

Who's in and who's out? A case study of multimodal CLIP-filtering in DataComp

Rachel Hong
hongrach@cs.washington.edu
University of Washington
Seattle, Washington, USA

Tadayoshi Kohno
University of Washington
Seattle, Washington, USA

William Agnew
Carnegie Mellon University
Pittsburgh, Pennsylvania, USA

Jamie Morgenstern
University of Washington
Seattle, Washington, USA

ABSTRACT

As training datasets become increasingly drawn from unstructured, uncontrolled environments such as the web, researchers and industry practitioners have increasingly relied upon data filtering techniques to “filter out the noise” of web-scraped data. While datasets have been widely shown to reflect the biases and values of their creators, in this paper we contribute to an emerging body of research that assesses the filters used to create these datasets. We show that image-text data filtering also has biases and is value-laden, encoding specific notions of what is counted as “high-quality” data. In our work, we audit a standard approach of image-text CLIP-filtering on the academic benchmark DataComp’s CommonPool by analyzing discrepancies of filtering through various annotation techniques across multiple modalities of image, text, and website source. We find that data relating to several imputed demographic groups — such as LGBTQ+ people, older women, and younger men — are associated with higher rates of exclusion. We also find prevalence of *Western bias*, where the CLIP filter is more likely to include data related to Western countries compared to that of non-Western countries. Moreover, we demonstrate cases of *exclusion amplification*: not only are certain marginalized groups already underrepresented in the unfiltered data, but CLIP-filtering excludes data from these groups at higher rates. The data-filtering step in the machine learning pipeline can therefore exacerbate representation disparities already present in the data-gathering step, especially when existing filters are designed to optimize a specifically-chosen downstream performance metric like zero-shot image classification accuracy. Finally, we show that the NSFW filter fails to remove sexually-explicit content from CommonPool, and that CLIP-filtering includes several categories of copyrighted content at high rates. Our conclusions point to a need for fundamental changes in dataset creation and filtering practices. **Content warning: This paper discusses societal stereotypes and sexually-explicit material that may be disturbing, distressing, and/or offensive to the reader.**

CCS CONCEPTS

• Information systems → Data cleaning; • Social and professional topics → User characteristics.

KEYWORDS

Multimodal filtering, Dataset collection, CLIP, Data ideology, Representation disparities, Filtering bias

1 INTRODUCTION

Text and image datasets in machine learning have grown to the scale of billions [42, 100, 111] in order to build larger text and image generative models [2, 59]. To reach such magnitudes, datasets are being increasingly scraped from the web [49, 101], through archives such as Common Crawl [37]. While data from these web dumps are expansive, they are not expertly curated and may reduce the efficacy of downstream models. Manually assessing every data point is infeasible, however, so this requires at-scale methods to remove data that curators deem undesirable.

To address these issues, machine learning practitioners and researchers have increasingly begun to rely on automated data filtering to improve training efficiency and discard what they refer to as “low-quality” data. Specifically, for multimodal data, researchers have applied a pretrained CLIP model to filter image-text data obtained from Common Crawl. As depicted in Figure 1, this *CLIP-filtering* method (referring to the use of the OpenAI CLIP model [99]) assesses similarity between an image and its corresponding text within the model’s embedding space. This filtering technique was used to create the LAION-400M [112], LAION-5B [111], and DataComp-1B [49] datasets with sizes of 400 million, 5 billion, and 1 billion respectively, which were then used to train other open-source CLIP models. Because these models obtained state-of-the-art performance on several zero-shot image classification tasks, researchers have concluded CLIP-filtering to be the “most performant” method, as opposed to other filtering techniques that use caption length or image size [47, 49].

However, OpenAI CLIP is neither intended for filtering nor built for deployment [17]. As stated in their model card [93]: “Any deployed use case of the model — whether commercial or not — is currently out of scope.” Here, using CLIP as a filter to build a deployed model is by definition a *use* case of CLIP in a deployed setting. Yet these datasets obtained via CLIP-filtering have subsequently been used to train text-to-image models like Stable Diffusion and Midjourney, which each has tens of millions of users [2, 82]. Recent work has shown that these models exhibit problematic behaviors, such as amplifying demographic stereotypes or generating violent and sexually-explicit content [11]. Much of these behaviors is attributed to the content of the training data [98], as prior work demonstrates that the LAION datasets contain high rates of hateful content and misogynistic stereotypes [15].

At the same time, it is unclear what role CLIP-filtering plays on downstream text-to-image models or whether CLIP-filtering

simply replicates the model’s training dataset, especially given the demographic biases embedded in OpenAI CLIP itself [1]. While there have been speculations on the demographic biases of CLIP-filtering [17], there has not been any exploration to assess exactly how the filter changes the demographic makeup of the training dataset. As such, we choose to examine this initial stage of dataset curation because of its potential impact on downstream models.

1.1 Contributions

In our work, we perform an audit of the standard CLIP-filtering approach to the DataComp CommonPool dataset, focusing specifically on the use of the OpenAI CLIP model as a filter, as a case study of the role data filtering plays within the broader machine learning pipeline. To guide our analysis, we seek to answer the following question: **In what ways does CLIP-filtering exclude or include various types of data in DataComp CommonPool?** We focus on this dataset specifically because of its usage as an academic benchmark, although our findings may extend to LAION-400M and LAION-5B since they were created in a similar manner.

We break down our investigation to address the following three questions: First: **Which demographic groups are disproportionately excluded by CLIP-filtering?** Different demographic groups being filtered at different rates implies a form of representational harm [114]. A filter that systematically omits data from a certain group may therefore worsen the performance of the resulting model on that group. As a result, we aim to quantify this breakdown between unfiltered and filtered data across various demographic dimensions. The need to assess sociodemographic attributes across multiple modalities motivates our next question: **How can we measure sociodemographic attributes at scale within an image-text dataset?** We examine a broad range of factors like gender, age, race, religion, sexuality, language, or geographic region. While by no means comprehensive, we infer attributes that appear in the text and image components that can easily be related to people. In addition to examining filtering discrepancies by demographics, we also examine discrepancies by data source: **What types of websites are considered “high-quality” by the CLIP filter and what are the implications of certain websites being included?** To the extent that it is valuable to know whether datasets are globally representative, we ask whether data from certain geographic regions or time periods of the internet are treated differently by CLIP-filtering, as well as what types of websites have data that are more likely to pass through the filter.

Driven by the above questions, we conduct a novel examination of demographic bias as a consequence of multimodal filtering. In doing so, we make the following contributions:

- (1) *Filter discrepancies:* Across multiple modalities, different imputed sociodemographic groups are filtered at different rates. For example, CLIP-filtering is more likely to exclude data relating to LGBTQ+ identities and non-Western regions.
- (2) *Exclusion amplifier:* We find that not only does CLIP-filtering disproportionately exclude data from certain imputed demographic groups, but that this form of exclusion *amplifies* existing representation disparities. Some imputed groups that are already underrepresented in the original dataset are disproportionately filtered out at higher rates compared to overrepresented groups.
- (3) *Website inclusion:* Data from stock photo websites and U.S. and British news sites, which may be subject to copyright restrictions, are included at much higher rates than average. Large quantities of images from websites serving sexually-explicit material are also present after both NSFW and CLIP-filtering.
- (4) *Evaluation tools:* To address our research questions, we also develop novel methodologies to effectively audit an image-text dataset at scale. We make our analysis code available at <https://github.com/hongrachel/clip-filtering-bias/>.

Overall, this existing data filtering method is a clear example of a value-embedded machine learning practice [14]. From our results, we see that CLIP-filtering does in fact substantially change the makeup of the final dataset. This reinforces more broadly the idea from Gururangan et al. [55] that data filtering must encode a specific ideology of what constitutes as “high-quality data.” Prior works on image-text filtering assume that “quality” is often intrinsic to the data and that this “noise” to be weeded out is vaguely defined [130]. The LAION and DataComp authors demonstrate the efficacy of the CLIP filter by selecting specific downstream metrics for specific datasets such as zero-shot ImageNet classification accuracy, rather than the societal impacts of downstream models [49, 112]. Here, we investigate a case where filtering is designed around particular assumptions of what data quality and performance should be without additional justification. *Despite the supposed objectivity of these decisions, our findings prove that filtering embeds societal norms as to what should and should not be excluded.* Following our analysis, we further elaborate on these implications and give recommendations for designing data filters in the future.

2 CLIP-FILTERING

We describe the filtering process first proposed as a high-performing filtering method to curate the LAION datasets [111, 112] and subsequently studied and expanded upon in DataComp [49].

Broadly, the goal of the image-text filtering step is to determine whether a piece of text accurately describes its corresponding image — in other words, to obtain image-text alignment. However, the term “alignment” is not well-defined nor well-measured, and it is unclear how using CLIP fits with the search for “high-quality” data.

2.1 Machine learning pipeline

Figure 1 provides an overview of the machine learning pipeline to build the LAION and DataComp datasets, including the CLIP-filtering step, in order to train large image-text models on these datasets. In our work, we primarily follow the DataComp pipeline as each step is well-documented, although this process is easily extendible to the LAION pipeline.

2.1.1 Raw dataset collection. Image-text data is first scraped from a snapshot of the web, CommonCrawl [37], through parsing HTML image tags with nonempty alt-text tags. Alt-text, or alternative text, is a description associated with an image when the image cannot be rendered. Primarily, alternative text is intended for accessibility purposes, such as being presented to users with screen readers [28]. However, extensive research demonstrates that a large proportion of

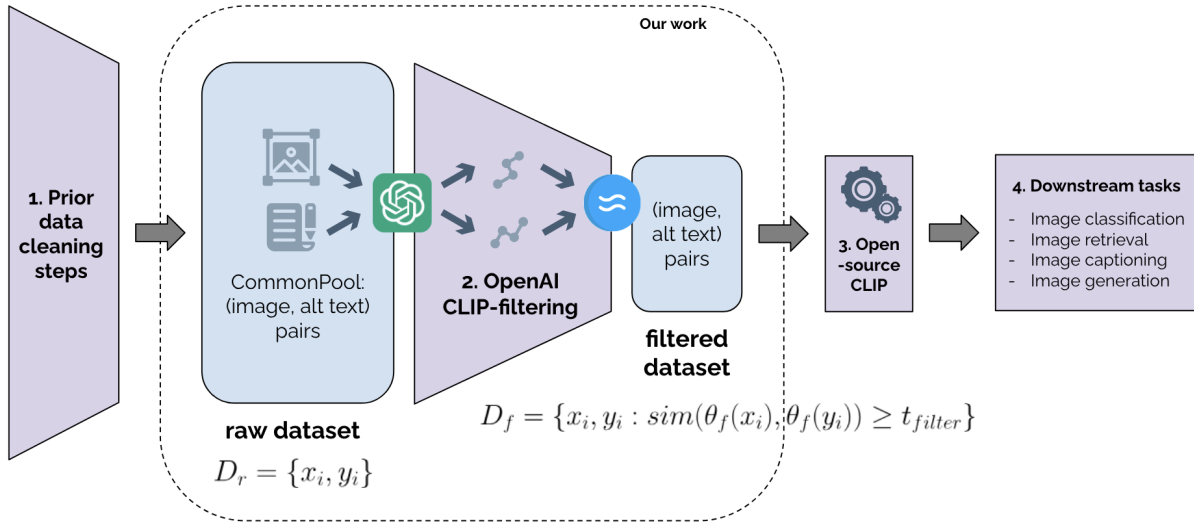


Figure 1: CLIP-filtering pipeline of LAION [111, 112] and DataComp [49] as described in Section 2.1. In step 1, a series of initial data cleaning techniques is applied to CommonCrawl [37] to form the raw dataset D_r with each sample as image x_i and corresponding alt-text tag y_i . The filtered dataset D_f is obtained by step 2, which applies the pre-trained OpenAI CLIP model θ_f to each image-text pair (x_i, y_i) . If the cosine similarity score between embeddings $\theta_f(x_i)$ and $\theta_f(y_i)$ is above some predefined threshold t_{filter} , then the pair is included in the filtered dataset. From this dataset, step 3 trains the open-source CLIP model, and step 4 applies the CLIP model to various downstream tasks. We scope our investigation to step 2 within the enclosed box.

websites do not follow standard accessibility guidelines for alt-text [12]. For example, among the one million most popular websites, 22.1% of home page images lack alt-text tags, and 10.9% of images with alternative text were non-descriptive including words like “text” or “blank” [129]. In addition, alt-text is a primary way search engines parse images, leading many alt-text tags being designed to improve website visibility on search engines rather than for accessibility or fidelity [53]. As a result, in order to build high-quality image-text datasets from the web, it becomes necessary to filter instances with questionable alt-text.

Once the dataset curators gather the image URLs and alt-text pairs, they attempt to download the images and keep the image-text pairs that are successfully downloaded. Some data cleaning methods are also applied here in the DataComp pipeline: they remove exact url-text duplicates within the dataset as well as potential image overlap with pre-selected evaluation sets. In addition, LAION removes short text and extremely small or large image sizes, while DataComp applies several toxicity classifiers to discard NSFW-detected images or text. This toxicity filtering is another key component of data filtering, and prior work has demonstrated that current hate speech detection models are more likely to label text written by African Americans as offensive [108]. However, we limit the scope of our work to the CLIP-filtering step (i.e., after the NSFW filter has been applied), although the toxicity filtering step may also disproportionately exclude data from certain identities.

After these preprocessing steps, the authors of DataComp name the subsequent dataset as **CommonPool** which consists of 12.8 billion pairs. For LAION-5B, this results in a dataset of 50 billion pairs before CLIP-filtering [111]. For the sake of clarity, we refer to

these image-text pairs before CLIP-filtering (CommonPool in the case of DataComp) as the **raw dataset**.

2.1.2 CLIP-filtering. After the creation of the raw dataset, the main filtering step uses a pre-trained CLIP model (ViT-B/32 or ViT-L/14) as follows: for every image-text pair, extract the CLIP image and text embeddings and obtain the cosine similarity score between these embeddings. The image-text pair passes the filter if the score is above a predefined threshold and discarded otherwise.

Each prior work uses differently-chosen similarity score thresholds — LAION-400M [112] has a threshold of 0.3, LAION-2B-en (the subset of LAION-5B with English-detected text) [111] has a threshold of 0.28, and DataComp has the threshold that separates the top 30% of data in CommonPool (0.281 for ViT-B/32 and 0.243 for ViT-L/14) [49]. Across all of these works, this discards a majority of the data from the raw dataset, and this threshold is chosen in order to optimize a specific set of downstream metrics. We refer to the downstream dataset obtained by applying the CLIP filter as the **filtered dataset**.

We note that for DataComp, the filtering task on the CommonPool benchmark is a track in their proposed competition in order encourage better filtering models [40]. CLIP-filtering is only one of the many filtering baseline experiments they conduct, but they find that their best baseline to improve their evaluation tasks, including ImageNet zero-shot accuracy, incorporates CLIP-filtering [49]. Thus, the CLIP-filtering step is a key component of the DataComp-1B dataset they release.

2.1.3 Pre-training. The filtered dataset then becomes the training dataset for any large model across a variety of tasks. The LAION

and DataComp researchers focus on training open-source CLIP models following the original CLIP paper [31, 99]. At the same time, these filtered datasets have also been used to train large text-to-image models like Stable Diffusion or Midjourney, despite OpenAI and LAION stating they are not ready to be used for commercial purposes [2, 82, 93]. Given that Gadre et al. [49] demonstrate that various filtering methods impact the resulting CLIP model’s performance, it becomes evident that filtering plays an important role in developing pre-trained models, whether for research or for deployed purposes. As a result, this CLIP-filtering approach has curated datasets to train models with a massive reach over tens of millions of users [2, 82].

2.1.4 Fine-tuning on downstream tasks. The final step before model deployment can include fine-tuning or training another model, such as a GAN [83], using the pre-trained CLIP model from the previous step. In the next section, we describe how CLIP can be used for various downstream tasks, such as image classification or image retrieval.

