloss, yet this filtering was imperfect and still resulted in inclusion of a number of synsets that are offensive or contain offensive synonyms. Going further, some synsets may not be inherently offensive but may be inappropriate for inclusion in a visual recognition dataset. This filtering of the concept vocabulary is the first problem that needs to be addressed.

⁴We use the words “concept” and “synset” interchangeably throughout the paper.

⁵Throughout the paper, we refrain from explicitly listing the offensive concepts associated with synsets and instead report only their synset IDs. For a conversion, please see wordnet.princeton.edu/documentation/wndb5wn.

4.1 Prior work on annotating offensiveness

Sociolinguistic research has explored the problem of offensiveness, largely focusing on studying profanity. Ofcom [56] and Sapolsky et al. [66] rate the offensiveness of words in TV programs. Dewaele [20] demonstrates offensiveness to be dependent on language proficiency by studying the ratings from English native speakers and non-native speakers. Beers Fägersten [8] designs two questionnaires to study the role of context in the level of offensiveness of profanity. In one questionnaire, subjects see a word together with a short dialogue in which the word appears. They are asked to rate the word on a 1-10 scale from “not offensive” to “very offensive”. In another questionnaire, the subjects see only the words without any context. The findings highlight the importance of context, as the perceived offensiveness depends heavily on the dialogue and on the gender and race of the subjects.

4.2 Methodology for filtering the unsafe synsets

We ask annotators to flag a synset as unsafe when it is either inherently “offensive,” e.g., containing profanity or racial or gender slurs, or “sensitive,” i.e., not inherently offensive but may cause offense when applied inappropriately, such as the classification of people based on sexual orientation and religion. In contrast to prior work, we do not attempt to quantify the “level” of offensiveness of a concept but rather exclude all potentially inappropriate synsets. Thus, we do not adopt a 1-5 or 1-10 scale [8, 20] for offensiveness ratings. Instead, we instruct workers to flag any synset that may be potentially unsafe, essentially condensing the rating 2-5 or 2-10 on the scale down to a single “unsafe” label.

4.3 Results and impact on ImageNet after removing the unsafe synsets

We conduct the initial annotation using in-house workers, who are 12 graduate students in the department and represent 4 countries of origin, male and female genders, and a handful of racial groups. The instructions are available in Appendix. So far out of 2,832 synsets within the person subtree, we have identified 1,593 unsafe synsets. The remaining 1,239 synsets are temporarily deemed safe. Table 1 gives some examples of the annotation results (with the actual content of offensive synsets obscured). The full list of synset IDs can be found in Appendix. The unsafe synsets are associated with 600,040 images in ImageNet. Removing them would leave 577,244 images in the safe synsets of the person subtree of ImageNet.

4.4 Limitations of the offensiveness annotation and future work

First, it is important to note that a “safe” synset only means that the label itself is not deemed offensive. It does not mean that it is possible, useful, or ethical to infer such a label from visual cues.

Second, despite the preliminary results we have, offensiveness is subjective and also constantly evolving, as terms develop new cultural context. Thus, we are opening up this question to the community. We are in the process of updating the ImageNet website to allow users to report additional synsets as unsafe. While the dataset may be large scale, the number of remaining concepts is

relatively small (here in the low thousands and further reduced in the next section), making this approach feasible.

5 PROBLEM 2: NON-VISUAL CONCEPTS

ImageNet attempts to depict each WordNet synset with a set of images. However, not all synsets can be characterized visually. For example, is it possible to tell whether a person is a philanthropist from images? This issue has been partially addressed in ImageNet’s annotation pipeline [18] (Sec. 3), where candidate images returned by search engines were verified by human annotators. It was observed that different synsets need different levels of consensus among annotators, and a simple adaptive algorithm was devised to determine the number of agreements required for images from each synset. Synsets harder to characterize visually would require more agreements, which led to fewer (or no) verified images.

Despite the adaptive algorithm, we find a considerable number of the synsets in the person subtree of ImageNet to be non-imageable—hard to characterize accurately using images. There are several reasons. One reason for them sneaking into ImageNet could be the large-scale annotation. Although non-imageable synsets require stronger consensus and have fewer verified images, they remain in ImageNet as long as there are some images successfully verified, which is likely given the large number of images. Another reason could be “positive bias”: annotators are inclined to answer “yes” when asked the question “Is there a <concept> in the image?” As a result, some images with weak visual evidence of the corresponding synset may be successfully verified.

The final and perhaps most compelling reason for non-imageable synsets to have been annotated in ImageNet is that search engines will surface the most *distinctive* images for a concept, even if the concept itself is not imageable. For example, identifying whether someone is Bahamian from a photograph is not always possible, but there will be some distinctive images (e.g., pictures of a people wearing traditional Bahamian costumes), and those will be the ones returned by the search engine. This issue is amplified by the presence of stock photography on the web, which contributes to and perpetuates stereotypes as discussed at length in e.g., [4, 29, 30]. Overall, this results in an undoubtedly biased visual representation of the categories, and while the issue affects all synsets, it becomes particularly blatant for categories that are inherently non-imageable. Thus in an effort to reduce the visual bias, we explicitly determine the imageability of the synsets in the person subtree and recommend that the community refrain from using those with low imageability when training visual recognition models.

5.1 Annotating imageability

Extensive research in psycholinguistics has studied the imageability (a.k.a. imagery) of words [5, 9, 17, 32, 57, 58], which is defined as “the ease with which the word arouses imagery” [58]. For annotating imageability, most prior works follow a simple procedure proposed by Paivio et al. [58]: The workers see a list of words and rate each word on a 1-7 scale from “low imagery (1)” to “high imagery (7)”. For each word, the answers are averaged to establish the final score.

We adopt this definition of imageability and adapt the existing procedure to annotate the imageability of synsets in the ImageNet person subtree. However, unlike prior works that use in-house

Table 1: Examples of synsets in the person subtree annotated as unsafe (offensive), unsafe (sensitive), safe but non-imageable, and simultaneously safe and imageable. For unsafe (offensive) synsets, we only show their synset IDs. The annotation procedure for distinguishing between unsafe and safe synsets is described in Sec. 4; the procedure for non-imageable vs. imageable is in Sec. 5. We recommend removing the synsets of the first two columns from ImageNet entirely, and refrain from using synsets from the third column when training visual recognition models.

Unsafe (offensive)	Unsafe (sensitive)	Safe non-imageable	Safe imageable
n10095420: <sexual slur>	n09702134: Anglo-Saxon	n10002257: demographer	n10499631: Queen of England
n10114550: <profanity>	n10693334: taxi dancer	n10061882: epidemiologist	n09842047: basketball player
n10262343: <sexual slur>	n10384392: orphan	n10431122: piano maker	n10147935: bridegroom
n10758337: <gendered slur>	n09890192: camp follower	n10098862: folk dancer	n09846755: beekeeper
n10507380: <criminative>	n10580030: separatist	n10335931: mover	n10153594: gymnast
n10744078: <criminative>	n09980805: crossover voter	n10449664: policyholder	n10539015: ropewalker
n10113869: <obscene>	n09848110: theist	n10146104: great-niece	n10530150: rider
n10344121: <pejorative>	n09683924: Zen Buddhist	n10747119: vegetarian	n10732010: trumpeter

workers to annotate imageability, we rely on crowdsourcing. This allows us to scale the annotation and obtain ratings from a diverse pool of workers [21, 63], but also poses challenges in simplifying the instructions and in implementing robust quality control.

In our crowdsourcing interface (in Appendix), we present the human subject with a list of concepts and ask them to identify how easy it is to form a mental image of each. To reduce the cognitive load on the workers, we provide a few examples to better explain the task, include the synonyms and the definition of each concept from WordNet, and change the rating to be a simpler 5-point (rather than 7-point) scale from “very hard (1)” to “very easy (5)”. The final imageability score of a synset is an average of the ratings.