2.2 CLIP model

Contrastive Language-Image Pretraining, or *CLIP*, was proposed by Radford et al. [99] as a pre-trained model that learns images from natural language text through mapping images and texts to the same representation space. The training objective follows contrastive representation learning [124]: to maximize the cosine similarity between image embeddings and their corresponding text embeddings, while minimizing similarity between all other incorrect pairings. The original CLIP model was trained on WebImageText, which is an unreleased dataset of 400 million English-text pairs curated with 500,000 web queries. Subsequent work has followed similar model architecture to train open-source versions of CLIP on publicly-available datasets. Xu et al. [135] replicate the training data for the original CLIP model by reconstructing the query metadata, whereas Cherti et al. [31] train on the LAION dataset, which imitates WebImageText but still relies on the original model as a filter.

The authors of the original CLIP model [99] and following work demonstrate that CLIP obtains state-of-the-art performance on various vision tasks without additional training, such as zero-shot ImageNet classification, image retrieval [58, 107], as well as using the embeddings to guide image captioning [83] and text-to-image models [38].

3 RELATED WORK

In this section, we highlight previous work on bias audits and data filtering methods that are relevant to our study. We aim to contribute to the existing discussion on data filtering bias through a novel case study of data filtering on a popular large-scale dataset.

3.1 CLIP bias

Since the release of the OpenAI CLIP model, multiple follow-up investigations have found that CLIP encodes harmful stereotypes and sociodemographic biases which may have been present in the training data. Notably, Agarwal et al. [1] find that CLIP misclassifies images of Black people as non-human at higher rates compared to images of other racial groups, as well as images of men as ‘executive’

or ‘doctor’ at higher rates compared to images of women. Wolfe and Caliskan [133] also find that CLIP embeddings of images of white people are more associated with being “American.” Other works evaluate CLIP-aided image retrieval and text-to-image models and demonstrate representation bias as well as amplification of stereotypes [4, 11]. As a result, we hypothesize that CLIP-filtering may further overrepresent content containing stereotypes or harmful biases and disproportionately exclude content that do not follow these problematic associations.

3.2 CLIP-filtered dataset audits

Some recent works have audited the LAION datasets in order to understand the nature of the content within these large-scale datasets. Recently, Thiel [120] uses hash-based detection and other computational techniques to identify thousands of CSAM (Child Sexual Abuse Material) images in LAION-5B, which resulted in the takedown of the LAION-5B dataset in December, 2023 [126]. In addition, Birhane et al. [15, 17] have found significant rates of hateful content, explicit images, stereotypes, and racial slurs in both the LAION-400M and LAION-2B-en datasets. Given the biases of CLIP, they illustrate hypothetical examples in which misogynist or racist descriptions would pass the CLIP filter but reasonable, benign descriptions would be excluded [17]. The question then remains: What, exactly, does CLIP-filtering include and exclude? We extend their findings through examining the impact of filtering on demographic representation at scale.

3.3 Data filtering bias

In NLP, filtering is a common practice to build large text datasets scraped from the internet. For instance, to curate the training dataset for the GPT3 model, researchers trained a quality filter to identify high-quality sources like Wikipedia [24]. Gururangan et al. [55] find that the GPT3 quality filter is more likely to classify text from wealthier, urban, and larger schools as high-quality. Dodge et al. [44] also examine the blocklist filter used to create the popular C4 dataset [101] and determine that mentions of identities from marginalized groups are more likely to be filtered out. Lucy et al. [77] investigate many English language and quality text filters applied to descriptions of website creators and show disparate rates of filtering based on topic and geographic region. For tabular data, common data preprocessing steps are found to remove data from historically disadvantaged groups and worsen downstream model fairness [18, 52].

There also has been much discussion on the systematic exclusion of data from marginalized identities in the data cleaning stage more broadly. Bender et al. [10] caution early on that filtering can suppress discourse from marginalized identities and compound power imbalances in text data collection practices. From a sociotechnical lens, Muller and Strohmayer [86] describe how data cleaning can be a method of erasure that is difficult to detect. Desai et al. [41] draw on an archival perspective to examine the social value judgements embedded in filtering choices. We are inspired by prior filtering discussions and extend their approaches to a commonly used multimodal data pipeline. To our knowledge, there has been no work assessing data filtering bias at scale in the context of image-text datasets.

3.4 Multimodal data filtering

Researchers have also investigated image-text data filtering in relation to performance or robustness to distribution shifts. Nguyen et al. [88] demonstrate that filtering noisy data with a pre-trained robust model like CLIP can lead to more robust models downstream, which may explain the state-of-the-art image classification accuracy by training on the LAION and DataComp datasets. Fang et al. [47] complicate this result by showing that using a higher-performing CLIP model as a filter counterintuitively does not always lead to better downstream models. While these findings reveal the intricacies of building effective multimodal filters, the ongoing work so far has not yet focused on societal implications.

3.5 Downstream text-to-image model harms

There has been ample work that evaluate the societal issues of downstream models trained on CLIP-filtered data. Luccioni et al. [74], for instance, find that Stable Diffusion (which trains on LAION-5B obtained via CLIP-filtering [105]), underrepresents marginalized identities in their image generations. Moreover, researchers have demonstrated that Stable Diffusion generates unsafe, hateful, or sexualized content at significantly high rates [98, 134]. Exposure to stereotypical and problematic imagery in general has been shown to shape people’s beliefs and behavior [23, 118], and this becomes more pressing as these models grow increasingly popular. As such, it thus becomes necessary to understand exactly where societal harms may arise in the ML pipeline.

4 APPROACH

We first replicate the CLIP-filtering step for the small version of CommonPool of 12.8 million image-text pairs. We use the CLIP ViT-L/14 model and set the threshold as 0.243 to select the top 30% of pairs according to CLIP similarity score as done in prior work [49, 111], resulting in a filtered dataset with 3.84 million samples. In this section, we describe how we obtain demographic annotations and measure filtering discrepancies.

4.1 Imputed demographic group annotations

Many large datasets, including multimodal datasets like CommonPool, contain little or no metadata about the people used in generating their data, which can include people in images, the authors of text, or patients whose medical records comprise a dataset. This makes auditing datasets and their downstream uses for disparate treatment along certain characteristics, like gender, race, and age, a challenging task. Because CommonPool does not contain demographic metadata, we cannot directly audit CLIP-filtering along these axes. Instead, we attempt to audit CLIP-filtering’s interaction with *imputed* demographic information, using existing models to evaluate filtering bias.

We determine demographic group or geographic region across multiple modalities in the CommonPool dataset: from the text content, the image content, and the source URL of the sample. As described in Section 5, this includes methods such as searching for mentions of keywords relating to demographic identity or applying an external face attribution predictor. Due to the sample size considered, we separate our analysis by modality, despite a sample containing both image and text, and leave the image-text

intersection for follow-up work. We emphasize these findings *do not represent behavior of filtering according to true demographics*; instead, they reveal how filtering interacts with imputed demographic attributes. While we elaborate on the limitations of imputed demographics in Section 4.3, findings here still have broad implications, as imputed demographics can correlate, although imperfectly, with socially salient populations [46].

4.2 Evaluation metrics

To assess the impact of filtering, we track the **pass rates** by imputed demographic group, which refer to the percentage of raw dataset that passes through the filter. In other words, a *higher pass rate* means more data is kept in the filtered dataset, and a *lower pass rate* means more data is excluded from the filtered dataset. This allows us to evaluate whether CLIP-filtering leads to representational harms in the filtered dataset if data relating to certain imputed demographic groups are considered “lower quality” on average.

4.3 Limitations

We recognize numerous limitations to the sociodemographic imputation and source analysis techniques that we use. Third-party gender predictors, for instance, have been demonstrated to have disparate error rates across sensitive attributes and can rely on biases present in the face detection step [25, 113]. Furthermore, these predictors only give binary labels as “Male” and “Female” which conflates gender with sex and ignores non-binary individuals [125]. Identity keywords in the text also are rough proxies for the presence of sociodemographic attributes in the image or text. We acknowledge that gender and race are not inherently visual components, but rather socioculturally defined [26, 73].

Evaluating filtering according to these imputed demographic groups is a form of measuring filtering’s group fairness properties. Group fairness is not sufficient especially when individuals belong to multiple groups or their group categories are not explicitly known [45]. At the same time, we recognize the difficulties of assessing sociodemographic information at scale – we believe that determining relationships between pass rates and these signals, although problematic or noisy, reveals behavior of the CLIP filter and the content of data in the LAION and DataComp-1B datasets.

Moreover, while any difference in pass rates that we find can be attributed to the impact of CLIP-filtering, it does not necessarily demonstrate a causal relationship that a certain demographic group will always be excluded by CLIP. There may be underlying differences in the raw dataset, where data relating to a particular group may come from websites that have established alt-text practices, yet this is infeasible to assess at scale across CommonPool. We argue regardless that disparities in filtering still matter in order to expose whether certain imputed demographic groups are being omitted in the construction of large-scale machine learning datasets. Thus, it still remains important to understand the impact filtering has on the resulting dataset given that these CLIP-filtered datasets like LAION and DataComp are used widely in both industry and research [2, 49, 82].

Table 1: List of regular expressions relating to identity keywords used in Dodge et al. [44].

african([-]americans)?	asian([-]american)?s?
bi-?sexuals?	blacks?
caucasians?	christians?
european([-]american)?s?	females?
gays?	heterosexuals?
homosexuals?	jew(s ish)?
latin[oax]s?	lesbians?
m[ae]n	males?
muslims?	non[-]?binary
straights?	trans(\+ gender)
whites?	wom[ae]n

5 DEMOGRAPHIC GROUP ANALYSIS

This section describes the demographic group imputation methods along both text and image modalities and presents their corresponding results. We include additional analysis in Appendix B.

5.1 Identity keyword

Method: We first identify mentions of demographic groups through simple keyword search, by following the same methodology as Dodge et al. [44], which focuses on blacklist filtering of text. Their list of regular expressions (shown in Table 1) investigates demographic dimensions related to gender, sexual orientation, race, and religion — while no means comprehensive, this initial exploration of filtering analysis allows us to compare our results to their findings on the blacklist filter for the C4 text dataset [101]. Next, we plot the pass rates (rate of inclusion in the filtered dataset) grouped by mentions of identity keywords in the text modality. Upon manual inspection, we dismiss analysis of the white and black keywords since text with these keywords typically describe clothing apparel rather than people.

5.1.1 The CLIP filter excludes data relating to LGBTQ+ identities at higher rates. In Figure 2, we see similar trends as Dodge et al. [44] found on text filtering, where mentions of keywords relating to LGBTQ+ identities (homosexual, lesbian, transgender, gay, bisexual) are excluded at much higher rates compared to mentions of other keywords. We extend existing findings on text to the multimodal setting: Cleaning methods often remove data relating to LGBTQ+ identities, although prior work specifically examines how hate speech classifiers are more likely to label this type of data as toxic [43]. This difference is notable given that CLIP is not trained explicitly as a toxicity classifier, yet as a filtering method still obtains the same trends.

5.1.2 Intersections in identity keywords reveal additional filtering discrepancies. We follow up this analysis and examine intersections between various dimensions of sociodemographic identity. This limits our sample size, so in Figure 3 we only include intersections with frequency in the raw dataset of at least 10. We observe that Latinx and Asian-related keywords have higher pass rates when intersected with woman-identifying terms compared to man-identifying terms, while the European keyword has higher pass

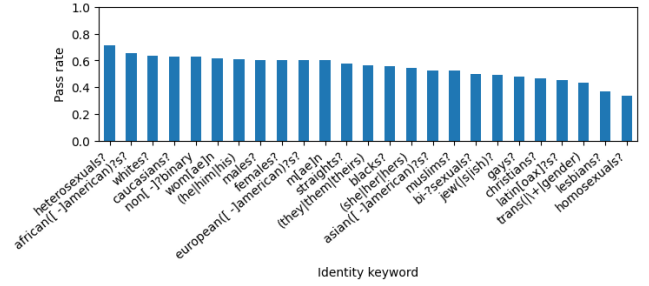
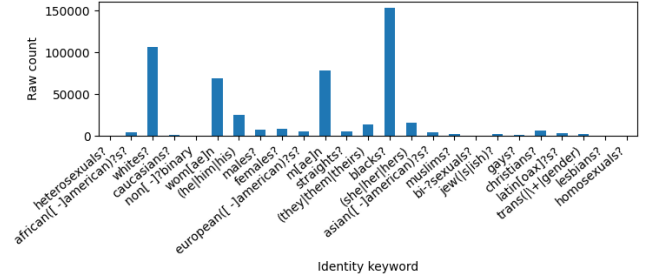
**(a) Pass rate by identity keyword sorted in descending order.****(b) Frequency in raw dataset by identity keyword.**

Figure 2: Pass rate (a) and raw dataset frequency (b) broken down by mentions of identity keywords from Table 1. A higher pass rate represents a higher proportion included in the resulting CLIP-filtered dataset. While white and black terms are commonly mentioned, we inspect samples manually and find they relate overwhelmingly to clothing items. Figure (a) shows that CLIP is more likely to exclude text samples containing LGBTQ+ keywords compared to other identity keywords.

rates when intersected with man-identifying terms compared to woman-identifying terms. For race and religion intersections, we find relatively low pass rates for the Asian and Christian intersection, and relatively high pass rates for the European and Christian intersection, compared to the average pass rates for a single keyword in the row or column of Figure 3.

5.1.3 Common words associated with certain genders embed gender stereotypes. To examine associations between words and mentions of gender, we isolate text samples that mention man or woman-related keywords (listed in Appendix A.1). Of those samples, we find a list of common words that appear in at least 100 samples, ignoring a set of stop words from the word_cloud library [84]. For each common word w , we calculate the pass rate for samples that contain both w and a woman-related keyword, and for samples that contain both w and a man-related keyword. In Figure 4, we plot the common words with the largest pass rate difference between woman and man-related samples — the top figure shows words with substantially higher woman-related pass rates, and the bottom figure shows words with substantially higher man-related pass rates. Common words that have a relatively higher woman-related pass rate are stereotypically associated with women: for instance, queen, asian, girl, and valentines. Common words that have a

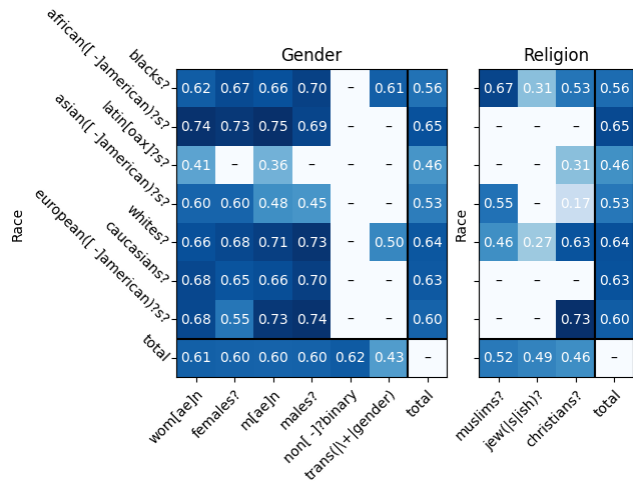


Figure 3: Heat map of pass rates by intersections of various demographic dimensions where a higher pass rate is a darker shade of blue. “-” indicates that the raw frequency of samples within that intersection is below 10 and therefore not reported. The “total” row and column in each map refer to the pass rate for all samples that mention one identity keyword. For instance, the total, wom[ae]n square shows that 61% of samples that mention keyword wom[ae]n pass the CLIP filter.

relatively higher man-related pass rate include technology, texas, person, career, comfortable, tall, and mature.

The aforementioned gender stereotype associations provide evidence that CLIP-filtering ranks data according to a stereotypical form of *alignment* between the text and images, rather than some inherent measure of their *quality*. In this example, a common word paired with one gender is included by the filter more often than the same word paired with a different gender. We see the feminization of the Asian identity [8] and the diminution of women as “girls” [79]. We also observe that text containing the word “person” is more likely to pass the filter when paired with words related to men than when paired with words related to women. Other common words with higher man-related pass rates illustrate the ongoing stereotypes of associating men with technology and careers, which have been shown to be present among both humans and machine learning models [29, 92, 137].

5.2 Image group

In this section, we describe results from analysis of the image component of the CommonPool dataset. Overall, we show that images of various imputed demographic groups according to gender, age, and race are filtered at different rates. Because image demographic imputation methods can be inaccurate for particular demographic groups [25], we apply two imputation techniques: (1) an external face attribute predictor Amazon Rekognition [5] and (2) a novel kNN-based clustering method of CLIP embeddings.

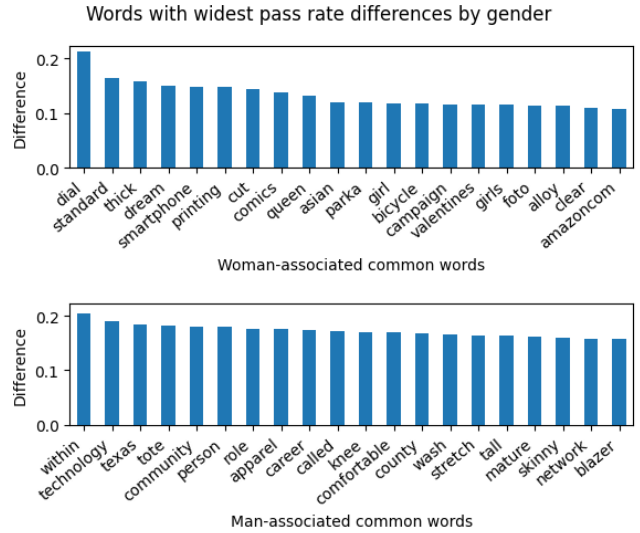


Figure 4: Common words with widest pass rate differences by gender. The top graph plots the top 20 words with the largest pass rate difference between mentions of women keywords versus mentions of men keywords and hence are more “woman-associated.” The bottom graph plots the top 20 words with the largest pass rate difference between men and women and hence are more “man-associated.”

Gender and age annotation method: We pass a subsample of 100,000 CommonPool images through Amazon Rekognition Amazon [5] in order to obtain face detections as well as imputed gender and age annotations on those detected faces. Out of this subsample, Rekognition detects faces contained in about 18,000 images, and to disambiguate between different people, who might have distinct imputed demographic attributes, we perform analysis on the 11,000 of those images that contain only one detected face. In this section, for consistency we refer to Rekognition gender labels of Male and Female but note that they conflate gender with sex [125]. We recognize the numerous issues and biases surrounding face detection and face attribute prediction as mentioned in Section 4.3.

5.2.1 Rekognition classifies more images as Male than Female in the raw dataset, and this gap widens after the CLIP filter. Out of the 11,000 images that are detected by Rekognition as containing a single face, there are more images with a Male-imputed face than images with a Female-imputed face, and the pass rate for the Male-imputed group is about 3.8 percentage points higher than the pass rate for the Female-imputed group. This difference in pass rate thus widens the representation gap after filtering – for Rekognition-annotated images, there initially are 42.1% more images in the Male-imputed group than images in the Female-imputed group in the raw dataset, but after filtering this disparity jumps to 63.9%.

This result demonstrates a notion of bias we call *exclusion amplification*. Data from imputed groups already underrepresented in the raw dataset are excluded at even higher rates (i.e. have a lower pass rate) compared to data from overrepresented groups – in other words, representation bias is amplified after filtering.

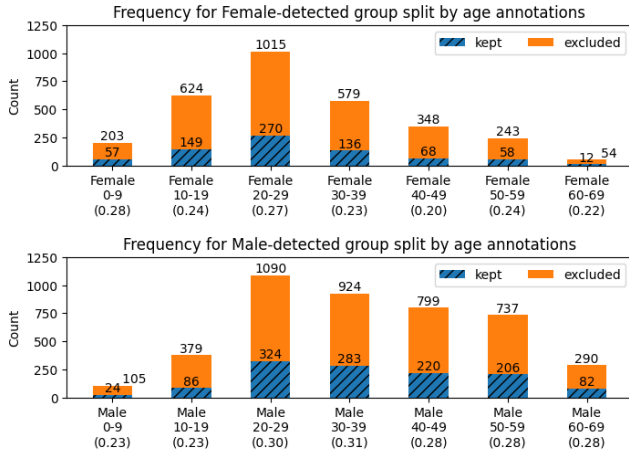


Figure 5: Frequency for samples detected by Rekognition as containing one face, split by Rekognition-detected age and gender groups. For each imputed group the pass rate is included in parentheses. We see that older detected ages in the Female-imputed group have substantially lower pass rates and lower representation than the corresponding Male-imputed group, and that younger detected ages in the Male-imputed group have lower pass rates and lower representation than the corresponding Female-imputed group.

5.2.2 *When split by imputed gender and age annotations, highly represented groups also pass through the filter at higher rates compared to less represented groups.* Figure 5 splits the Rekognition annotations by imputed gender-age group and displays the pass rate and frequency associated with each group. The age distribution concentrates around ages 20–29 for the raw and filtered datasets. Controlling for age, we observe higher representation of images in the Female-imputed groups for young ages compared to that of the corresponding Male-imputed groups. These Female-imputed groups are also more likely to pass the CLIP filter, which again reinforces the correlation between representation and pass rate. This trend then reverses for ages above 20 — for example, there are more images from the Male,60–69-imputed group as well as a higher pass rate compared to images from the Female,60–69-imputed group. This finding again supports the notion of exclusion amplification in multimodal filtering: an imputed gender-age group with higher frequency in the raw dataset is roughly associated with a higher pass rate.

Gender and race annotation method: Given the established biases and shortcomings of pre-trained third-party attribute predictors [25], in addition to using Amazon Rekognition, we introduce another technique to complement the previously presented results. We follow methodology very similar to Bianchi et al. [11] and Lucioni et al. [74], which uses embedding clustering techniques on a reference database (in our case, the Chicago Face Database using their self-reported gender and race annotations [78]) in order to assess a model’s internal representation of gender and race. We defer to Appendix A.2 the additional details about the implementation

and validation of this method in comparison to Amazon Rekognition. This method requires face box annotations, so we apply this method to the same set of 11,000 images detected by Rekognition as containing a single face. Going forward, we refer to this technique as CLIP representation-based kNN, or *CLIP kNN* for short, and we follow the same annotation group names as used in the Chicago Face Database.

At a high level, we note that CLIP kNN is not equivalent to a demographic group predictor, but rather a way to map an image to the space where CLIP encodes demographic information according to the embeddings of the Chicago Face Database. Therefore, we interpret annotations from this CLIP kNN technique as how the CLIP model internally associates images with self-reported demographic information. Because we aim to assess how CLIP acts as a filter, this allows us to group embeddings from CommonPool and evaluate the impact of filtering by each associated group.

5.2.3 *The CLIP kNN technique confirms the Rekognition findings on imputed gender, where images from the Female-imputed group are less represented in the raw dataset and excluded at slightly higher rates.* We confirm similar trends of exclusion amplification by gender on the same subsample, as shown in Figure 6. For the Female-imputed group, 608 out of 2,444 samples pass the CLIP filter for a pass rate of 0.25. For the Male-imputed group, 847 out of 3070 samples pass the CLIP filter for a pass rate of 0.28.

5.2.4 *There are differences in how the CLIP filter treats the CLIP kNN clusters by race, and White-imputed images consist of a majority of these images in both the raw and filtered datasets.* By the CLIP kNN technique, out of the same subsample of 100 thousand images that are detected by Rekognition as containing exactly one face, Figure 6 demonstrates that there are pass rate discrepancies between various race-related clusters. Asian-imputed images are filtered out the most, then White, then Latino, then Black. These relative comparisons still hold when split by gender annotations, although there is a substantial pass rate gap between Black Male (0.35) and Black Female CLIP clusters (0.29). We also find highest frequency of White-imputed images (66.8%) in the raw dataset, which also holds as most frequent in the filtered dataset.

6 LANGUAGE AND GEOGRAPHY ANALYSIS

Examination of filtering by language, country domain, or geolocation reveals that the filtered dataset overrepresents data associated with Western regions when compared to the raw dataset. As a result, we extend prior findings on Western bias in text filtering [77] to the multimodal setting. When performing the same analysis on English text, we find that these discrepancies between Western and non-Western countries still sometimes hold although the disparity is smaller, which we defer to Appendix B.2.

6.1 Language

Method: We use langdetect library [116] on a random set of 100,000 samples to determine the language of the text. While prior work demonstrates that language detection predictions may be affected by the text content [77], this method enables us to assess filtering differences by language at a large scale.

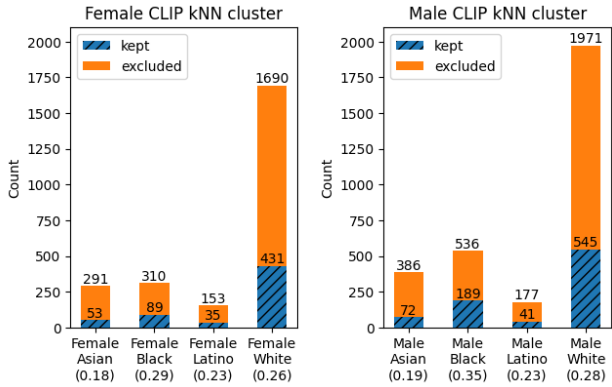


Figure 6: Frequency split by CLIP kNN according to race annotations with pass rates included in parentheses. A majority of the images in both the filtered and raw datasets belong to the White-imputed group.

6.1.1 *English is the most frequent language and is included at significantly higher rates compared to other languages.* We find that the most frequently detected language is English (35.5% of the raw dataset) and that English also has the highest pass rate (0.48). This is to be expected, given that CLIP is trained on English image-text pairs [99] and therefore would have more valid embedding representations for English text compared to other languages.

6.1.2 *Non-Western languages are not well-represented in the raw dataset, and filtering amplifies this further.* For languages outside of English, Figure 7 shows the pass rates of each language (with at least 1,000 samples) against their frequency in CommonPool. We observe that the next highest pass rates include Western languages like Spanish, Dutch, French, Catalan, German, and Portuguese, potentially due to their similarity to the English language. Furthermore, we demonstrate a similar notion of *exclusion amplification* as the trend from Section 5.2: data from detected languages that have lower representation in the raw dataset are discarded by the filter at higher rates (more details, along with a more comprehensive set of languages, are presented in Appendix B.1).

6.1.3 *Our findings reveal that multilingual data obtained after CLIP-filtering skews representation towards Western languages.* LAION-2B-multi is a subset of LAION-5B containing text in non-English languages, which is also obtained after CLIP-filtering. Researchers argue that this particular dataset is one of the largest multilingual datasets to date and therefore can fuel new research on “low-resource” languages [111]. Our findings of Western language bias in CLIP-filtering complicate LAION-2B-multi’s contribution — filtering with a CLIP model trained on English text overrepresents text written in Western languages and does not result in a uniform representation of the distribution of multilingual text on the internet.

6.2 Country domain

Method: We examine rates of filtering across country domains as a rough proxy to assess how data relating to certain geographic

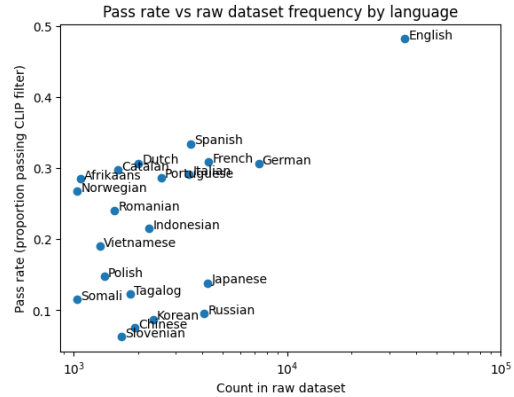


Figure 7: Pass rate versus raw dataset frequency for common languages with at least 1,000 samples in the raw dataset. Western languages have higher pass rates, and text detected as English comprise of 35,514 samples (out of 100,000 samples) in the raw dataset with a pass rate of 0.48.

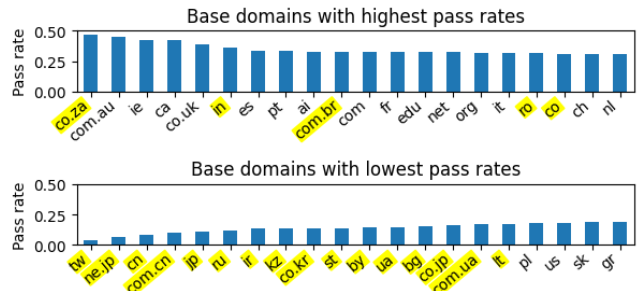


Figure 8: Base domains with at least 10,000 samples sorted by the highest (top) and lowest pass rates (bottom). Highlighted base domains associate with countries not considered in the Rich West [102]. We see that most of the base domains with lowest pass rates correspond to non-Western regions.

regions are treated by the CLIP filter. For this analysis, we use all 12.8 million samples from the small version of CommonPool and examine base domains with at least 10,000 samples from the raw dataset.

6.2.1 *Country domains from non-Western regions are excluded at higher rates than country domains from Western regions.* Figure 8 reinforces the presence of Western bias along the source modality. The base domains of websites with the highest pass rates mainly come from Western countries, while the base domains with the lowest pass rates are predominantly associated with non-Western countries.

6.3 IP address geolocation

Method: In addition to country domain, we refer to an IP address geolocation database in order to determine the country related to a source domain. Similar to the approach in Dodge et al. [44], we

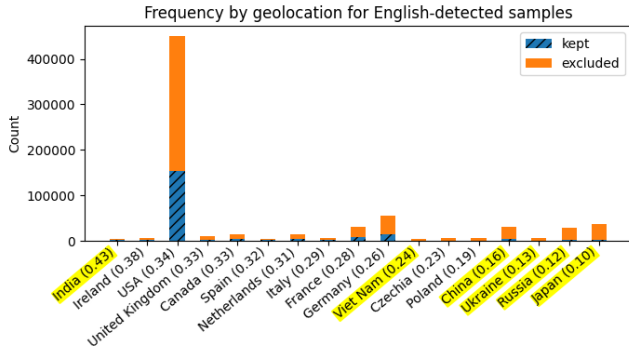


Figure 9: Raw dataset frequency by country of IP-based location sorted in descending order by pass rate (included in parentheses). We display countries with at least 5,000 samples, and most of these non-Western countries (highlighted) have low pass rates.

use the IP2Location Lite IP-Country Database [61], which has been demonstrated to have high country-level accuracy [66, 96]. We obtain the IP address set using the python socket library for a random selection of 1 million CommonPool samples, which results in 200 thousand unique websites. IP2Location then successfully identifies the country of 95.7% of the sampled website domains, which corresponds to 770 thousand CommonPool samples. We recognize that websites can be hosted in various data servers around the world, or based on user location, so IP addresses may not accurately represent the location of a website [72]. As such, we use inferred geolocation to complement the other methods we implement to obtain geographic region.

6.3.1 The inferred geolocation of most websites come from the United States, and samples with IP address locations from Western countries bypass the filter at higher rates. Figure 9 demonstrates that for both the raw dataset and filtered dataset, the inferred geolocation of an overwhelming majority of data come from the United States, although this may be due to the origin location used to ping the website. We also find that the IP address locations of countries with higher pass rates correspond to those in Western regions, and that data with inferred geolocations from non-Western regions are filtered out at higher rates.

7 SOURCE DOMAIN ANALYSIS

In this section we highlight results that examine the source domain of DataComp image-text samples and analyze what attributes correspond to being more or less likely to pass the CLIP filter.

7.1 Websites with high pass rates

Method: We determine which common websites (with at least 10,000 samples in the raw dataset) have the highest pass rates. We then manually examine samples from the websites considered “high-quality” by CLIP.