For quality control, we manually select 20 synsets as gold standard questions (in Appendix); half of them are obviously imageable (should be rated 5), and the other half are obviously non-imageable (should be rated 1). They are used to evaluate the quality of workers. If a worker has a high error on the gold standard questions, we remove all the ratings from this worker. We also devise a heuristic algorithm to determine the number of ratings to collect for each synset. Please refer to Appendix for details.

5.2 Results and impact on ImageNet after removing the non-imageable synsets

We annotate the imageability of 1,239 synsets in the person subtree which have been marked as safe synsets in the previous task. Fig. 1 shows the imageability ratings for a selected set of synsets. Synsets such as *irreligionist* and *nurse* have well-accepted imageability (*irreligionist* is deemed to be decidedly non-imageable, *nurse* is deemed to be clearly imageable). In contrast, it is much harder to reach a consensus on the imageability of *host* and *waltzer*. Fig. 2 shows the distribution of the final imageability scores for all of the 1,239 safe synsets. The median is 2.60; only 158 synsets have imageability greater than or equal to 4. Table 1 shows some examples of non-imageable synsets. The complete list is in Appendix.

After manually examining the results, we suggest that all synsets in the person subtree with imageability less than 4 be considered “non-imageable” and not be used for training models. There would be 443,547 images and 1,081 synsets flagged, including *hobbyist* (1.20), *job candidate* (2.64), and *bookworm* (3.77); there would be

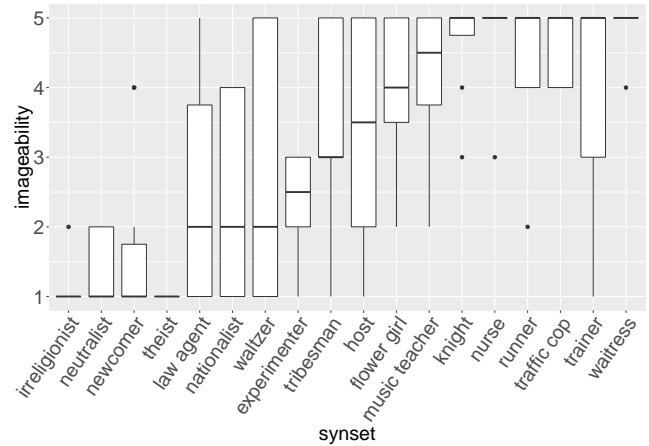


Figure 1: The distribution of raw imageability ratings for selected synsets. *irreligionist* and *nurse* have more well-accepted imageability than *host* and *waltzer*.

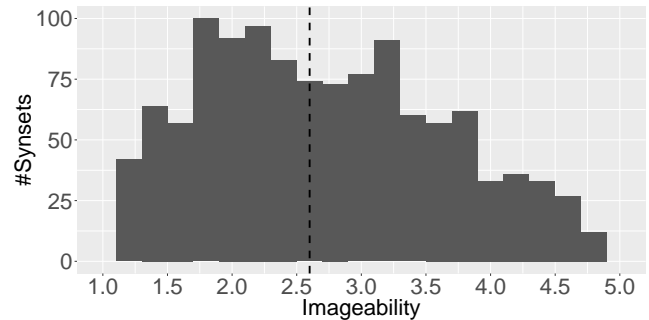


Figure 2: The distribution of the final imageability scores for all of the 1,239 safe synsets. The median is 2.60.

133,697 images and 158 synsets remaining, including *rock star* (4.86), *skier* (4.50), and *cashier* (4.20). More examples are in Table 1. Future researchers are free to adjust this threshold as needed.

5.3 Limitations of the imageability annotation

By manually examining a subset of the synsets, we find the imageability results to be reasonable overall, but we also observe a few interesting exceptions. Some synsets with high imageability are actually hard to characterize visually, e.g., `daughter` (5.0) and `sister` (4.6); they should not have any additional visual cues besides being female. Their high imageability scores could be a result of the mismatch between “the ease to arouse imagery” and “the ease to characterize using images”. `daughter` and `sister` are hard to characterize visually, but they easily arouse imagery if the annotator has a daughter or a sister. The definition based on ease of characterization with visual cues is more relevant to computer vision datasets, but we adopt the former definition as a surrogate since it is well-accepted in the literature, and there are mature procedures for annotating it using human subjects.

Another interesting observation is that workers tend to assign low imageability to unfamiliar words. For example, `cotter` (a peasant in the Scottish Highlands) is scored 1.70 while the generic `peasant` is scored 3.36. Prior works have demonstrated a strong correlation between familiarity and imageability [9, 32, 79], which explains the low imageability of the less frequent `cotter`. However, low familiarity with a concept is anyway an important factor to consider in crowdsourcing dataset annotation, as unfamiliar terms are more likely to be misclassified by workers. This suggests that removing synsets identified as less imageable by our metric may also have the additional benefit of yielding a more accurate dataset.⁶

When analyzing Table 1, we further wonder whether even the synsets that are both safe and imageable should remain in ImageNet. For example, is the `Queen of England` an acceptable category for visual recognition? Would `basketball player` be better replaced with `person interacting with a basketball and captured as a human-object-interaction` annotation? Would `bridegroom` be rife with cultural assumptions and biases? As always, we urge downstream users to exercise caution when training on the dataset.

And finally, we observe that even the remaining imageable synsets may contain biased depictions as a result of search engine artifacts. For example, the synset `mother` (imageability score 4.3) primarily contains women holding children; similarly, the synset `beekeeper` (imageability score 4.6) predominantly contains pictures of people with bees. Even though one remains a mother when not around children, or a beekeeper when not around bees, those images would rarely be surfaced by a search engine (and, to be fair, would be difficult for workers to classify even if they had been).

Despite these concerns about the imageability annotations and about the lingering search engine bias, one thing that is clear is that at least the non-imageable synsets are problematic. Our annotation of imageability is not perfect and not the final solution, but an important step toward a more reasonable distribution of synsets in the ImageNet person subtree.

⁶Further, unfamiliar terms are also likely to be less common and thus less relevant to the downstream computer vision task, making their inclusion in the dataset arguably less important.

5.4 Relationship between imageability and visual recognition models

To conclude the discussion of imageability, we ask one final question: What is the relationship between the imageability of synset and the accuracy of a corresponding visual recognition model? Concretely, are the imageable synsets actually easier to recognize because they correspond to visual concepts? Or, on the flip side, is it perhaps always the case that non-imageable synsets contain an overly-simplified stereotyped representation of the concept and thus are easy for models to classify? If so, this would present additional evidence about the dangers of including such categories in a dataset since their depicted stereotypes are easily learned and perpetuated by the models.

Computer vision experiment setup. To evaluate this, we run a simple experiment to study the relationship between the imageability of a synset and the ability of a modern deep learning-based image classifier to recognize it. We pick a subset of 143 synsets from the 1,239 safe synsets so that each synset has at least 1,000 images. The selected synsets are leaf nodes in the WordNet hierarchy, meaning that they cannot be ancestors of each other and they represent disjoint concepts. We randomly sample 1,000 images from each synset, 700 for training, 100 for validation, and 200 for testing. We use a standard ResNet34 network [35] to classify the images as belonging to one of the 143 synsets. During training, the images are randomly cropped and resized to 224×224 ; we also apply random horizontal flips. During validation and testing, we take 224×224 crops at the center. The network is trained from scratch for 90 epochs, which takes two days using a single GeForce GTX 1080 GPU. We minimize the cross-entropy loss using stochastic gradient descent; the learning rate starts at 0.05 and decreases by a factor of 10 every 30 epochs. We also use a batch size of 256, a momentum of 0.9, and a weight decay of 0.001.