7.1.1 The majority of websites with highest-quality data according to CLIP are stock photo and E-commerce sites. Among the top 20

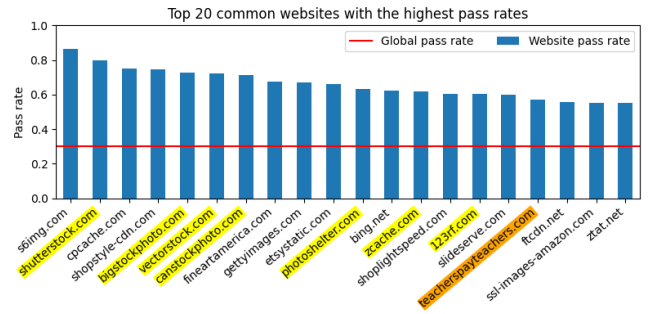


Figure 10: The top 20 websites with the highest pass rates sorted in descending order. Stock photo websites are highlighted in yellow, and Teachers Pay Teachers, an online marketplace where educators can buy and sell education curriculum, is highlighted in orange. Many of these sites where the majority of their data bypass the CLIP filter are platforms where the image assets themselves are considered copyrighted material intended to be purchased.

websites with the highest pass rates, Figure 10 shows a large number of websites that upon manual examination we categorize as e-commerce or stock photo platforms. This coincides with results from a recent LAION-Aesthetics v2.6+ audit [9], which examines frequency in the filtered dataset (rather than pass rates). We note the types of websites that are common from the prior audit but do not have high pass rates, such as user-generated content platforms like Pinterest or Wordpress.

7.1.2 Some stock photo images that are included in the final dataset are thumbnail images that do not have watermarks. Manual inspection of image-text samples from stock photo sites reveals high-resolution images with watermarks and thumbnail images without watermarks, which indicate that these images may have been scraped without the corresponding license for these websites. Shutterstock, for example, requires a license purchase before use of their images [115], and it is unclear if distributing links containing these images requires a license.

7.1.3 Intellectual property from educators intending to be purchased is also scraped as data without compensation. Additionally, Teachers Pay Teachers, another website with a high pass rate, is an online education content platform, where teachers create their own curriculum to sell to other teachers. These images in DataComp are often worksheets, which therefore disseminate the exact products that are the intellectual property of the curriculum creator [119]. These findings have potential implications in the ongoing copyright discussion of training large models [51, 57, 70], which we expound upon in Section 8.1.3.

7.2 News sites

Method: We look at pass rates of popular news sites [32] that each contain at least 200 samples.

7.2.1 The CLIP filter considers this collection of news sites (especially from the US and UK) as higher quality than average. We find on

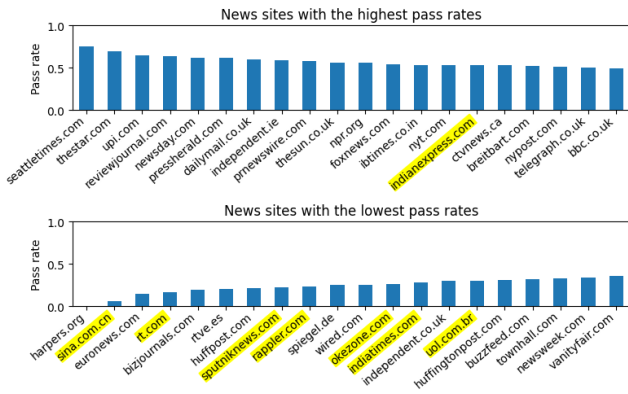


Figure 11: New sites with the highest pass rates (top) and lowest pass rates (bottom). News platforms that correspond to non-Western regions are highlighted in yellow, and we see high non-Western presence among news sites with lowest pass rates.

average higher CLIP similarity scores for samples from news sites compared to overall samples. However, Figure 11 illustrates that news sites from the United States and the United Kingdom have much higher pass rates compared to news sites from other countries, which reinforces the Western bias highlighted in Section 6. Given the ongoing lawsuit between The New York Times and OpenAI [51], according to Figure 11, we note that data from The New York Times has substantially higher pass rates than average.

7.3 Website category

Method: Inspired from prior work on text data filtering [77], we also investigate whether CLIP acts as a topical domain filter at the website level. Because prior work shows that the Cloudflare Domain Intelligence API [33] has high accuracy across categories [106], we apply this API to categorize a random subset of 100 thousand websites.

The Cloudflare API [33] returns category predictions for 94,428 out of the 100,000 website domains, where a website can correspond to multiple categories. This corresponds to 94 categories across roughly 400,000 samples in CommonPool, and we examine categories with at least 1,000 samples. Since some of these categories overlap in meaning or are confusingly defined, we merge the categories in the same manner as Ruth et al. [106] and extend it to new Cloudflare categories, which we defer to Appendix A.3.

7.3.1 E-commerce, image-hosting, and user-generated websites are among the most popular Cloudflare-detected categories in the raw dataset. Figure 12a reveals a breakdown by website category in the raw dataset. The most common categories are E-commerce, Content Servers (sites that host static images), Technology, Personal Blogs (user-generated content), and Education. This again matches prior analysis of a subset of LAION [9] with a high representation of shopping-related and user-generated websites. In our case, however, we find there are relatively fewer samples from the Stock Photos category compared to more popular categories.

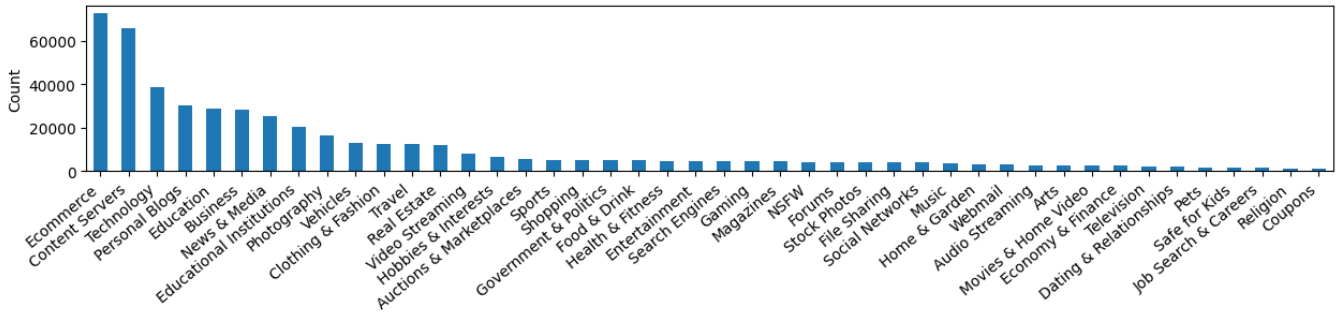
7.3.2 Data from stock photo sites are more likely to be included by the CLIP filter. The most common categories are not necessarily included at higher rates, as Figure 12b shows that the Stock Photos category has the substantially highest pass rate (20 percentage points higher than the next category, Photograph). We manually verify that the websites that fall under this category consist of mainly stock photo sites, which confirms the findings by website domain from Figure 10. In contrast to Section 7.2, we find that the News & Media category has a lower pass rate than average, possibly because this category contains a broader range of websites than the popular news dataset we use.

7.3.3 We find presence of sexually-explicit text content from NSFW-categorized websites that passes both the NSFW and CLIP filters. Figures 12a and 12b reveal that some websites are categorized by Cloudflare as NSFW, which includes content relating to pornography, nudity, extremism, and violence. This corresponds to 4,104 samples before CLIP-filtering (out of the 400,000 CommonPool samples with websites categorized) and 1,038 kept after CLIP-filtering. Upon this finding, because the image-text samples themselves may not contain NSFW, we manually examine the text of these samples to flag sexually-explicit data. Of the 4,000 or so text samples, we find 211 samples with sexually-explicit text content in CommonPool, and of these, 11 samples that pass the CLIP filter. Five sexually-explicit text samples contain words like “teen” and “schoolgirl,” and we have chosen not to visually examine these images.

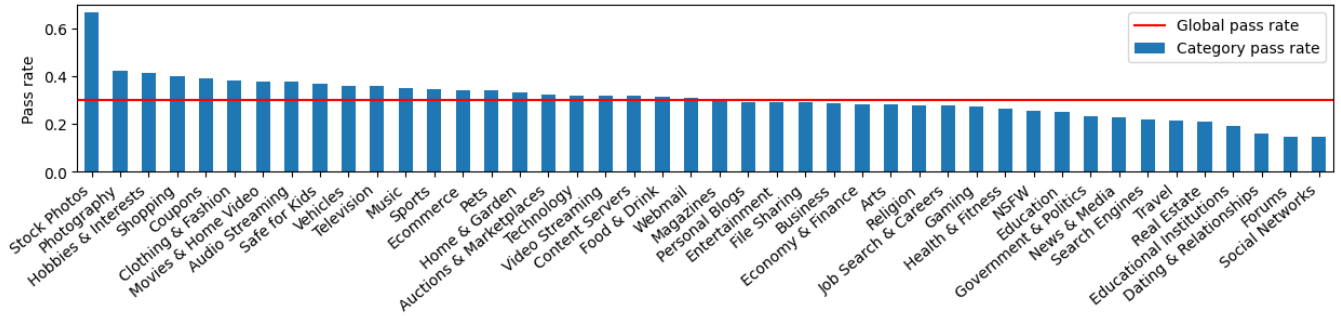
It is not surprising that sexually-explicit text passes through the CLIP filter as CLIP is not trained to detect NSFW-related content. However, in the formation of CommonPool, an NSFW text filter was applied to the initial web dump in order to remove sexually-explicit text [49]. Given we only categorize the websites of 400,000 random samples due to API limits, we find that 1.03% of samples of this small subset come from NSFW-categorized websites. Extending to all 12.8 billion samples of CommonPool, at a 95% confidence interval of this pass rate estimate, the number of samples that come from websites categorized as NSFW by Cloudflare is between 129 million and 130 million, of which 32 to 34 million samples would pass the CLIP filter. In the case of the proportion of sexually-explicit text content in our subsample, at a 95% confidence interval, the number of samples in CommonPool with sexually-explicit text is between 6.3 million and 7.2 million, of which 250,000 to 450,000 samples would pass the CLIP filter. We also note that these numbers are a substantial lower bound as we examine only English text, ignore image content, and consider samples solely from websites categorized as NSFW by Cloudflare — more investigation is needed here, but given that NSFW-filtering is outside of our original evaluation scope, we leave follow-up analysis to future work.

7.4 Origin date

Method: To track the creation date of a URL address, we apply the Wayback Machine API [7] to a subsample of 1 million image URL addresses to track the earliest-indexed date in the Internet Archive. We acknowledge that many webpages may not have snapshots on the Internet Archive and that the earliest-indexed date is only an upper bound for when the webpage was actually created.



(a) Category frequencies sorted in descending order.



(b) Category pass rates sorted in descending order. The global pass rate (0.3) is plotted in red.

Figure 12: Website categorization by Cloudflare API [33] for websites with at least 1000 samples.

7.4.1 *Most data comes from the last five years, according to earliest-indexed Internet Archive dates.* The Wayback Machine API [7] responses indicate that the Internet Archive does not record snapshots for a majority of the exact URL addresses in our subsample — we only track index dates for 22.7% of the one million CommonPool samples. Of those tracked, we find that these samples overwhelmingly come from sources that were first indexed within the last ten years. We also find that pass rates tend to be somewhat higher the more recent the year. We present the graphs and details in Appendix B.4.

8 DISCUSSION

Our analysis demonstrates that CLIP-filtering on the DataComp CommonPool dataset removes data nonuniformly when broken down by demographic group imputation, geographic region, or source domain. First, we find that the CLIP filter disproportionately excludes data relating to various proxies for demographic groups which include LGBTQ+ identities, women, and non-Western regions. Second, these filter discrepancies often amplify representation biases present in the raw dataset, where underrepresented groups are filtered out at higher rates compared to overrepresented groups. Third, we find that data from sites where the images themselves are the intellectual property (e.g., stock photos and educational curriculum) are considered high-quality by CLIP and kept at higher rates. Finally, the presence of NSFW samples after an initial cursory search of a small subset of CommonPool reveals that the combination of applying NSFW and CLIP filters fails to catch all instances of sexually-explicit content.

In this section, we elaborate on how the above findings imply that CLIP-filtering and similar filtering practices may result in societal harms and potential legal implications. Using the CLIP filter as a case study, we then dive into several key assumptions and failures of existing data cleaning methods when they are used to build large image-text models. For instance, we discuss how the CLIP filter may be more likely to pass samples similar to its own training data in Section 8.2.1. Finally, we use our analysis to inform recommendations for designing, implementing, and evaluating data filtering techniques and data curation methods.

8.1 Societal consequences of data filtering

We highlight three aspects of CLIP-filtering CommonPool data that have societal ramifications: the amplification of demographic bias, the inclusion of unmoderated inappropriate content, and the upweighting of websites that host licensed or copyrighted images. Altogether these dimensions illustrate an example of how existing data filtering practices have failed to mitigate certain downstream harms and, hence, that stochastic neural-network-based filtering may not be the end-all-be-all technique to capture *all* “impurities” in web-scraped data.

8.1.1 *Disproportionate exclusion from filtering may worsen downstream models.* As mentioned in Section 4.3, we cannot necessarily jump to conclusions that all else held equal, varying the demographic group will change the filter likelihood, and we leave this type of causal investigation [67] for our future work. It is, for example, possible that more data relating to a certain group can come

from websites with established alt-text practices, which would change filtering rates even without using a biased model like CLIP. At the same time, however, the finding of correlations between imputed demographic groups and filtering rates is important absent this causal measurement. For CommonPool, the CLIP filter modifies the original demographic distribution and thus does not result in a uniform sampling of web-scraped data. *Our analysis in Section 5.1 shows that CLIP-filtering can create training datasets more unrepresentative of certain populations and containing higher proportions of stereotypes than the unfiltered datasets.*

Because we find that the CLIP filter excludes data relating to certain imputed demographic subgroups, a lack of data from these subgroups can therefore impact downstream models trained on CLIP-filtered datasets. Prior work demonstrates that diverse representation in training is necessary to ensure model accuracy for minority subgroups [30, 104]. The exclusion amplification trend we find can further compound these model inequities: web-scraped data already does not encompass diverse views [55], yet CLIP-filtering may worsen the amount of data homogeneity in the resulting dataset. As a consequence of this form of representational harm [19], machine learning models trained on data obtained via biased filtering can underperform even more for certain marginalized groups.

In addition, we observe that the CLIP filter is more likely to include stereotypical associations between mentions of certain words and genders (Section 5.1.3). This indicates that CLIP-filtering may propagate and thus embed discriminatory associations into downstream datasets. We hypothesize this bias can lead to generative models defaulting to stereotypes of marginalized groups [103, 109].

8.1.2 Continued presence of lightly or unmoderated sexually-explicit material raises CSAM and NCII concerns. Based on the findings in Section 7.3.3, out of all 12.8 billion samples of CommonPool, we estimate approximately 130 million samples from websites categorized as NSFW by Cloudflare make it past NSFW-filtering into CommonPool, and approximately 33 million of those samples pass CLIP-filtering. While additional investigation is needed to assess CommonPool’s image content, we observe that sexually-explicit text samples are not caught by DataComp’s NSFW-filtering step. Given that the direct predecessor to CommonPool, LAION-5B, was found to contain thousands of images of child sexual abuse material (CSAM) and that some of these images were not tagged by LAION’s NSFW classifier [120], we cannot assume that both CommonPool and CLIP-filtered CommonPool do not contain CSAM, even with the prior application of an NSFW filter.