Computer vision experiment results. The network has an overall testing accuracy of 55.9%. We are more interested in the breakdown accuracies for each synset and how they correlate with the imageability. The network’s testing accuracy on the easily imageable synsets (score ≥ 4) is 63.8%, which is higher than the accuracy of 53.0% on the synsets deemed non-imageable (score < 4). Overall there is a positive correlation between imageability and accuracy (Pearson correlation coefficient $r = 0.23$ with a p-value of 0.0048) as depicted in Fig. 3 (left). To better understand this, we analyze four representative examples, also depicted in Fig. 3 (right), which highlight the different aspects at play here:

- *Imageable, easy to classify:* A category such as `black belt` is both deemed imageable (score of 4.4) and is easy to classify (accuracy of 92%). The retrieved images contain visually similar results that are easy to learn by the model and easy to distinguish from other people categories.
- *Non-imageable, hard to classify:* On the other end of the spectrum, `conversational partner` is deemed non-imageable (score of only 1.8) as it doesn’t evoke a prototypical visual example. The images retrieved from search engines contain groups of people engaged in conversations, so the annotators verifying these images in the ImageNet pipeline correctly labeled these images as containing a `conversation partner`.

However, the resulting set of images is too diverse, and the visual cues are too weak to be learned by the model (accuracy only 20.5%).

- *Imageable, hard to classify*: *Bridegroom* (synset ID n10147935) is an example of a category with a mismatch between imageability and accuracy. It is annotated as imageable (perfect score of 5.0), because it easily arouses imagery (albeit highly culturally-biased imagery). The retrieved search result images are as expected culturally biased, but correctly verified for inclusion in ImageNet. However, the accuracy of the classifier in this case is low (only 40%) partially because of the visual diversity of the composition of images but primarily because of confusion with a closely related synset n10148035, which also corresponds to the term *bridegroom* but with a slightly different definition (n10147935: a man who has recently been married, versus n10148035: a man participant in his own marriage ceremony). This highlights the fact that classification accuracy is not a perfect proxy for visual distinctiveness, as it depends not only on the intra-synset visual cues but also on the inter-synset variability.
- *Non-imageable, easy to classify*: Finally, *Ancient* (person who lived in ancient times) is deemed non-imageable (score of 2.5), because the imageability annotators have never seen such a person, so it is difficult to properly imagine what they might look like. However, the image search results are highly biased to ancient artifacts, including images that are not even people. The annotators agreed that these images correspond to the word *ancient*, at times making mistakes in failing to read the definition of the synset and annotating ancient artifacts as well. In the resulting set of images, visual classifiers would have no difficulty distinguishing this set of images with distinctive color patterns and unusual objects from the other people categories (accuracy 89%).

The findings highlight the intricacies of image search engine results, of the ImageNet annotation pipeline, of the imageability annotations, and of evaluating visual distinctiveness using visual classifiers. A deeper analysis is needed to understand the level of impact of each factor, and we leave that to future research. Until then, we suggest that the community refrain from using synsets deemed non-imageable when training visual recognition models, and we will update ImageNet to highlight that.

6 PROBLEM 3: LACK OF IMAGE DIVERSITY

So far we have considered two problems: the inclusion of potentially offensive concepts (which we will remove) and the illustration of non-imageable concepts with images (which we will clearly identify in the dataset). The last problem we consider is insufficient representation among ImageNet images. ImageNet consists of Internet images collected by querying image search engines [18], which have been demonstrated to retrieve biased results in terms of race and gender [14, 41, 55]. Taking gender as an example, Kay et al. [41] find that when using occupations (e.g., banker) as keywords, the image search results exhibit exaggerated gender ratios compared to the real-world ratios. In addition, bias can also be introduced during the manual cleanup phase when constructing ImageNet, as people

are inclined to give positive responses when the given example is consistent with stereotypes [41].

ImageNet has taken measures to diversify the images, such as keywords expansion, searching in multiple languages, and combining multiple search engines. Filtering out non-imageable synsets also mitigates the issue: with stronger visual evidence, the workers may be less prone to stereotypes. Despite these efforts, the bias in protected attributes remains in many synsets in the person subtree. It is necessary to study how this type of bias affects models trained for downstream vision tasks, which would not be possible without high-quality annotation of image-level demographics.

6.1 Prior work on annotating demographics

Image-level annotation of demographics is valuable for research in machine learning fairness. However, it is difficult to come up with a categorization of demographics, especially for gender and race. Buolamwini and Gebru [10] adopt a binary gender classification and the Fitzpatrick skin type classification system [28]. Zhao et al. [83] and Kay et al. [41] also adopt a binary gender classification. Besides *Male* and *Female*, Burns et al. [11] add another category *Neutral* to include people falling out of the binary gender classification. Ryu et al. [65] do not explicitly name the gender and race categories, but they have discrete categories nevertheless: five race categories (*S1, S2, S3, S4, Other*) and three gender categories (*G1, G2, Other*).

6.2 Methodology for annotating demographics

Annotated attributes. To evaluate the demographics within ImageNet and propose a more representative subset of images, we annotate a set of protected attributes on images in the person subtree. We consider U.S. anti-discrimination laws, which name race, color, national origin, religion, sex, gender, sexual orientation, disability, age, military history, and family status as protected attributes [1–3]. Of these, the only imageable attributes are color, gender, and age, so we proceed to annotate these.

- (1) *Gender*. We annotate perceived gender rather than gender identity, as someone’s gender identity may differ from their gender expression and thus not be visually prominent. It is debatable what a proper categorization of gender is and whether gender can be categorized at all. Rather than addressing the full complexity of this question, we follow prior work [10, 11, 41, 65, 83] and use a set of discrete categories: *Male*, *Female*, and *Unsure*, in which *Unsure* is used to both handle ambiguous visual cues as well as to include people with diverse gender expression.
- (2) *Skin color*. We annotate skin color according to an established dermatological metric—individual typology angle (ITA) [16]. It divides the spectrum of skin color into 6 groups, which is too fine-grained for our purpose. Instead, we combine the groups into *Light*, *Medium*, and *Dark*. Melanin index [72] is another metric for skin color, which is used by the Fitzpatrick skin type classification system. However, we opt to use the more modern ITA system. Similar to prior work [10], skin color is used as a surrogate for race membership because it is more visually salient.
- (3) *Age*. We annotate perceived age groups according to discrimination laws, which led to the categories of *Child* or *Minor*

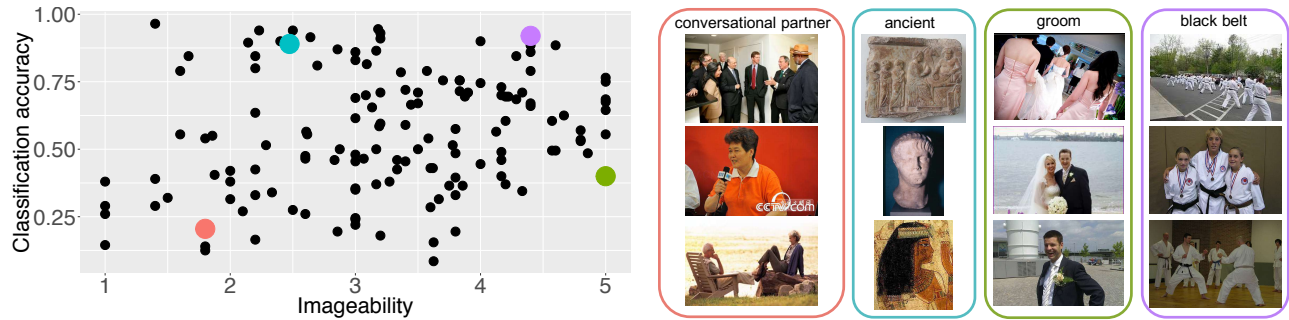


Figure 3: (Left) The computer vision model’s classification accuracy vs. synset imageability for 143 safe synsets which contain at least 1000 images. More imageable synsets are not necessarily easier for models to recognize, with Pearson correlation coefficient $r = 0.23$. **(Right)** Example images from synsets that are non-imageable and hard to classify (conversational partner); non-imageable but easy to classify (ancient); imageable but hard to classify (groom); imageable and easy to classify (black belt).