In addition to the prevalence of nonconsensual intimate imagery (NCII) on the web [68], we repeat the argument Birhane et al. [17] made against LAION-400M in 2021: neural network-based filtering cannot provide guarantees of removing all NCII and CSAM. Methods other than neural-network-based filtering are likely to be required in order to remove NCII and CSAM from web-scraped image datasets. Merely changing or updating filters, as has been done since 2021, is insufficient, as our analysis shows. These ongoing, well-documented failures are causing real, severe harms at scale: text-to-image models trained on these “filtered” datasets are being used to generate NCII of countless celebrities and random, regular people [22, 128] as well as CSAM [54, 62], with teenage girls being

targeted in particular [117]. Because of the AI field’s rapid embrace of these datasets and models despite widespread documentation of the deep NCII and CSAM harms they are causing, we urge AI industry, academia, and other players to prioritize these issues.

8.1.3 Certain types of copyrighted data are considered more “valuable” in the machine learning pipeline. There has been extensive discussion on the copyright implications of generative models trained on large-scale datasets scraped from the web. In the recent lawsuit from The New York Times against OpenAI, for instance, the plaintiff has argued that The New York Times articles are considered more “valuable” during the training process [51], where “value” here refers to the amount of content contributed to training a model. OpenAI in response has claimed that these articles are just “a tiny slice of overall training data” [94]. Our analysis in Section 7.2 provides additional evidence that data from certain websites, including The New York Times, are implicitly upweighted after filtering, which implies that systems trained through this kind of filtering do in fact place more value on certain types of copyrighted works, although it is unknown if OpenAI implements a similar filtering mechanism in constructing its training data.

In Sections 7.1 and 7.2, we show that certain kinds of typically copyrighted work pass the CLIP filter at significantly higher rates — namely stock photo sites, some news sites, and an education curriculum marketplace — all of which rely on their image and textual content as part of their business model. We also find presence of stock photo thumbnail images without watermarks, despite no known licensing agreement from the dataset curators to distribute these photos [115]. This finding indicates that data from these platforms are on average considered “high-quality” sources by CLIP, which means that the CLIP filter overrepresents and implicitly targets data from these copyrighted websites.

In order to make a case for copyright infringement in the United States, according to the fair use doctrine, the defendant must demonstrate that use of copyrighted material follows four factors [90]. According to legal analysis of fair use in the context of generative machine learning models, the fourth factor of fair use, which examines the potential market effect, draws the most attention from legal scholars [57, 70]. Alhadeff et al. [3], for example, argue that the fourth factor weighs against a finding of fair use because the output generative AI tools may be a substitution for market harm.

Text-to-image models that train on this data without purchasing licenses can therefore replicate these images which can lead to potential business losses. Users, for instance, may no longer license from stock photo sites if they can produce images on their own, or educators may generate worksheets rather than buying from curriculum developers. Our findings thus provide more context on the potential market harm as described in fair use copyright doctrine [3, 57, 70], especially if these data creators are not compensated. Dependent on additional legal analysis, this raises questions on whether plaintiffs in copyright infringement lawsuits, such as The New York Times [51], may be able to seek higher damages as filtering inadvertently deems their data as higher value.

8.2 Assumptions in data filtering practices

To understand how the above issues were overlooked in the design of CLIP-filtering, we examine how prior work justified certain

filtering choices. This enables us to form takeaways about OpenAI CLIP as a filtering model and about data filtering more generally.

8.2.1 OpenAI CLIP is not intended to assess image-text alignment.

In our work, we demonstrate instances of stereotypes perpetuated by the CLIP filter when examining common words or intersections of demographic dimensions (Sections 5.1.2 and 5.1.3), as well as instances of Western bias (Section 6). These trends support the notion that the OpenAI CLIP model is not intended to assess *image-text alignment*, as pointed out in prior work on image captioning CLIP evaluations [97] and on CLIP classification bias [1]. Especially because OpenAI CLIP was never designed to be a data filter, much less a filter used to build deployed models [93], we strongly caution against using CLIP as a filter to build future training datasets.

It is entirely plausible that OpenAI CLIP is trained on a distribution of data similar to the filtered dataset makeup, which means the CLIP filter may implicitly attempt to replicate its own training dataset. Our findings on exclusion amplification confirm this notion, as well-represented classes of data seem to be included at higher rates. However, none of this can be confirmed given that the training data for OpenAI CLIP has not been released, nor have the data filtering steps to obtain said training data been formally described [99]. In other words, the CLIP model is generated by a black-box pipeline, trained on an unknown dataset. This results in a black-box filtering model whose behavior is difficult to predict without prohibitively many interactions with it. The main advantage of CLIP-filtering here is automation, which embodies the existing notion of how “scale beats noise” as characterized in Birhane et al. [17].

8.2.2 Data filters will always encode ideology as to what is considered “high-quality.”

Our study corroborates past work on text filtering bias [44, 55, 77] in which image-text filtering also encodes ideology of what sociodemographic identities are associated with “high-quality” data. In the case of CLIP, we determine that the CLIP filter considers data relating to overrepresented groups as supposedly high-quality. Regardless of what a data filter lets through, the act of discarding data involves judgement calls of what is considered “good” or “bad” data.

Prior works on data filtering focus on improving the quality of the resulting dataset [24, 130], and there has been subsequent work examining the tradeoffs between quality and quantity [50, 88]. This notion of “quality,” however, is not well-defined and is often treated as an inherent component of the data that must somehow be discovered. As a result, filtering is depicted as a passive act, where issues of the data are blamed on the state of the world. In this manner, machine learning practitioners and researchers are able to avoid the responsibility of addressing societal inequities or problematic content [39, 56].

Investigating the design of CLIP-filtering, we find assumptions ingrained into the objectives of the filter. The LAION-400M developers optimize for a cosine-similarity threshold to improve the performance on the downstream CLIP model [112], but this seemingly neutral objective is encoded with implicit priorities. “Performance” in the DataComp benchmark refers to improvement on a specific suite of evaluation image classification and retrieval tasks, which mainly measure object classification or distribution shift [47, 49].

Birhane et al. [14] state that performance is often considered “intrinsic,” yet current measures of what constitutes as success for data filtering ignore societally-relevant concepts like fairness, toxicity generation, or privacy preservation, all of which are risks of text-to-image models [13]. CLIP-filtering is argued to be effective, but only because it aligns with a specific notion of performance, in which justification behind the selection of these metrics is unstated.

8.3 Recommendations

After questioning the assumptions of current data filtering practices, we subsequently form several recommendations. While most of these recommendations apply to data filtering methods specifically, we argue that data filters cannot be treated as a panacea [17]. As such, we extend some of our recommendations to the data curation process more generally.

8.3.1 Intentionally consider filtering criteria to account for diversity.

All these lines of inquiry demonstrate the need to build filters that are explicitly designed to include representative and non-stereotypical data from marginalized groups. Like prior works on data filtering biases [44, 55], we also recognize the need to incorporate inclusive data collection practices [10, 63], especially concentrating on ethical practices for human-centric data [6, 110]. If these large-scale web-scraped image-text datasets continue to be built, it becomes necessary to investigate how to build filters that do not perpetuate representational harms. Creating fairer filters that allow for more diverse data is an open question we leave for future work, and extending the data filtering network framework from Fang et al. [47] to downstream model bias evaluations can be a promising first step.

8.3.2 Provide justification for data filter design.

In addition to intentional data curation, we also argue that it is important to justify and unpack assumptions in filtering design choices. Given the vague definition of “performance” in machine learning literature [14], we recommend that work proposing new data filtering methods critically examine choices in the evaluation and optimization of a filtering model. Filtering is bound to make judgement calls on data, so one should understand the different impacts of various design decisions and provide justification as to what should be considered “quality” data. Because models trained on filtered data may be deployed in high-stakes settings, we also encourage considering user-centered design practices to build data filters [127].

8.3.3 Report and evaluate data filters.

We recognize the importance of documentation for filtering techniques. Academic benchmarks like DataComp’s filtering track and the open release of CommonPool [49] enable our evaluation of image-text filtering bias. Many commercially-deployed models do not release their data collection processes, much less their filtering methods, and it is unclear how many pre-trained filter models like CLIP are applied to large-scale datasets. We believe that making datasets like LAION and DataComp publicly available for independent audits is vital to the consequences of data filtering practices [15]. However, public release must include stronger safeguards against misuse. CLIP and LAION state they are not meant for real-world production contexts or applications, yet they have been used to create products with millions of users [2, 82]. DataComp CommonPool is not intended

for production-ready products, but grants permission for anyone to train models on their datasets, whether deployed or not [49]. We encourage dataset curators to use licenses that restrict against misuse, such as the RAIL license [36].

Moving forward, we argue that to even begin to build fairer filters, one must evaluate them first. We recommend future work on proposed filters to evaluate their filter models for bias building upon our approach. Our analysis reveals that the filtering step is a stage where harmful stereotypes can be injected and amplified, especially as filters can be applied to create multiple datasets. Therefore, to build multimodal models that exhibit less biased and problematic behavior, one should assess the strengths and limitations of a proposed filtering approach to account for long-term implementations on future datasets. We also recognize that it is important to conduct bias evaluations in each step in the broader machine learning pipeline [136] and to continue to do so as the pipeline evolves.

8.3.4 Account for prevention and detection of problematic content, as filters are insufficient to create legal or ethical datasets. Our work demonstrates that existing data filters do not capture all instances of problematic content. Existing NSFW or toxicity detection mechanisms are often flawed [44, 108], which indicates more research is necessary in this area. Nonetheless, additional prevention and response mechanisms are needed because filtering cannot guarantee the removal of all undesirable content [69]. Drawing from a defense-in-depth security approach [85], we recommend that researchers and dataset curators combine multiple mechanisms, which includes providing methods of recourse to take down problematic data and update downstream models accordingly [89]. In Appendix C, we highlight some potential interventions to existing datasets that contain NSFW content and discuss their implications.

Moreover, our findings corroborate the need to incorporate third-party audits into the machine learning development pipeline. For instance, our audit reveals that removing samples from domains categorized NSFW by Cloudflare is a simple practical method to reduce problematic content from one avenue, although we acknowledge that this method is not a comprehensive remedy. The presence of NCII and CSAM in even filtered web-scale datasets, along with the inability of automated tools to provide rigorous guarantees of their removal, poses a fundamental challenge to web-scale AI.

9 CONCLUSION

In this work, we audit the CLIP-filtering step in the DataComp and LAION pipelines, where to construct “high-quality” large-scale training datasets from web-scraped data, an existing pre-trained CLIP model is used to determine image-text alignment. We find that current forms of image-text filtering, similar to prior work [44, 55, 77], embed societal judgements in determining what kinds of data should be discarded. Specifically, we find that content relating to marginalized identities or non-Western regions are more likely to be left out of the final training dataset. Moreover, we demonstrate examples of *exclusion amplification* – data from certain imputed demographic groups already underrepresented in the unfiltered dataset are filtered out at higher rates. We find that copyrighted data from certain types of websites are judged to be higher quality and that NSFW filters fail to remove large quantities of sexually-explicit text, raising CSAM and NCII concerns. These findings raise

new issues on the downstream societal effects of filtering methods within the broader machine learning pipeline, especially as large multimodal models become increasingly deployed in high-stakes settings [4, 21].

ACKNOWLEDGMENTS

We would like to thank Max Del Real, Kentrell Owens, Miranda Wei, and Christina Yeung for their helpful feedback, as well as Lucy Li for the occupation dataset. We are grateful to the Cloudflare Domain Intelligence team, as well as Kimberly Ruth and Sudheesh Singanamalla, for their help with the Cloudflare API. We also thank Alex Fang, Thao Nguyen, Mitchell Wortsman, Pang Wei Koh, and Ludwig Schmidt for our initial discussions on data filtering.

The first author is supported by the NSF Graduate Research Fellowship Program. This work was supported in part by U.S. National Science Foundation awards CNS-2205171 and CCF-2045402, the Carnegie Bosch Postdoctoral Fellowship, and a grant from the Simons Foundation.

REFERENCES

- [1] Sandhini Agarwal, Gretchen Krueger, Jack Clark, Alec Radford, Jong Wook Kim, and Miles Brundage. 2021. Evaluating CLIP: Towards Characterization of Broader Capabilities and Downstream Implications. <http://arxiv.org/abs/2108.02818> arXiv:2108.02818 [cs].
- [2] Stability AI. 2024. <https://stability.ai/stable-image>
- [3] Jacob Alhadeff, Cooper Cuene, and Max Del Real. 2024. Limits of Algorithmic Fair Use. *Wash. J. L. Tech. & Arts* 19 (2024), 1.
- [4] Junaid Ali, Matthäus Kleindessner, Florian Wenzel, Kailash Budhathoki, Volkan Cevher, and Chris Russell. 2023. Evaluating the Fairness of Discriminative Foundation Models in Computer Vision. In *Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society*. ACM, Montreal, QC, Canada, 809–833. <https://doi.org/10.1145/3600211.3604720>
- [5] Amazon. 2024. Amazon Rekognition. <https://docs.aws.amazon.com/rekognition/latest/dg/what-is.html>
- [6] Jerone Andrews, Dora Zhao, William Thong, Apostolos Modas, Orestis Papyriakopoulos, and Alice Xiang. 2024. Ethical Considerations for Responsible Data Curation. *Advances in Neural Information Processing Systems* 36 (2024), 55320–55360.
- [7] Internet Archive. 2022. Wayback CDX Server API documentation. <https://archive.org/developers/wayback-cdx-server.html>
- [8] Sameena Azhar, Antonia RG Alvarez, Anne SJ Farina, and Susan Klumpner. 2021. “You’re so exotic looking”: An intersectional analysis of Asian American and Pacific Islander stereotypes. *Affilia* 36, 3 (2021), 282–301.
- [9] Andy Baio. 2022. Exploring 12 million of the 2.3 billion images used to train stable diffusion’s image generator. Retrieved July 6 (2022), 2023. <https://waxy.org/2022/08/exploring-12-million-of-the-images-used-to-train-stable-diffusions-image-generator/>
- [10] Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the dangers of stochastic parrots: Can language models be too big?. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*. ACM, Online, 610–623.
- [11] Federico Bianchi, Pratyusha Kalluri, Esin Durmus, Faisal Ladhak, Myra Cheng, Debora Nozza, Tatsunori Hashimoto, Dan Jurafsky, James Zou, and Aylin Caliskan. 2023. Easily Accessible Text-to-Image Generation Amplifies Demographic Stereotypes at Large Scale. In *2023 ACM Conference on Fairness, Accountability, and Transparency*. ACM, Chicago, IL, USA, 1493–1504. <https://doi.org/10.1145/3593013.3594095>
- [12] Jeffrey P Bigham, Ryan S Kaminsky, Richard E Ladner, Oscar M Danielsson, and Gordon L Hempton. 2006. WebInSight: making web images accessible. In *Proceedings of the 8th International ACM SIGACCESS Conference on Computers and Accessibility*. ACM, Portland, OR, USA, 181–188.
- [13] Charlotte Bird, Eddie Ungless, and Atoosa Kasirzadeh. 2023. Typology of risks of generative text-to-image models. In *Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society*. ACM, Montreal, Canada, 396–410.
- [14] Abeba Birhane, Pratyusha Kalluri, Dallas Card, William Agnew, Ravit Dotan, and Michelle Bao. 2022. The values encoded in machine learning research. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*. ACM, Seoul, South Korea, 173–184.

- [15] Abeba Birhane, Vinay Prabhu, Sang Han, Vishnu Naresh Boddeti, and Alexandra Sasha Luccioni. 2023. Into the LAION's Den: Investigating Hate in Multimodal Datasets. *Advances in Neural Information Processing Systems* 36 (2023), 21268–21284.
- [16] Abeba Birhane and Vinay Uday Prabhu. 2021. Large image datasets: A pyrrhic win for computer vision?. In *2021 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE, Online, 1536–1546.
- [17] Abeba Birhane, Vinay Uday Prabhu, and Emmanuel Kahembwe. 2021. Multimodal datasets: misogyny, pornography, and malignant stereotypes.
- [18] Sumon Biswas and Hridesh Rajan. 2021. Fair preprocessing: towards understanding compositional fairness of data transformers in machine learning pipeline. In *Proceedings of the 29th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. ACM, Athens, Greece, 981–993. <https://doi.org/10.1145/3468264.3468536>
- [19] Su Lin Blodgett, Solon Barocas, Hal Daumé III, and Hanna Wallach. 2020. Language (technology) is power: A critical survey of “bias” in NLP.
- [20] Su Lin Blodgett, Lisa Green, and Brendan O'Connor. 2016. Demographic dialectal variation in social media: A case study of African-American English.
- [21] Brandon M Booth, Louis Hickman, Shree Krishna Subburaj, Louis Tay, Sang Eun Woo, and Sidney K D'Mello. 2021. Bias and fairness in multimodal machine learning: A case study of automated video interviews. In *Proceedings of the 2021 International Conference on Multimodal Interaction*. ACM, Montreal, QC, Canada, 268–277.
- [22] Matthieu Bourel. 2024. Fake Photos, Real Harm: AOC and the Fight Against AI Porn. <https://www.rollingstone.com/culture/culture-features/aoc-deepfake-ai-porn-personal-experience-defiance-act-1234998491/>
- [23] Dawn Beverley Branley and Judith Covey. 2017. Is exposure to online content depicting risky behavior related to viewers' own risky behavior offline? *Computers in Human Behavior* 75 (2017), 283–287.
- [24] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems* 33 (2020), 1877–1901.
- [25] Joy Buolamwini and Timnit Gebru. 2018. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on Fairness, Accountability, and Transparency*. PMLR, New York, NY, USA, 77–91.
- [26] Judith Butler. 2013. Gender as performance. In *A critical sense*. Routledge, New York, NY, USA, 109–125.
- [27] C3P. 2024. Canadian Centre for Child Protection. <https://www.protectchildren.ca/en/>
- [28] Ben Caldwell, Michael Cooper, Loretta Guarino Reid, Gregg Vanderheiden, Wendy Chisholm, John Slatin, and Jason White. 2008. Web content accessibility guidelines (WCAG) 2.0. *WWW Consortium (W3C)* 290 (2008), 1–34.
- [29] Aylin Caliskan, Pimparkar Parth Ajay, Tessa Charlesworth, Robert Wolfe, and Mahzarin R Banaji. 2022. Gender bias in word embeddings: A comprehensive analysis of frequency, syntax, and semantics. In *Proceedings of the 2022 AAAI/ACM Conference on AI, Ethics, and Society*. ACM, Oxford, England, 156–170.
- [30] Irene Chen, Fredrik D Johansson, and David Sontag. 2018. Why is my classifier discriminatory? *Advances in neural information processing systems* 31 (2018).
- [31] Mehdi Chertit, Romain Beaumont, Ross Wightman, Mitchell Wortsman, Gabriel Ilharco, Cade Gordon, Christoph Schuhmann, Ludwig Schmidt, and Jenia Jitsev. 2023. Reproducible scaling laws for contrastive language-image learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. IEEE, Vancouver, Canada, 2818–2829.
- [32] B. Clemm von Hohenberg, E. Menchen-Trevino, A. Casas, and M. Wojcieszak. 2021. A list of over 5000 US news domains and their social media accounts. <https://doi.org/10.5281/zenodo.7651047>
- [33] Cloudflare. 2024. Cloudflare API v4 documentation: Get multiple domain details. <https://developers.cloudflare.com/api/operations/domain-intelligence-get-multiple-domain-details>
- [34] Samantha Cole. 2023. Largest Dataset Powering AI Images Removed After Discovery of Child Sexual Abuse Material. <https://www.404media.co/laion-datasets-removed-stanford-csam-child-abuse/>
- [35] Creative Commons. 2024. CC BY 4.0. <https://creativecommons.org/licenses/by/4.0/deed.en>
- [36] Danish Contractor, Daniel McDuff, Julia Katherine Haines, Jenny Lee, Christopher Hines, Brent Hecht, Nicholas Vincent, and Hanlin Li. 2022. Behavioral Use Licensing for Responsible AI. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency* (, Seoul, Republic of Korea.) (FAccT '22). Association for Computing Machinery, New York, NY, USA, 778–788. <https://doi.org/10.1145/3531146.3533143>
- [37] Common Crawl. 2024. <https://commoncrawl.org/>
- [38] Katherine Crowson, Stella Biderman, Daniel Kornis, Dashiell Stander, Eric Hallahan, Louis Castricato, and Edward Raff. 2022. Vqgan-clip: Open domain image generation and editing with natural language guidance. In *European Conference on Computer Vision*. Springer, Tel Aviv, Israel, 88–105.
- [39] Allan Dafoe. 2015. On technological determinism: A typology, scope conditions, and a mechanism. *Science, Technology, & Human Values* 40, 6 (2015), 1047–1076.
- [40] DataComp. 2024. DataComp Tracks. <https://www.datacomp.ai/#tracks>
- [41] Meera Desai, Abigail Jacobs, and Dallas Card. 2023. An Archival Perspective on Pretraining Data.
- [42] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding.
- [43] Lucas Dixon, John Li, Jeffrey Sorensen, Nithum Thain, and Lucy Vasserman. 2018. Measuring and mitigating unintended bias in text classification. In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*. ACM, Santa Clara, CA, 67–73.
- [44] Jesse Dodge, Maarten Sap, Ana Marasović, William Agnew, Gabriel Ilharco, Dirk Groeneveld, Margaret Mitchell, and Matt Gardner. 2021. Documenting large webtext corpora: A case study on the colossal clean crawled corpus.
- [45] 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, Cambridge, MA, USA, 214–226.
- [46] Marc N Elliott, Peter A Morrison, Allen Fremont, Daniel F McCaffrey, Philip Pantoja, and Nicole Lurie. 2009. Using the Census Bureau's surname list to improve estimates of race/ethnicity and associated disparities. *Health Services and Outcomes Research Methodology* 9 (2009), 69–83.
- [47] Alex Fang, Albin Madappally Jose, Amit Jain, Ludwig Schmidt, Alexander Toshev, and Vaishaal Shankar. 2023. Data filtering networks.
- [48] Hany Farid. 2021. An overview of perceptual hashing. *Journal of Online Trust and Safety* 1, 1 (2021), 22 pages.
- [49] Samir Yitzhak Gadre, Gabriel Ilharco, Alex Fang, Jonathan Hayase, Georgios Smyrnis, Thao Nguyen, Ryan Marten, Mitchell Wortsman, Dhruva Ghosh, Jieyu Zhang, et al. 2024. DataComp: In search of the next generation of multimodal datasets. *Advances in Neural Information Processing Systems* 36 (2024), 27092–27112.
- [50] Sachin Goyal, Pratyush Maini, Zachary C. Lipton, Aditi Raghunathan, and J. Zico Kolter. 2024. Scaling Laws for Data Filtering – Data Curation cannot be Compute Agnostic. [arXiv:2404.07177 \[cs.LG\]](https://arxiv.org/abs/2404.07177)
- [51] Michael M Grynbaum and Ryan Mac. 2023. The Times Sues OpenAI and Microsoft. , 1 pages.
- [52] Shubha Guha, Falaah Arif Khan, Julia Stoyanovich, and Sebastian Schelter. 2023. Automated Data Cleaning Can Hurt Fairness in Machine Learning-based Decision Making. In *2023 IEEE 39th International Conference on Data Engineering (ICDE)*. IEEE, Anaheim, CA, USA, 3747–3754. <https://doi.org/10.1109/ICDE55515.2023.00303>
- [53] Darren Guinness, Edward Cutrell, and Meredith Ringel Morris. 2018. Caption crawler: Enabling reusable alternative text descriptions using reverse image search. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*. ACM, Montreal, QC, Canada, 1–11.
- [54] Ritvik Gupta. 2024. LAION and the Challenges of Preventing AI-Generated CSAM. <https://www.techpolicy.press/laion-and-the-challenges-of-preventing-ai-generated-csam/>
- [55] Suchin Gururangan, Dallas Card, Sarah K Dreier, Emily K Gade, Leroy Z Wang, Zeyu Wang, Luke Zettlemoyer, and Noah A Smith. 2022. Whose language counts as high quality? measuring language ideologies in text data selection.
- [56] Alex Hanna and Tina M Park. 2020. Against scale: Provocations and resistances to scale thinking.
- [57] Peter Henderson, Xuechen Li, Dan Jurafsky, Tatsunori Hashimoto, Mark A. Lemley, and Percy Liang. 2023. Foundation Models and Fair Use. <https://doi.org/10.48550/arXiv.2303.15715> arXiv:2303.15715 [cs].
- [58] Mariya Hendriksen, Maurits Bleeker, Svitlana Vakulenko, Nanne van Noord, Ernst Kuiper, and Maarten de Rijke. 2022. Extending CLIP for Category-to-image Retrieval in E-commerce. In *European Conference on Information Retrieval*. Springer, Stavanger, Norway, 289–303.
- [59] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. 2022. Training compute-optimal large language models.
- [60] Ben Hutchinson, Andrew Smart, Alex Hanna, Emily Denton, Christina Greer, Oddur Kjartansson, Parker Barnes, and Margaret Mitchell. 2021. Towards accountability for machine learning datasets: Practices from software engineering and infrastructure. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*. ACM, Online, 560–575.
- [61] IP2Location. 2024. IP2Location Lite IP-Country IPv6 Database. <https://lite.ip2location.com/ip2location-lite>
- [62] IWF. 2023. How AI is being abused to create child sexual abuse imagery. https://www.iwf.org.uk/media/q4zll2ya/iwf-ai-csam-report_public-oct23v1.pdf
- [63] Eun Seo Jo and Timnit Gebru. 2020. Lessons from archives: Strategies for collecting sociocultural data in machine learning. In *Proceedings of the 2020 conference on fairness, accountability, and transparency*. ACM, Barcelona, Spain, 306–316.
- [64] Khari Johnson. 2020. MIT takes down 80 Million Tiny Images data set due to racist and offensive content. <https://venturebeat.com/ai/mit-takes-down-80->

- [million-tiny-images-data-set-due-to-racist-and-offensive-content/](#)
- [65] Mehtab Khan and Alex Hanna. 2022. The subjects and stages of ai dataset development: A framework for dataset accountability. *Ohio St. Tech. LJ* 19 (2022), 171.
- [66] Dan Komosny, Miroslav Voznak, and Saeed Ur Rehman. 2017. Location accuracy of commercial IP address geolocation databases. *Information technology and control* 46, 3 (2017), 333–344.
- [67] Matt J Kusner, Joshua Loftus, Chris Russell, and Ricardo Silva. 2017. Counterfactual fairness. *Advances in neural information processing systems* 30 (2017), 11 pages.
- [68] Ganaele Langlois and Andrea Slane. 2017. Economies of reputation: The case of revenge porn. *Communication and Critical/Cultural Studies* 14, 2 (2017), 120–138.
- [69] Hee-Eun Lee, Tatiana Ermakova, Vasilis Ververis, and Benjamin Fabian. 2020. Detecting child sexual abuse material: A comprehensive survey. *Forensic Science International: Digital Investigation* 34 (2020), 301022.
- [70] Katherine Lee, A. Feder Cooper, and James Grimmelmam. 2023. Talkin' 'Bout AI Generation: Copyright and the Generative-AI Supply Chain. <https://doi.org/10.2139/ssrn.4523551>
- [71] Jonathan D Levine and Frederick L Oswald. 2013. O* NET: The occupational information network. In *The Handbook of Work Analysis*. Routledge, New York, NY, USA, 312–332.
- [72] Ioana Livadariu, Thomas Dreiholz, Anas Saeed Al-Selwi, Haakon Bryhni, Olav Lysne, Steinar Bjornstad, and Ahmed Elmokashfi. 2020. On the accuracy of country-level IP geolocation. In *Proceedings of the applied networking research workshop*. ACM, Online, 67–73.
- [73] Ian F Haney Lopez. 1995. *The social construction of race*. Harvard Civil Rights-Civil Liberties Law Review, Cambridge, MA, USA.
- [74] Alexandra Sasha Luccioni, Christopher Akiki, Margaret Mitchell, and Yacine Jernite. 2024. Stable bias: Evaluating societal representations in diffusion models. *Advances in Neural Information Processing Systems* 36 (2024), 56338–56351.
- [75] Alexandra Sasha Luccioni, Frances Corry, Hamsini Sridharan, Mike Ananny, Jason Schultz, and Kate Crawford. 2022. A framework for deprecating datasets: Standardizing documentation, identification, and communication. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*. ACM, Seoul, South Korea, 199–212.
- [76] Alexandra Sasha Luccioni and Joseph Viviano. 2021. What's in the box? An analysis of undesirable content in the Common Crawl corpus. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*. ACL, Bangkok, Thailand, 182–189.
- [77] Li Lucy, Suchin Gururangan, Luca Soldaini, Emma Strubell, David Bamman, Lauren Klein, and Jesse Dodge. 2024. AboutMe: Using Self-Descriptions in Webpages to Document the Effects of English Pretraining Data Filters.
- [78] Debbie S Ma, Joshua Correll, and Bernd Wittenbrink. 2015. The Chicago face database: A free stimulus set of faces and norming data. *Behavior research methods* 47 (2015), 1122–1135.
- [79] Susan R Madsen. 2021. Why Calling Women'Girls' Is A Bigger Deal Than You May Think.
- [80] Emanuel Maiberg. 2024. Tech Companies Promise to Try to Do Something About All the AI CSAM They're Enabling. <https://www.404media.co/tech-companies-promise-to-try-to-do-something-about-all-the-ai-csam-theyre-enabling/>
- [81] Microsoft. 2024. PhotoDNA. <https://www.microsoft.com/en-us/photodna>
- [82] Midjourney. 2024. <https://www.midjourney.com/home>
- [83] Ron Mokady, Amir Hertz, and Amit H Bermano. 2021. Clipcap: Clip prefix for image captioning.
- [84] Andreas Mueller. 2023. word_cloud. https://github.com/amueller/word_cloud
- [85] Arif Ali Mughal. 2018. The Art of Cybersecurity: Defense in Depth Strategy for Robust Protection. *International Journal of Intelligent Automation and Computing* 1, 1 (2018), 1–20.
- [86] Michael Muller and Angelika Strohmayer. 2022. Forgetting practices in the data sciences. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*. ACM, New Orleans, LA, USA, 1–19.
- [87] NCMC. 2024. National Center for Missing and Exploited Children. <https://www.missingkids.org/home>
- [88] Thao Nguyen, Gabriel Ilharco, Mitchell Wortsman, Sewoong Oh, and Ludwig Schmidt. 2022. Quality not quantity: On the interaction between dataset design and robustness of clip. *Advances in Neural Information Processing Systems* 35 (2022), 21455–21469.
- [89] Thanh Tam Nguyen, Thanh Trung Huynh, Phi Le Nguyen, Alan Wee-Chung Liew, Hongzhi Yin, and Quoc Viet Hung Nguyen. 2022. A survey of machine unlearning.
- [90] Copyright Law of the United States. 1976. Section 107. Limitations on exclusive rights: Fair use. <https://www.copyright.gov/title17/92chap1.html#107>
- [91] Maya Okawa, Ekdeep S Lubana, Robert Dick, and Hidenori Tanaka. 2024. Compositional abilities emerge multiplicatively: Exploring diffusion models on a synthetic task. *Advances in Neural Information Processing Systems* 36 (2024), 23 pages.
- [92] Ruth Oldenziel. 1999. *Making technology masculine: men, women and modern machines in America, 1870-1945*. Amsterdam University Press, Amsterdam, Netherlands.
- [93] OpenAI. 2022. Model Card: CLIP. <https://github.com/openai/CLIP/blob/main/model-card.md>
- [94] OpenAI. 2024. OpenAI and journalism. <https://openai.com/blog/openai-and-journalism>
- [95] Kenny Peng, Arunesh Mathur, and Arvind Narayanan. 2021. Mitigating dataset harms requires stewardship: Lessons from 1000 papers.
- [96] Ingmar Poesse, Steve Uhlig, Mohamed Ali Kaafar, Benoit Donnet, and Bamba Gueye. 2011. IP geolocation databases: Unreliable? *ACM SIGCOMM Computer Communication Review* 41, 2 (2011), 53–56.
- [97] Haoyi Qiu, Zi-Yi Dou, Tianlu Wang, Asli Celikyilmaz, and Nanyun Peng. 2023. Gender Biases in Automatic Evaluation Metrics for Image Captioning. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, Houda Bouamor, Juan Pino, and Kalika Bali (Eds.). Association for Computational Linguistics, Singapore, 8358–8375. <https://doi.org/10.18653/v1/2023.emnlp-main.520>
- [98] Yiting Qu, Xinyue Shen, Xinlei He, Michael Backes, Savvas Zannettou, and Yang Zhang. 2023. Unsafe diffusion: On the generation of unsafe images and hateful memes from text-to-image models. In *Proceedings of the 2023 ACM SIGSAC Conference on Computer and Communications Security*. ACM, Copenhagen, Denmark, 3403–3417.
- [99] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*. PMLR, Online, 8748–8763.
- [100] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog* 1, 8 (2019), 9.
- [101] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research* 21, 1 (2020), 5485–5551.
- [102] World Population Review. 2024. Western Countries 2024. <https://worldpopulationreview.com/country-rankings/western-countries>
- [103] Reece Rogers. 2024. Here's How Generative AI Depicts Queer People.
- [104] Esther Rolf, Theodora T Worledge, Benjamin Recht, and Michael Jordan. 2021. Representation matters: Assessing the importance of subgroup allocations in training data. In *International Conference on Machine Learning*. PMLR, Online, 9040–9051.
- [105] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-Resolution Image Synthesis With Latent Diffusion Models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, New Orleans, LA, USA, 10684–10695.
- [106] Kimberly Ruth, Aurore Fass, Jonathan Azose, Mark Pearson, Emma Thomas, Caitlin Sadowski, and Zakir Durumeric. 2022. A world wide view of browsing the world wide web. In *Proceedings of the 22nd ACM Internet Measurement Conference*. ACM, Nice, France, 317–336.
- [107] Aneeshan Sain, Ayan Kumar Bhunia, Pinaki Nath Chowdhury, Subhadeep Koley, Tao Xiang, and Yi-Zhe Song. 2023. Clip for all things zero-shot sketch-based image retrieval, fine-grained or not. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. ACM, Vancouver, Canada, 2765–2775.
- [108] Maarten Sap, Dallas Card, Saadia Gabriel, Yejin Choi, and Noah A Smith. 2019. The risk of racial bias in hate speech detection. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. ACL, Florence, Italy, 1668–1678.
- [109] Mia Sato and Emillia David. 2024. I'm still trying to generate an AI Asian man and white woman.
- [110] Morgan Klaus Scheuerman, Alex Hanna, and Emily Denton. 2021. Do datasets have politics? Disciplinary values in computer vision dataset development. *Proceedings of the ACM on Human-Computer Interaction* 5, CSCW2 (2021), 1–37.
- [111] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. 2022. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in Neural Information Processing Systems* 35 (2022), 25278–25294.
- [112] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. 2021. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs.
- [113] Carsten Schwemmer, Carly Knight, Emily D Bello-Pardo, Stan Oklobdzija, Martijn Schoonvelde, and Jeffrey W Lockhart. 2020. Diagnosing gender bias in image recognition systems. *Socius* 6 (2020), 2378023120967171.
- [114] Renee Shelby, Shalaleh Rismani, Kathryn Henne, AJung Moon, Negar Rostamzadeh, Paul Nicholas, N'Mah Yilla-Akbari, Jess Gallegos, Andrew Smart, Emilio Garcia, et al. 2023. Sociotechnical harms of algorithmic systems: Scoping