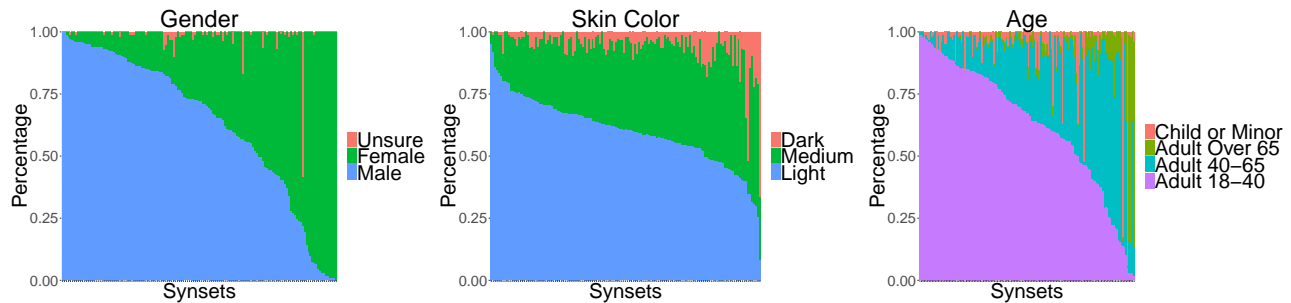


Figure 4: The distribution of demographic categories across the 139 safe and imageable synsets which contain at least 100 images. The size of the different color areas reveal the underrepresentation of certain groups.

(under 18 years old), Adult (18-40), Over 40 Adult (40-65), and Over 65 Adult (65+).

Annotation instructions. We use crowdsourcing to annotate the attributes on Amazon Mechanical Turk. We downloaded all ImageNet images from safe synsets in the person subtree whose imageability score is 4.0 or higher. An image and the corresponding synset form a task for the workers, and consists of two parts. First, the worker sees the image, the synset (including all words in it), and its definition in WordNet [53] and is asked to identify *all* persons in the image who look like members of the synset. If at least one person is identified, the worker proceeds to annotate their gender, skin color, and age. The labels are image-level, rather than specific to individual persons. There can be multiple labels for each attribute. For example, if the worker identified two persons in the first phase, they may check up to two labels when annotating the gender. The user interface is in Appendix.

The task is less well-defined when multiple persons are in the image. It can be difficult to tell which person the synset refers to, or whether the person exists in the image at all. We have tried to use automatic methods (e.g., face detectors) to detect people before manually annotating their demographics. However, the face detector is a trained computer vision model and thus also subject to dataset bias. If the face detector is only good at detecting people

from a particular group, the annotation we get will not be representative of the demographic distribution in ImageNet. Therefore, we opt to let workers specify the persons they annotate explicitly.

Quality control. For quality control, we have pre-annotated a set of gold-standard questions (in Appendix) for measuring the quality of workers. The worker’s accuracy on a gold standard question i is measured by intersection-over-union (IOU):

$$IOU_i = \frac{|A_i \cap G_i|}{|A_i \cup G_i|} \quad (1)$$

where A_i is the set of categories annotated by the worker, and G_i is the set of ground truth categories. For example, for an image containing a black female adult and a white female child, $G_i = \{Dark, Light, Female, Adult, Child\}$. If a worker mistakenly take the child to be an adult and annotates $A_i = \{Dark, Light, Female, Adult\}$, the annotation quality is computed as $IOU_i = 4/5 = 0.8$. We exclude all responses from workers whose average IOU is less than 0.5. After removing high-error workers, we aggregate the annotated categories of the same image from independent workers. Each image is annotated by at least two workers. For any specific category (e.g. *Adult*), we require consensus from $\max\{2, \lceil n_i/2 \rceil\}$ workers, where n_i is the number of workers for this image. For any image, we keep collecting annotations from independent workers until the consensus is reached. In the annotation results, the consensus

is reached with only two workers for 70.8% of the images; and 4 workers are enough for 97.3% images.

6.3 Results of the demographic analysis

We annotated demographics on the 139 synsets that are considered both safe (Sec. 4) and imageable (Sec. 5) and that contain at least 100 images. We annotated 100 randomly sampled images from each synset, summing up to 13,900 images. Due to the presence of multiple people in an image, each image may have more than one category for each attribute. We ended up with 43,897 attribute categories annotated (14,471 annotations for gender; 14,828 annotations for skin; and 14,598 annotations for age). This was the result of obtaining and consolidating 109,545 worker judgments.

Fig. 4 shows the distribution of categories for different synsets, which mirrors real-world biases. For gender, there are both male-dominated synsets and female-dominated synsets; but the overall pattern across all synsets reveals underrepresentation of female, as the blue area in Fig. 4 (Left) is significantly larger than the green area. Relatively few images are annotated with the *Unsure* category except a few interesting outliers: *birth* (58.2% images labeled *Unsure*) and *scuba diver* (16.5%). The gender cues in these synsets obscured because *birth* contain images of newborn babies, and *scuba diver* contains people wearing diving suits and helmets.

The figure for skin color (Fig. 4 Middle) also presents a biased distribution, highlighting the underrepresentation of people with dark skin. The average percentage of the *Dark* category across all synsets is only 6.2%, and the synsets with significant portion of *Dark* align with stereotypes: *rapper* (66.4% images labeled *Dark*) and *basketball player* (34.5%). An exception is *first lady* (51.9%), as most images in this synset are photos of Michelle Obama, the First Lady of the United States when ImageNet was being constructed.

6.4 Limitations of demographic annotation

Given the demographic analysis, it is desired to have a constructive solution to improve the diversity in ImageNet images. Publicly releasing the collected attribute annotations would be a natural next step. This would allow the research community to train and benchmark machine learning algorithms on different demographic subsets of ImageNet, furthering the work on machine fairness.

However, we have to consider that the potential mistakes in demographics annotations are harmful not just for the downstream visual recognition models (as all annotation mistakes are) but to the people depicted in the photos. Mis-annotating gender, skin color, or age can all cause significant distress to the photographed subject. Gender identity and gender expression may not be aligned (similarly for skin color or age), and thus some annotations may be incorrect despite our best quality control efforts. So releasing the image-level annotations may not be appropriate in this context.

6.5 Methodology for increasing image diversity

We aim for an alternative constructive solution, one that strikes a balance between advancing the community’s efforts and preventing additional harm to the people in the photos. One option we considered is internally using the demographics for targeted data collection, where we would find and annotate additional images to re-balance each synset. However, with the known issues of bias

in search engine results [55] and the care already taken by the ImageNet team to diversify the images for each synset (Sec. 3), this may not be the most fruitful route.

Instead, we propose to release a Web interface that automatically re-balances the image distribution within each synset, aiming for a target distribution of a single attribute (e.g., gender) by removing the images of the overrepresented categories. There are two questions to consider: first, what is an appropriate target distribution, and second, what are the privacy implications of such balancing.

First, identifying the appropriate target distribution is challenging and we leave that to the end users of the database. For example, for some applications it might make sense to produce a uniform gender distribution, for example, if the goal is to train an activity recognition model with approximately equal error rates across genders. In other cases, the goal might be to re-balance the data to better mimic the real-world distribution of gender, race or age in the category (as recorded by census data for example) instead of using the distribution exaggerated by search engines. Note that any type of balancing is only feasible on synsets with sufficient representation within each attribute category. For example, the synset *baby* naturally does not contain a balanced age distribution. Thus, we allow the user to request a subset of the categories to be balanced; for example, the user can impose equal representation of the three adult categories while eliminating the *Child* category.