- a taxonomy for harm reduction. In *Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society*. ACM, Montreal, QC, Canada, 723–741.
- [115] Shutterstock. 2024. Can I use images on my website? https://support.shutterstock.com/s/article/Can-I-use-Images-on-my-website?language=en_US
- [116] Nakatani Shuyo. 2014. langdetect. <https://github.com/Mimino666/langdetect>
- [117] Natasha Singer. 2024. Teen Girls Confront an Epidemic of Deepfake Nudes in Schools.
- [118] Morgan P Slusher and Craig A Anderson. 1987. When reality monitoring fails: The role of imagination in stereotype maintenance. *Journal of Personality and Social Psychology* 52, 4 (1987), 653.
- [119] Teachers Pay Teachers. 2022. How do I obtain a copyright in my work? Should I register my copyright? <https://help.teacherspayteachers.com/hc/en-us/articles/360042535652-How-do-I-obtain-a-copyright-in-my-work-Should-I-register-my-copyright>
- [120] David Thiel. 2023. Identifying and Eliminating CSAM in Generative ML Training Data and Models.
- [121] David Thiel, Melissa Stroebel, and Rebecca Portnoff. 2023. Generative ML and CSAM: Implications and Mitigations.
- [122] Thorn. 2024. Safer. <https://get.safer.io/csam-detection-tool-for-child-safety>
- [123] Thorn. 2024. Safety by Design for Generative AI: Preventing Child Sexual Abuse. <https://info.thorn.org/hubfs/thorn-safety-by-design-for-generative-AI.pdf>
- [124] Yonglong Tian, Dilip Krishnan, and Phillip Isola. 2020. Contrastive multiview coding. In *European Conference on Computer Vision*. Springer, Online, 776–794.
- [125] Francisco Valdes. 1996. Unpacking hetero-patriarchy: tracing the conflation of sex, gender & (and) sexual orientation to its origins. *Yale JL & Human.* 8 (1996), 161.
- [126] Pranshu Verma and Drew Harwell. 2023. Exploitive, illegal photos of children found in the data that trains some AI. <https://www.washingtonpost.com/technology/2023/12/20/ai-child-pornography-abuse-photos-laion/>
- [127] Karel Vredenburg, Ji-Ye Mao, Paul W Smith, and Tom Carey. 2002. A survey of user-centered design practice. In *Proceedings of the SIGCHI conference on Human factors in computing systems*. ACM, Minneapolis, MN, USA, 471–478.
- [128] Jess Weatherbed. 2024. Trolls have flooded X with graphic Taylor Swift AI fakes.
- [129] WebAIM. 2024. The WebAIM Million: An annual accessibility analysis of the top 1,000,000 home pages. <https://webaim.org/projects/million/#alttext>
- [130] Guillaume Wenzek, Marie-Anne Lachaux, Alexis Conneau, Vishrav Chaudhary, Francisco Guzmán, Armand Joulin, and Edouard Grave. 2019. CCNet: Extracting high quality monolingual datasets from web crawl data.
- [131] Kaylee Williams. 2023. Exploring Legal Approaches to Regulating Nonconsensual Deepfake Pornography. <https://www.techpolicy.press/exploring-legal-approaches-to-regulating-nonconsensual-deepfake-pornography/>
- [132] Benjamin Wilson, Judy Hoffman, and Jamie Morgenstern. 2019. Predictive inequity in object detection.
- [133] Robert Wolfe and Aylin Caliskan. 2022. American == White in Multimodal Language-and-Image AI. In *Proceedings of the 2022 AAAI/ACM Conference on AI, Ethics, and Society*. ACM, Oxford, United Kingdom, 800–812. <https://doi.org/10.1145/3514094.3534136>
- [134] Robert Wolfe, Yiwei Yang, Bill Howe, and Aylin Caliskan. 2023. Contrastive language-vision ai models pretrained on web-scraped multimodal data exhibit sexual objectification bias. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*. ACM, Chicago, IL, USA, 1174–1185.
- [135] Hu Xu, Saining Xie, Xiaoqing Ellen Tan, Po-Yao Huang, Russell Howes, Vasu Sharma, Shang-Wen Li, Gargi Ghosh, Luke Zettlemoyer, and Christoph Feichtenhofer. 2023. Demystifying clip data.
- [136] Ke Yang, Biao Huang, Julia Stoyanovich, and Sebastian Schelter. 2020. Fairness-Aware Instrumentation of Preprocessing Pipelines for Machine Learning.
- [137] Jiping Zuo and Shengming Tang. 2000. Breadwinner status and gender ideologies of men and women regarding family roles. *Sociological perspectives* 43, 1 (2000), 29–43.

A METHODOLOGY DETAILS

In this section, we highlight additional details on the demographic group imputation techniques and website categorization method.

A.1 Gender keywords

Table 2 describes the keywords relating to mentions of women and men in the text samples, in order to perform analysis in pass rates by common words in Section 5.1.3.

Table 2: List of regular expressions relating to gender keywords for women and men.

Women-related	Man-related
wom[ae]n	m[ae]n
females?	males?
(she her hers)	(he him his)

A.2 CLIP kNN

The CLIP representation-based kNN method (used in Section 5.2.3) is an audit technique to evaluate CLIP associations with demographic markers like gender and race without ground-truth demographic annotations. Based on the method from Bianchi et al. [11], we use the Chicago Face Database [78] as a reference database to cluster CLIP embeddings of CommonPool images. The steps are defined as follows:

- (1) We obtain the CLIP embeddings for images in the Chicago Face Database, or CFD, which contain self-reported gender and race attributes corresponding to the images of 597 unique individuals [78].
- (2) We then consider image samples from our CommonPool subsample that have a single face detected by Rekognition [5]. We crop these images to the bounding box.
- (3) Given a facial image, we extract the CLIP embedding and run k -nearest neighbors with respect to the CFD embeddings. Note that we use separate kNN classifiers for gender and race. We select $k = 7$ for gender and $k = 5$ for race to maximize validation accuracy on held-out CFD images, and use Minkowski distance as the distance metric.
- (4) The distribution of the nearest neighbors becomes the probability scores of the gender or race of the DataComp image, and the argmax is the final group annotation.

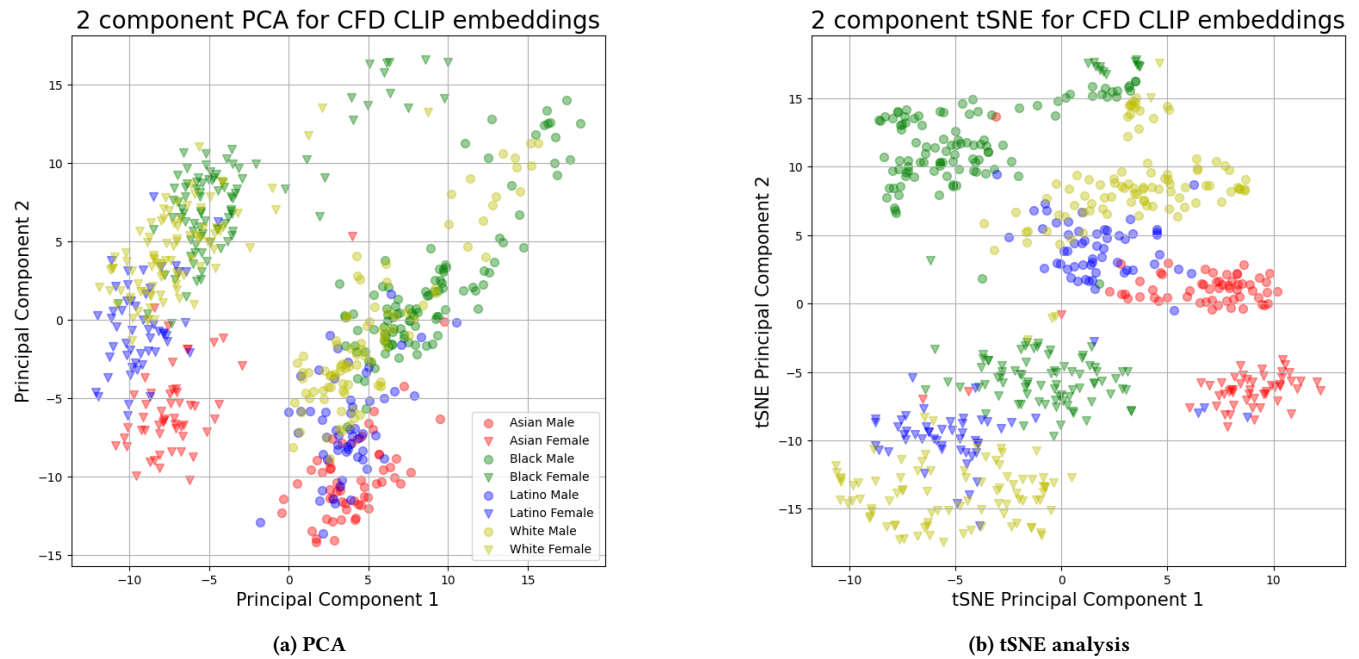


Figure 13: Visualizations of CLIP embeddings of Chicago Face Database images [78], grouped by gender-race annotations. We observe clear clustering by gender and race, although within the same group, embeddings fall into multiple clusters.

This process exactly follows the method in Bianchi et al. [11] with CFD as the reference dataset, except for steps 3 and 4, in which their work calculates the average of all embeddings grouped by gender or race in order to obtain an “archetypal vector representation” for a demographic category. Then, for a given test image embedding they select the group with the closest average embedding in cosine distance.

Initially, we find that the CFD embeddings are clustered by demographic group as shown in Figure 13. However, within the same gender and race group, there are multiple clusters at different regions in the embedding space. This indicates that averaging may not be as effective as a clustering-based classification method, which motivates our use of k-nearest neighbors instead.

A.2.1 Validation. On 160 randomly held-out CFD images, the gender annotation method obtains a validation accuracy rate of 100% (across all gender-race groups). The race annotation method obtains a validation accuracy of 90.6%, but we note that the Latino Female group obtains the lowest accuracy of 65.0%.

We find that the CLIP kNN method closely agrees with the gender predictions of Amazon Rekognition when applied to the same subset of 11,000 images that Rekognition detects as containing a single face. With a confidence threshold set to include only samples where all k -nearest neighbors (in CFD) have agreeing ground-truth gender annotations, we find that 95.4% of the 5,514 CLIP kNN annotations match the Rekognition annotations. Of the 254 samples that disagree, manual examination reveals that the CLIP kNN annotation more often aligns with human annotation. Again, we note that the CLIP kNN method is not meant to be interpreted as a perfect prediction of gender, especially as gender is not a visual construct [26], but rather to assess CLIP’s internalized encoding of the gender attribute. This technique allows us to evaluate whether the CLIP filter treats these internal associations differently. Figure 14 show that both methods support trends of *exclusion amplification* by imputed gender.

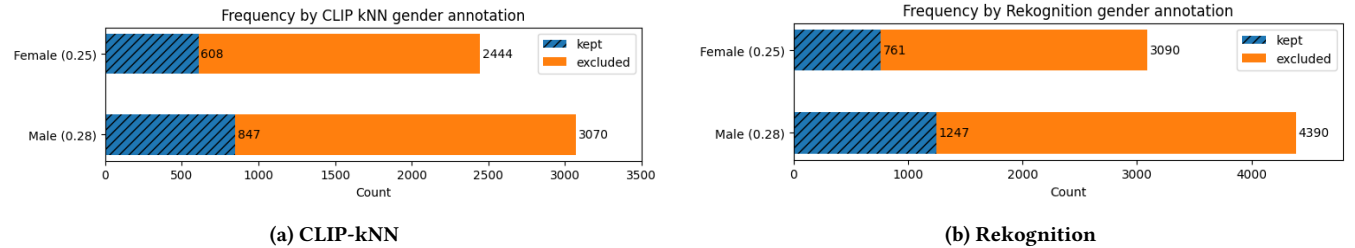


Figure 14: Frequency by imputed gender for both CLIP-kNN and Rekognition methods applied to a set of 11,000 images that Rekognition detects as containing a face. We see similar trends of exclusion amplification where there are more images from the Male-imputed group than images from the Female-imputed group. Moreover, the pass rate for the Male-imputed group is higher than that of the Female-imputed group. We note that CLIP-kNN annotates fewer total images than Rekognition due to the confidence threshold.

A.2.2 Limitations. We recognize that the kNN annotation method relies on face detection boxes as input, which requires an accurate face detection model. This limits the extension of our method to a larger set of images because in our initial analysis we observe that the DataComp face-bounding boxes provided in DataComp metadata often do not contain human faces.