Second, with regard to privacy, there is a concern that the user may be able to use this interface to infer the demographics of the *removed* images. For example, it would be possible to visually analyze a synset, note that the majority of people within the synset appear to be female, and thus infer that any image removed during the gender-balancing process are annotated as female. To mitigate this concern, we always only include 90% of images from the minority category in the balanced images and discard the other 10%. Further, we only return a balanced distribution of images if at least 2 attribute categories are requested (e.g., the user cannot request a female-only gender distribution) and if there are at least 10 images within each requested category.

While we only balance the distribution of a single attribute (e.g., gender), it is desirable to balance across multiple attributes. However, it will result in too few images per synset after re-balancing. For example, if we attempt to balance both skin color and gender, we will end up with very few images. This creates potential privacy concerns with regard to being able to infer the demographic information of the people in the individual photos.

6.6 Results and estimated impact of the demographic balancing on ImageNet

Fig. 5 provides one example of the effect of our proposed demographic balancing procedure on the synset *programmer*. Based on our analysis and statistics so far, and under the restrictions described in Sec. 6.5, we could offer such a balancing on 131 synsets for gender (ignoring the highly skewed *Unsure* category and posing uniform distribution among *Male* and *Female*), 117 synsets for skin color (uniform distribution for the three categories), and 81 synsets for age (removing the *Child* category and posing a uniform



Figure 5: The distribution of images in the ImageNet synset programmer before and after balancing to a uniform distribution.

distribution for the other three age categories). Users can create customized balancing results for each synset by choosing the attribute categories to balance on.

6.7 Limitations of the balancing solution

The downside of this solution is that balancing the dataset instead of releasing the image-level attribute annotations makes it impossible to evaluate the error rates of machine learning algorithms on demographic subsets of the data, as is common in the literature [10, 36, 82, 83]. Nevertheless, this strategy is a better alternative than using the existing ImageNet person subtree (strong bias), releasing the image-level annotations (ethically problematic), or collecting additional images (technically impractical).

7 DISCUSSION

We took the first steps towards filtering and balancing the distribution of the person subtree in the ImageNet hierarchy. The task was daunting, as with each further step of annotation and exploration, we discovered deeper issues that remain unaddressed. However, we feel that this is a significant leap forward from the current state of ImageNet. We demonstrate that at most 158 out of the 2,832 existing synsets should remain in the person subtree, as others are inappropriate categories for visual recognition and should be filtered out. Of the remaining synsets, 139 have sufficient data (at least 100 images) to warrant further exploration. On those, we provide a detailed analysis of the gender, skin color and age distribution of the corresponding images, and recommend procedures for better balancing this distribution.

While 139 categories may seem small in comparison to the current set, it is nevertheless sufficiently large-scale to remain interesting to the computer vision community: e.g., the PASCAL dataset has only 20 classes [26], CelebA has 40 attributes [50], COCO has

80 object categories [48], the fine-grained CUB-200 dataset has 200 bird species [77]. Further, note that the most commonly used subset of ImageNet is the set of 1,000 categories in the ImageNet Large Scale Visual Recognition Challenge (ILSVRC) [64], which remains unaffected by our filtering: the three ILSVRC synsets of the person subtree are bridegroom (n10148035; safe, imageability 5.0), ballplayer (n09835506; safe, imageability 4.6) and scuba diver (n09835506; safe, imageability 5.0).

There is still much remaining to be done outside the person subtree, as incidental people occur in photographs in other ImageNet synsets as well, e.g., in synsets of pets, household objects, or sports. It is likely that the density and scope of the problem is smaller in other subtrees than within this one, so the filtering process should be simpler and more efficient. We are releasing our annotation interfaces to allow the community to continue this work.

ACKNOWLEDGMENTS

Thank you to Arvind Narayanan and Timnit Gebru for thoughtful discussions, to the graduate students who annotated the unsafe synsets, and to the Amazon Mechanical Turk workers for the imageability and demographic annotations. This work is supported by the National Science Foundation under Grant No. 1763642. The project began in fall 2018, the ImageNet site was temporarily placed under maintenance in January 2019, and the work was submitted for peer review in August 2019. This research was initiated, conceptualized, and executed solely by the authors.

A APPENDIX

We include the annotation interfaces and additional results as promised in the main paper. We organize the the appendix according to the sections of the main paper for ease of reference.

A.1 PROBLEM 1: STAGNANT CONCEPT VOCABULARY

The instructions used in-house to annotate the offensiveness of synsets are shown in Fig. A. We attach the synset IDs of the “unsafe” and “safe” synsets we have annotated. As before, we avoid explicitly naming the synsets, but the conversion from synset IDs to names can be found at wordnet.princeton.edu/documentation/wndb5wn.

Offensive synsets (1,593 in total).

image-net.org/filtering-and-balancing/unsafe_synsets.txt

Safe synsets (1,239 in total).

image-net.org/filtering-and-balancing/safe_synsets.txt

A.2 PROBLEM 2: NON-VISUAL CONCEPTS

Instructions. Fig. C shows the user interface for crowdsourcing imageability scores.

Quality control. Table A lists the gold standard questions for quality control; half of them are obviously imageable (should be rated 5), and the other half are obviously non-imageable (should be rated 1). For a worker who answered a set of gold standard questions Q , we calculate the root mean square error of the worker as:

$$Error = \sqrt{\frac{1}{|Q|} \sum_{i \in Q} (\hat{x}_i - x_i)^2} \quad (2)$$

where \hat{x}_i is the rating from the worker and x_i is the ground truth imageability for question i ($\hat{x}_i \in \{1, 2, 3, 4, 5\}, x_i \in \{1, 5\}$). If $Error \geq 2.0$, we exclude all ratings of the worker.

Even after removing the answers from high-error workers, the raw ratings can still be noisy, which is partially attributed to the intrinsic subjectiveness in the imageability of synsets. We average multiple workers’ ratings for each synset to compute a stable estimate of the imageability. However, it is tricky to determine the number of ratings to collect for a synset [68]; more ratings lead to a more stable estimate but cost more. Further, the optimal number of ratings may be synset-dependent; more ambiguous synsets need a larger number of ratings. We devise a heuristic algorithm to determine the number of ratings dynamically for each synset.

Intuitively, the algorithm estimates a Gaussian distribution using the existing ratings, and terminates when three consecutive new ratings fall into a high-probability region of the Gaussian. It automatically adapts to ambiguous synsets by collecting more ratings. Concretely, abusing notation from above (for simplicity of exposition), let $\hat{\mathbf{x}} = [\hat{x}_1, \hat{x}_2, \hat{x}_3, \dots, \hat{x}_m]$ now be the sequence of ratings for a single synset from workers 1, 2, 3, \dots m . After collecting $m \geq 4$ ratings, we partition the sequence into the last 3 ratings $\hat{\mathbf{x}}_{new} = [\hat{x}_{m-2}, \hat{x}_{m-1}, \hat{x}_m]$ and the rest $\hat{\mathbf{x}}_{old} = [\hat{x}_1, \hat{x}_2, \dots, \hat{x}_{m-3}]$. We compute the mean and standard deviation of $\hat{\mathbf{x}}_{old}$ as μ_{old} and σ_{old} , and we check whether the following holds:

$$\forall x \in \hat{\mathbf{x}}_{new}, \mu_{old} - \sigma_{old} \leq x \leq \mu_{old} + \sigma_{old} \quad (3)$$

Table A: Gold standard questions for quality control in imageability annotation.

Synset ID	Synset	Ground truth imageability
n10101634	football player, footballer	5
n10605253	skier	5
n09834885	ballet master	5
n10366966	nurse	5
n10701644	tennis pro, professional tennis player	5
n09874725	bride	5
n10772092	weatherman, weather forecaster	5
n10536416	rock star	5
n09624168	male, male person	5
n10087434	fighter pilot	5
n10217208	irreligionist	1
n10743356	Utopian	1
n09848110	theist	1
n09755788	abecedarian	1
n09794668	animist	1
n09778927	agnostic	1
n10355142	neutral	1
n10344774	namer	1
n09789898	analogist	1
n10007787	delegate	1

If it holds, the imageability annotations are approximately converging and we compute the final imageability score as the average of all ratings. Otherwise we collect more ratings and check again.

Fig. B shows the number of ratings collected for the synsets. The average number is 8.8, and 72% synsets need no more than 10 ratings. The file image-net.org/filtering-and-balancing/imageability_scores.txt includes the complete list of imageability scores for the 1,239 safe synsets in the person subtree.

A.3 PROBLEM 3: LACK OF IMAGE DIVERSITY

The user interface used to annotate the protected attributes is shown in Fig. D.

REFERENCES

- [1] U.S. House. 101st Congress, 2nd Session. 101 H.R. 2273. 1990. Americans with Disabilities Act of 1990. Washington: Government Printing Office.
- [2] U.S. House. 88th Congress, 1st Session. 88 H.R. 6060. 1963. Equal Pay Act of 1963. Washington: Government Printing Office.
- [3] U.S. House. 98th Congress, 2nd Session. 98 H.R. 5490. 1984. Civil Rights Act of 1984. Washington: Government Printing Office.
- [4] Giorgia Aiello and Anna Woodhouse. 2016. When corporations come to define the visual politics of gender. *Journal of Language and Politics* 15, 3 (2016), 351–366.
- [5] Jeanette Altarriba, Lisa M Bauer, and Claudia Benvenuto. 1999. Concreteness, context availability, and imageability ratings and word associations for abstract, concrete, and emotion words. *Behavior Research Methods, Instruments, & Computers* 31, 4 (1999), 578–602.
- [6] Solon Barocas, Elizabeth Bradley, Vasant Honavar, and Foster Provost. 2017. Big data, data science, and civil rights. *arXiv preprint arXiv:1706.03102* (2017).
- [7] Solon Barocas and Andrew D Selbst. 2016. Big data’s disparate impact. *Calif. L. Rev.* 104 (2016), 671.
- [8] Kristy Beers Fägersten. 2007. *A sociolinguistic analysis of swear word offensiveness*. Universität des Saarlands.
- [9] Helen Bird, Sue Franklin, and David Howard. 2001. Age of acquisition and imageability ratings for a large set of words, including verbs and function words. *Behavior Research Methods, Instruments, & Computers* 33, 1 (2001), 73–79.
- [10] Joy Buolamwini and Timnit Gebru. 2018. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on fairness, accountability and transparency*. 77–91.
- [11] Kaylee Burns, Lisa Anne Hendricks, Kate Saenko, Trevor Darrell, and Anna Rohrbach. 2018. Women also snowboard: Overcoming bias in captioning models. In *European Conference on Computer Vision*.

Instructions:

We are aiming to label the synsets in the person subtree of ImageNet as “safe” (label of 2) or “sensitive” (label of 1), “offensive” (label of 0). Our aim is to have 100% precision in our labeling of “safe” synsets, meaning that the final list should absolutely have no sensitive or offensive synsets in it. So if you are on the fence about a synset, err on the side of caution and label it sensitive or offensive.

Sensitive synsets (label of 1) are not inherently offensive, but can cause offense if applied inappropriately.

Examples of a sensitive synset:

- classification of people based on culture, race, ethnicity, nationality, E.g. “Russian”, “Asian”
- classification based on politics or religion: “Republican”, “Christian”.
- classification based on sex, gender, or sexual orientation: “Woman”, “Transgender”
- synsets with negative connotations (“thief”, “underling”, “plutocrat”), i.e. you wouldn’t like it if the label is applied to you.

Offensive synsets (label of 0) are inherently offensive.

Examples of an offensive synset:

- racial slurs of any sort
- profanity

In general, if you can think of any reason why it would be bad to include a synset, mark it as sensitive or offensive.

Imageability: here you do NOT need to worry about whether a label is “imageable”, i.e. whether it is possible to predict the label from just pixels. Just assume all labels are imageable.

Please make sure to read the definition of a synset, because the meaning may not be the one you assume. E.g. “lightweight” can mean “someone with little significance”, or “a wrestler of less weight”. The former is unsafe but the latter is safe.

For your given sheet, you should see three columns corresponding to a synset id, the word, and the glossary and definition of the word, respectively. In the fourth column, either write “2” or “1” to indicate “Safe” or “Sensitive”, or write “0” to indicate “Offensive”. *The labels column is initialized to all 0’s.*

Example of safe synset:

n10252547	lecturer	someone who lectures professionally	
-----------	----------	-------------------------------------	--

n10369317	oboist	a musician who plays the oboe	
-----------	--------	-------------------------------	--

Example of sensitive synset:

n09727440	Filipino	a native or inhabitant of the Philippines	
-----------	----------	---	--

n10519494	religious leader	leader of a religious order	
-----------	------------------	-----------------------------	--

Example of offensive synset:

n10401204	parricide	someone who kills his or her parent	
-----------	-----------	-------------------------------------	--

n10722965	traitor, treasonist	someone who betrays his country by committing treason	
-----------	---------------------	---	--

Figure A: The instructions for annotating the offensiveness of synsets. The annotation was done in-house rather than using crowdsourcing, thus the user interface is kept simple.

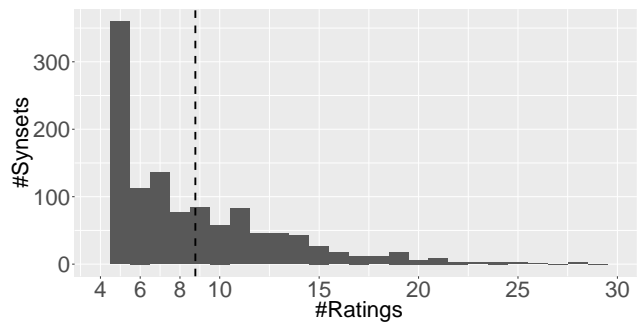


Figure B: The distribution of the number of raw imageability ratings collected for each synset. On average, the final imageability score of a synset is an average of 8.8 ratings.

[12] Toon Calders, Faisal Kamiran, and Mykola Pechenizkiy. 2009. Building classifiers with independency constraints. In *2009 IEEE International Conference on Data Mining Workshops*. IEEE, 13–18.

[13] Toon Calders and Sicco Verwer. 2010. Three naive Bayes approaches for discrimination-free classification. *Data Mining and Knowledge Discovery* 21, 2 (2010), 277–292.

[14] L Elisa Celis and Vijay Keswani. 2019. Implicit Diversity in Image Summarization. *arXiv preprint arXiv:1901.10265* (2019).

[15] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. 2015. Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012* (2015).

[16] A Chardon, I Cretois, and C Hourseau. 1991. Skin colour typology and suntanning pathways. *International journal of cosmetic science* 13, 4 (1991), 191–208.

[17] Michael J Cortese and April Fugett. 2004. Imageability ratings for 3,000 monosyllabic words. *Behavior Research Methods, Instruments, & Computers* 36, 3 (2004), 384–387.

[18] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. Imagenet: A large-scale hierarchical image database. In *IEEE conference on computer vision and pattern recognition*.

[19] Terrance DeVries, Ishan Misra, Changhan Wang, and Laurens van der Maaten. 2019. Does Object Recognition Work for Everyone? *arXiv preprint arXiv:1906.02659* (2019).

[20] Jean-Marc Dewaele. 2016. Thirty shades of offensiveness: L1 and LX English users’ understanding, perception and self-reported use of negative emotion-laden words. *Journal of Pragmatics* 94 (2016), 112–127.

[21] Djellel Difallah, Elena Filatova, and Panos Ipeirotis. 2018. Demographics and dynamics of mechanical Turk workers. In *Proceedings of the eleventh acm international conference on web search and data mining*. ACM, 135–143.

[22] Chris Dulhanty and Alexander Wong. 2019. Auditing ImageNet: Towards a Model-driven Framework for Annotating Demographic Attributes of Large-Scale Image Datasets. *arXiv preprint arXiv:1905.01347* (2019).

[23] Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. 2012. Fairness through awareness. In *Proceedings of the 3rd innovations in theoretical computer science conference*. ACM, 214–226.

[24] Harrison Edwards and Amos Storkey. 2016. Censoring representations with an adversary. In *ICLR*.

[25] Eran Eidinger, Roei Enbar, and Tal Hassner. 2014. Age and gender estimation of unfiltered faces. *IEEE Transactions on Information Forensics and Security* 9, 12

Given a word. How easy is it to form an image (in your mind) of the word?

Examples of "very easy" (select 5):

- police woman
- ballet dancer
- swimmer

Examples of "very hard" (select 1):

- liar
- perfectionist
- atheist

Some words do not fall into these two categories. For these words, select a score between 2 and 4 using your best judgement.

Examples:

- professor (Some features may be shared among many professors, but different professors can also look very different.)

In the rare case that you cannot understand a given word, please Google it.

sorcerer, magician, wizard, necromancer, thaumaturge, thaumaturgist: one who practices magic or sorcery

1 - very hard 2 - somewhat hard 3 - medium 4 - somewhat easy 5 - very easy

nonparticipant: a person who does not participate

1 - very hard 2 - somewhat hard 3 - medium 4 - somewhat easy 5 - very easy

Figure C: User interface for crowdsourcing the imageability annotation.

- (2014), 2170–2179.
- [26] Mark Everingham, Luc Gool, Christopher K. Williams, John Winn, and Andrew Zisserman. 2010. The Pascal Visual Object Classes (VOC) Challenge. *Int. J. Comput. Vision* 88, 2 (June 2010), 303–338.
- [27] Michael Feldman, Sorelle A Friedler, John Moeller, Carlos Scheidegger, and Suresh Venkatasubramanian. 2015. Certifying and removing disparate impact. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 259–268.
- [28] Thomas B Fitzpatrick. 1988. The validity and practicality of sun-reactive skin types I through VI. *Archives of dermatology* 124, 6 (1988), 869–871.
- [29] Paul Frosh. 2001. Inside the image factory: stock photography and cultural production. *Media, Culture & Society* 23, 5 (2001), 625–646.
- [30] Paul Frosh. 2002. Rhetorics of the Overlooked: On the communicative modes of stock advertising images. *Journal of Consumer Culture* 2, 2 (2002), 171–196.
- [31] Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna M. Wallach, Hal Daumé III, and Kate Crawford. 2018. Datasheets for Datasets. In *Workshop on Fairness, Accountability, and Transparency in Machine Learning*.
- [32] Ken J Gilhooly and Robert H Logie. 1980. Age-of-acquisition, imagery, concreteness, familiarity, and ambiguity measures for 1,944 words. *Behavior research methods & instrumentation* 12, 4 (1980), 395–427.
- [33] Ross Girshick. 2015. Fast r-cnn. In *Proceedings of the IEEE international conference on computer vision*. 1440–1448.
- [34] Moritz Hardt, Eric Price, Nati Srebro, et al. 2016. Equality of opportunity in supervised learning. In *Advances in neural information processing systems*. 3315–3323.
- [35] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
- [36] Lisa Anne Hendricks, Kaylee Burns, Kate Saenko, Trevor Darrell, and Anna Rohrbach. 2018. Women also snowboard: Overcoming bias in captioning models. In *European Conference on Computer Vision*. Springer, 793–811.
- [37] Minyoung Huh, Pulkit Agrawal, and Alexei A Efros. 2016. What makes ImageNet good for transfer learning? *arXiv preprint arXiv:1608.08614* (2016).
- [38] Ben Hutchinson and Margaret Mitchell. 2019. 50 Years of Test (Un)fairness: Lessons for Machine Learning. In *ACM Conference on Fairness, Accountability and Transparency*.
- [39] Faisal Kamiran and Toon Calders. 2009. Classifying without discriminating. In *2009 2nd International Conference on Computer, Control and Communication*. IEEE, 1–6.
- [40] Toshihiro Kamishima, Shotaro Akaho, and Jun Sakuma. 2011. Fairness-aware learning through regularization approach. In *2011 IEEE 11th International Conference on Data Mining Workshops*. IEEE, 643–650.
- [41] Matthew Kay, Cynthia Matuszek, and Sean A Munson. 2015. Unequal representation and gender stereotypes in image search results for occupations. In *Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems*. ACM, 3819–3828.
- [42] Brendan F Klare, Ben Klein, Emma Taborsky, Austin Blanton, Jordan Cheney, Kristen Allen, Patrick Grother, Alan Mah, and Anil K Jain. 2015. Pushing the frontiers of unconstrained face detection and recognition: Iarpa janus benchmark a. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 1931–1939.
- [43] Jon Kleinberg, Sendhil Mullainathan, and Manish Raghavan. 2017. Inherent trade-offs in the fair determination of risk scores. In *Proc. 8th Conf. on Innovations in Theoretical Computer Science (ITCS)*.
- [44] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A. Shamma, Michael S. Bernstein, and Li Fei-Fei. 2017. Visual Genome: Connecting Language and Vision Using Crowdsourced Dense Image Annotations. *International Journal of Computer Vision* 123, 1 (01 May 2017), 32–73.
- [45] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. 2012. ImageNet classification with deep convolutional neural networks. In *Advances in neural information processing systems*. 1097–1105.
- [46] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Mallocci, Tom Duerig, et al. 2018. The Open Images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. *arXiv preprint arXiv:1811.00982* (2018).
- [47] Sam Levin. 2016. A beauty contest was judged by AI and the robots didn't like dark skin.
- [48] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft COCO: Common objects in context. In *European conference on computer vision*.

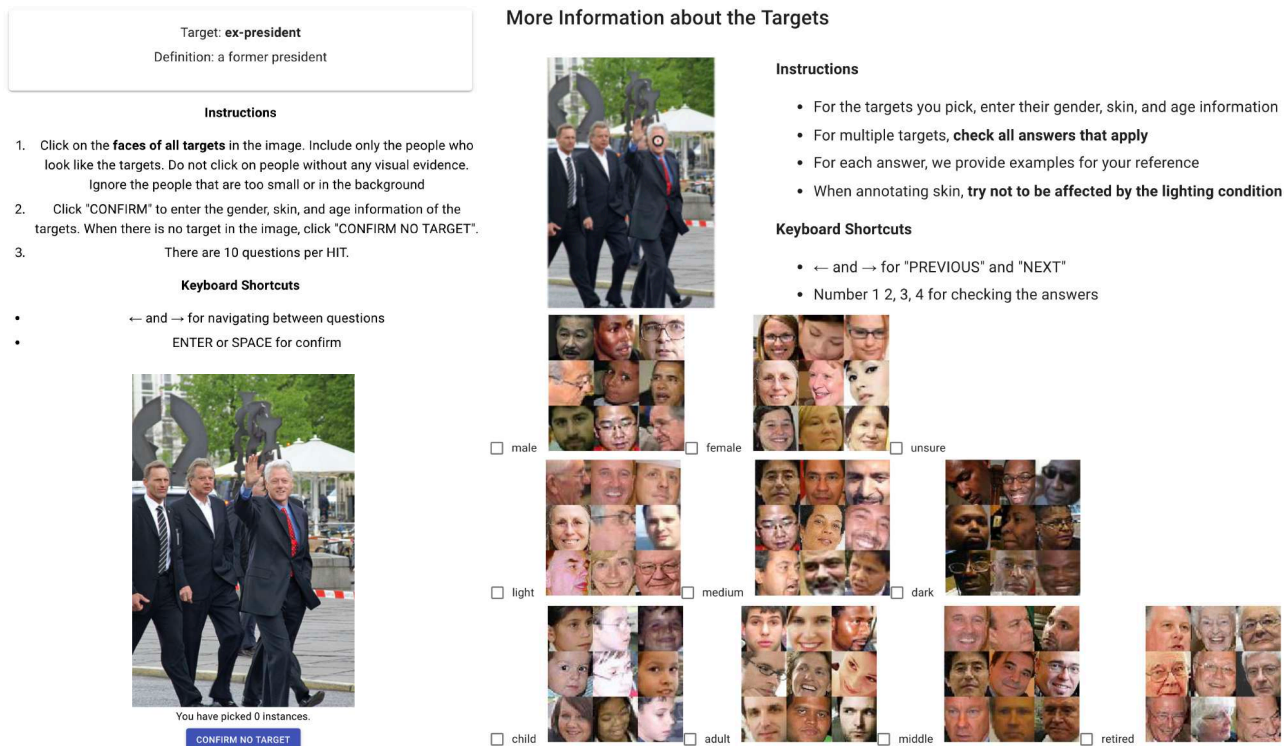


Figure D: User interface for crowdsourcing the demographics annotation.

[49] Fayao Liu, Chunhua Shen, and Guosheng Lin. 2015. Deep convolutional neural fields for depth estimation from a single image. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 5162–5170.

[50] Ziwei Liu, Ping Luo, Xiaoang Wang, and Xiaoou Tang. 2015. Deep Learning Face Attributes in the Wild. In *ICCV*.

[51] Ziwei Liu, Zhongqi Miao, Xiaohang Zhan, Jiayun Wang, Boqing Gong, and Stella X Yu. 2019. Large-Scale Long-Tailed Recognition in an Open World. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2537–2546.

[52] David Madras, Elliot Creager, Toniann Pitassi, and Richard Zemel. 2018. Learning adversarially fair and transferable representations. *arXiv preprint arXiv:1802.06309* (2018).

[53] George A Miller. 1998. *WordNet: An electronic lexical database*. MIT press.

[54] Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchinson, Elena Spitzer, Inioluwa Deborah Raji, and Timnit Gebru. 2018. Model Cards for Model Reporting. In *ACM Conference on Fairness, Accountability and Transparency*.

[55] Safiya Umoja Noble. 2018. *Algorithms of oppression: How search engines reinforce racism*. nyu Press.

[56] The Office of Communications (Ofcom). 2016. *Attitudes To Potentially Offensive Language and Gestures on TV and Radio*. Technical Report. <https://www.ofcom.org.uk/research-and-data/tv-radio-and-on-demand/tv-research/offensive-language-2016>

[57] Allan Paivio. 1965. Abstractness, imagery, and meaningfulness in paired-associate learning. *Journal of Verbal Learning and Verbal Behavior* 4, 1 (1965), 32–38.

[58] Allan Paivio, John C Yuille, and Stephen A Madigan. 1968. Concreteness, imagery, and meaningfulness values for 925 nouns. *Journal of experimental psychology* 76, 1p2 (1968), 1.

[59] Dino Pedreshi, Salvatore Ruggieri, and Franco Turini. 2008. Discrimination-aware data mining. In *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 560–568.

[60] Geoff Pleiss, Manish Raghavan, Felix Wu, Jon Kleinberg, and Kilian Q Weinberger. 2017. On fairness and calibration. In *Advances in Neural Information Processing Systems*. 5680–5689.

[61] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. 2016. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 779–788.

[62] Lauren Rhue. 2018. Racial Influence on Automated Perceptions of Emotions. (2018). <https://ssrn.com/abstract=3281765>

[63] Joel Ross, Andrew Zaldivar, Lilly Irani, and Bill Tomlinson. 2009. Who are the turkers? worker demographics in amazon mechanical turk. *Department of Informatics, University of California, Irvine, USA, Tech. Rep* (2009).

[64] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. 2015. ImageNet Large Scale Visual Recognition Challenge. *International journal of computer vision* 115, 3 (2015), 211–252.

[65] Hee Jung Ryu, Hartwig Adam, and Margaret Mitchell. 2018. Inclusivefacenet: Improving face attribute detection with race and gender diversity. In *Proceedings of FATML*.

[66] Barry S Sapolsky, Daniel M Shafer, and Barbara K Kaye. 2010. Rating offensive words in three television program contexts. *Mass Communication and Society* 14, 1 (2010), 45–70.

[67] Shreya Shankar, Yoni Halpern, Eric Breck, James Atwood, Jimbo Wilson, and D Sculley. 2017. No classification without representation: Assessing geodiversity issues in open data sets for the developing world. In *NeurIPS workshop: Machine Learning for the Developing World*.

[68] Victor S Sheng, Foster Provost, and Panagiotis G Ipeirotis. 2008. Get another label? improving data quality and data mining using multiple, noisy labelers. In *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 614–622.

[69] Karen Simonyan and Andrew Zisserman. 2014. Two-stream convolutional networks for action recognition in videos. In *Advances in neural information processing systems*. 568–576.

[70] Karen Simonyan and Andrew Zisserman. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014).

[71] Pierre Stock and Moustapha Cisse. 2018. Convnets and ImageNet beyond accuracy: Understanding mistakes and uncovering biases. In *Proceedings of the European Conference on Computer Vision (ECCV)*. 498–512.

[72] Hirotugu Takiwaki et al. 1998. Measurement of skin color: practical application and theoretical considerations. *Journal of Medical Investigation* 44 (1998), 121–126.

[73] Antonio Torralba, Alexei A Efros, et al. 2011. Unbiased look at dataset bias.. In *CVPR*, Vol. 1. Citeseer, 7.

[74] Antonio Torralba, Rob Fergus, and William T. Freeman. 2008. 80 Million Tiny Images: A Large Data Set for Nonparametric Object and Scene Recognition. *IEEE*

- Trans. Pattern Anal. Mach. Intell.* 30, 11 (Nov. 2008), 1958–1970.
- [75] Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. 2015. Show and tell: A neural image caption generator. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 3156–3164.
- [76] Yilun Wang and Michal Kosinski. 2018. Deep neural networks are more accurate than humans at detecting sexual orientation from facial images. *Journal of Personality and Social Psychology (JPSP)* (2018).
- [77] P. Welinder, S. Branson, T. Mita, C. Wah, F. Schroff, S. Belongie, and P. Perona. 2010. *Caltech-UCSD Birds 200*. Technical Report CNS-TR-2010-001. California Institute of Technology.
- [78] Meredith Whittaker, Kate Crawford, Roel Dobbe, Genevieve Fried, Elizabeth Kazianas, Varoon Mathur, Sarah Myers West, Rashida Richardson, Jason Schultz, and Oscar Schwartz. 2018. AI Now Report 2018. https://ainowinstitute.org/AI_Now_2018_Report.pdf.
- [79] Lydia TS Yee. 2017. Valence, arousal, familiarity, concreteness, and imageability ratings for 292 two-character Chinese nouns in Cantonese speakers in Hong Kong. *PloS one* 12, 3 (2017), e0174569.
- [80] Muhammad Bilal Zafar, Isabel Valera, Manuel Gomez Rodriguez, and Krishna P Gummadi. 2015. Fairness constraints: Mechanisms for fair classification. *arXiv preprint arXiv:1507.05259* (2015).
- [81] Rich Zemel, Yu Wu, Kevin Swersky, Toni Pitassi, and Cynthia Dwork. 2013. Learning fair representations. In *International Conference on Machine Learning*. 325–333.
- [82] Brian Hu Zhang, Blake Lemoine, and Margaret Mitchell. 2018. Mitigating unwanted biases with adversarial learning. In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*. ACM, 335–340.
- [83] Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2017. Men also like shopping: Reducing gender bias amplification using corpus-level constraints. In *EMNLP*.
- [84] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. 2017. Places: A 10 million Image Database for Scene Recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2017).