A.3 Cloudflare website categorization

In Table 3, we show the merged category names or renamings based on the names returned by the Cloudflare Domain Intelligence API [33].

B ADDITIONAL RESULTS

Here we present additional results from our analysis in the manuscript.

B.1 Detected languages

Figure 15 plots the pass rate versus the raw dataset frequency for all detected languages that appear in at least 100 samples (out of the 100,000 randomly-chosen samples). Again, we find that the detected languages with the highest pass rates are Western languages, and a statistically-significant positive relationship between pass rate and raw dataset frequency. This illustrates the *exclusion amplification* trend we find, where underrepresented languages are excluded at higher rates.

B.2 Western bias on English data

When examining filter discrepancies on data relating to different geographic regions, we isolate our analysis to English-detected data (via langdetect [116]). Figure 16a, Figure 16b, and Figure 16c show the pass rates by news site, IP address geolocation, and country domain, respectively. We find that similar trends of Western bias hold, although much lower in magnitude and not across all dimensions of analysis.

Table 3: A mapping of category names to the corresponding Cloudflare categories. All other category names follow Cloudflare naming [33]. We follow categorization merging from Ruth et al. [106] and extend to new or confusingly-named Cloudflare categories.

Merged category name	Cloudflare category names
Arts	Arts; Fine Art
Audio Streaming	Audio Streaming; Radio
Business	Business; Professional Networking
Chat & Messaging	Chat; Instant Messengers; Messaging
Clothing & Fashion	Clothing; Fashion; Lingerie & Bikini; Swimsuits
Drugs & Alcohol	Alcohol; Drugs; Tobacco
File Sharing	File Sharing; Photo Sharing
Government & Politics	Government; Military; Politics, Advocacy, and Government-Related
Movies & Home Video	Home Video/DVD; Movies
NSFW	Adult Themes; CIPA Filter; Militancy, Hate & Extremism; Nudity; Pornography; Violence; Weapons
Science	Science; Space & Astronomy
Stock Photos	News, Portal & Search
Technology	APIs; Artificial Intelligence; Information Security; Information Technology; Technology
Video Streaming	P2P; Video Streaming

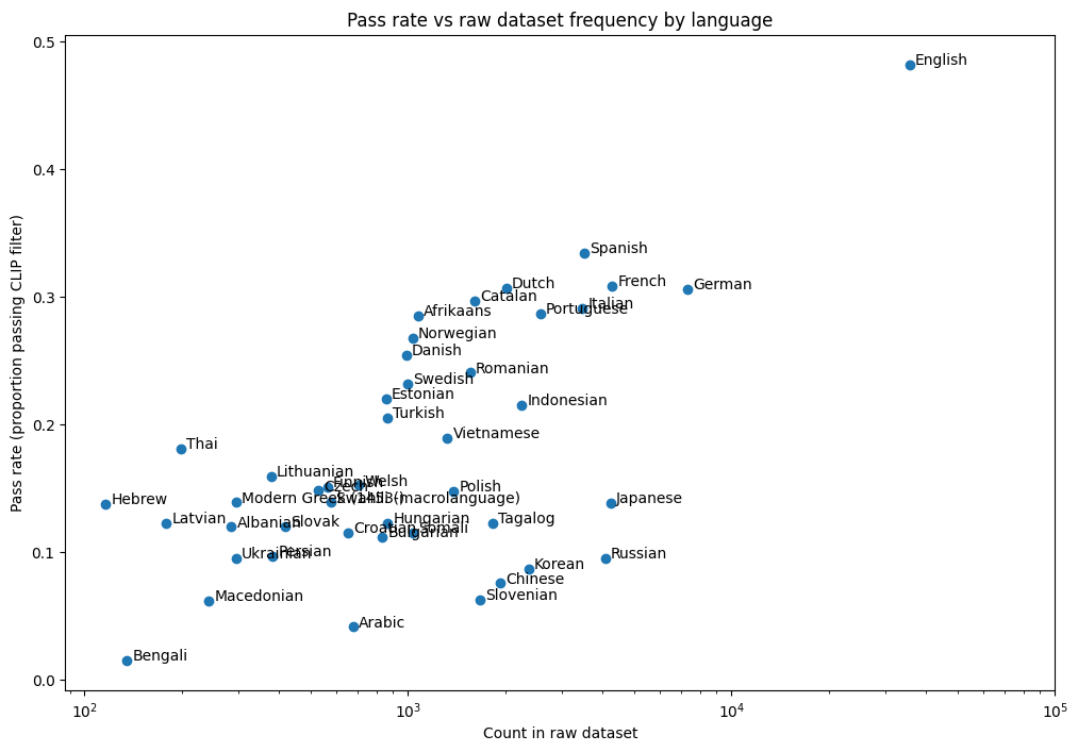
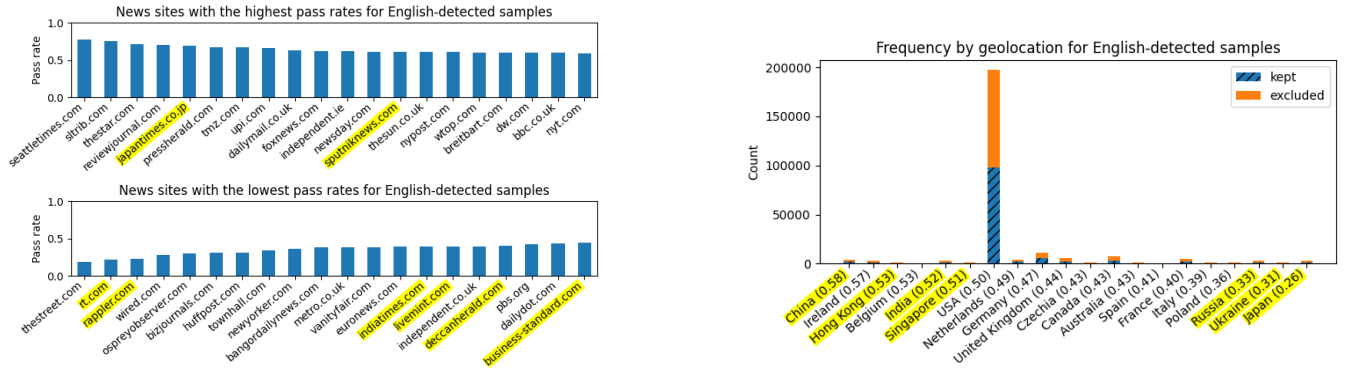


Figure 15: Pass rate vs raw dataset frequency for a more comprehensive set of languages with at least 100 samples in the raw dataset. We observe a positive trend ($p < 0.001$) in which the more represented a language is in the raw dataset then the higher the pass rate.

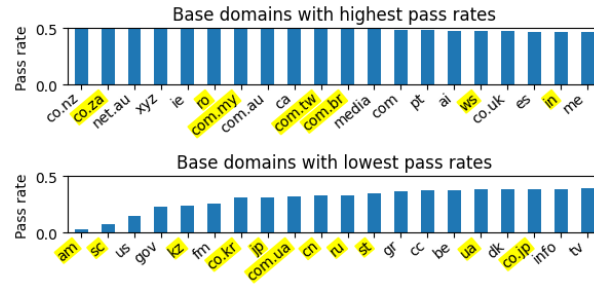
B.3 Examples of sexually-explicit text

Table 4 displays samples of sexually-explicit text that we manually find among the data from websites categorized as NSFW by Cloudflare. These samples are not caught by the NSFW filter, and some of them pass the CLIP filter as well.



(a) New sites with the highest pass rates (top) and lowest pass rates (bottom) for English-detected samples.

(b) Frequency of English-detected data by country of IP address sorted by pass rates in descending order.



(c) Pass rates for base domains with the highest (top) and lowest pass rates (bottom) for English-detected data.

Figure 16: Pass rate by news site (a), IP address geolocation (b), and country domain (c) on English-detected data. Groups that correspond to non-Western regions are highlighted in yellow. We observe that some non-Western regions now have high pass rates, although data relating to many non-Western regions are included by the CLIP filter at low rates.

Table 4: Examples of sexually-explicit text in CommonPool (i.e. passes initial NSFW filter). Names and locations are redacted.

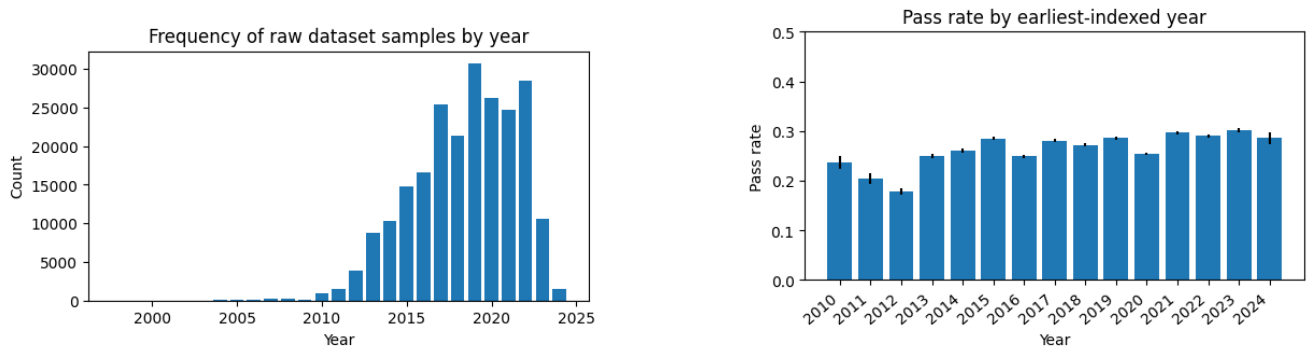
Pass CLIP filter	Text
Yes	**** Cleavage
Yes	**** caught by my spy cam taking a shower
Yes	... you_gonna_get_raped ...
No	Teen at gyno Thumbnail
No	Sexy **** teen girls
No	**** Nude Leaks

B.4 Internet Archive earliest-indexed date

In Figure 17, we demonstrate that most samples in the Internet Archive are first indexed in the last five years. Moreover, when we limit to years with at least 500 samples (i.e. 2010–2024), we observe a slight increasing trend in pass rates. More recent years have slightly higher pass rates than prior years and at the same time are more popular in the raw dataset.

B.5 Occupation

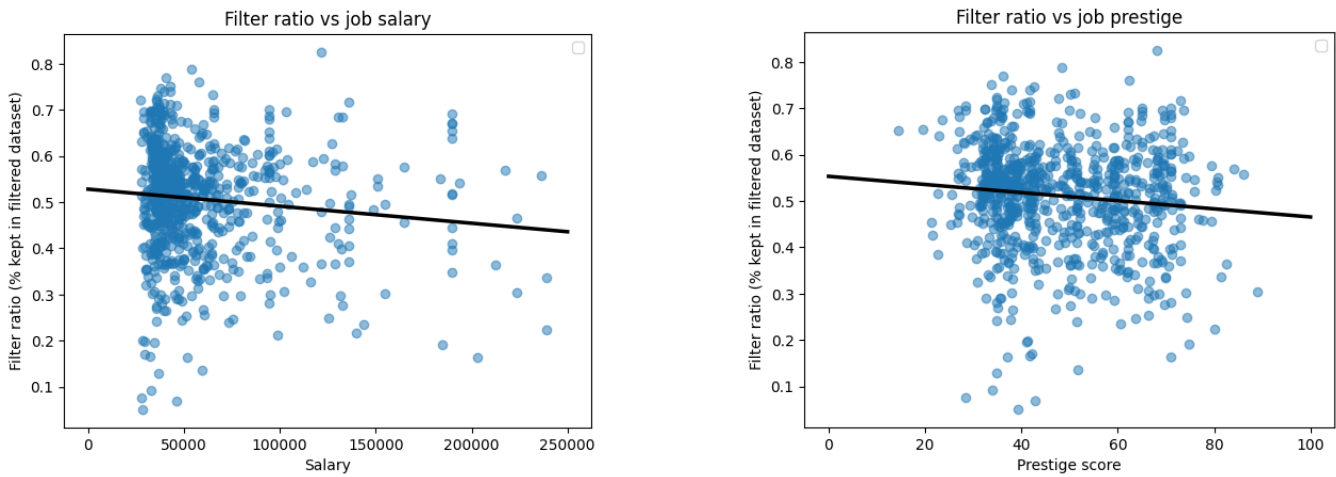
We examine text samples in CommonPool that mention occupations. With the O*NET job titles, salary, and prestige [71], we find in Figure 18a and Figure 18b slight statistically-significant negative relationships. In other words, mentions of occupations with higher salary or higher prestige are more likely to be excluded. While this counters the findings from obtained from Lucy et al. [77], these negative relationships are slight, however, so it is difficult to form a meaningful observation.



(a) Frequency by year in raw dataset. A majority of the samples examined come from the years 2019 - 2024.

(b) Pass rate by year for years with at least 500 associated samples. From 2010 - 2024 there is a slight positive trend in pass rate.

Figure 17: Raw dataset frequency (a) and pass rate (b) by earliest-indexed year in the Internet Archive. A higher pass rate represents a higher proportion included in the resulting CLIP-filtered dataset.



(a) Pass rates vs job salary for O*NET job titles. $p = 0.001$

(b) Pass rates vs job prestige for O*NET job titles. $p = 0.002$

Figure 18: Pass rates by occupations according to O*NET job titles [71] following analysis from Lucy et al. [77].

B.6 Dialect

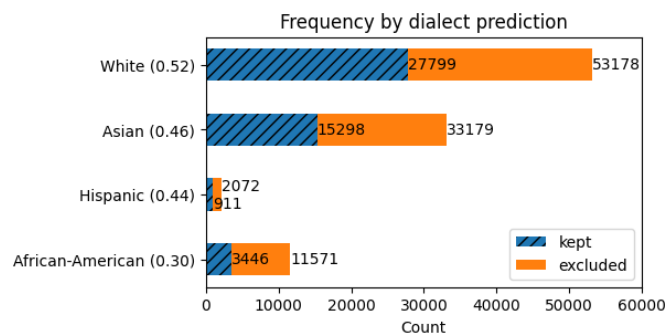


Figure 19: Frequency by imputed dialect applied to a random subset of 100,000 samples detected with English text.

Figure 19 illustrates the frequency and pass rates by the dialect prediction model from Blodgett et al. [20] on a random subset of 100,000 samples detected as English by Langdetect [116]. However, because the dialect prediction model was trained on Twitter data, upon manual examination we find that the dialect predictions are difficult to interpret with alt-text as input.

C POTENTIAL MITIGATIONS FOR DATASET CREATORS

In this section, we list several actions that maintainers of web-scraped datasets like CommonPool may consider in order to mitigate the potential harms of NSFW, NCII, and CSAM content present in these datasets, as defined in Section 8.1.2. These interventions may be components of the defense-in-depth approach discussed in Section 8.3.4 and thus may be combined together, but cannot be considered complete solutions as they each contain tradeoffs. For each potential mitigation we detail their limitations and implications.

While we focus on the actions of dataset creators here, we also note that determining who is responsible for the presence of illegal or harmful data [60, 65] remains outside the scope of this work. Given that CommonPool sources from Common Crawl [37], it is unclear how Common Crawl can address some of these NSFW concerns [76]. In the case of dataset users (in other words, model developers), any update to the dataset may not impact models already trained on the original dataset. For instance, models trained on CSAM are not considered child sexual abuse material in US law [54], which explains why Midjourney and Stable Diffusion were not taken down when their training data was revealed to contain CSAM [34]. In addition, legal regulation of nonconsensual deepfake sexual imagery is still patchy, with many jurisdictions having no specific laws about deepfake NCII [131].

C.1 Remove images from websites labeled as NSFW by Cloudflare

In Section 7.3.3, our audit finds that 1.03% of CommonPool samples come from websites that Cloudflare labeled as NSFW (including pornography). These samples could easily be detected and removed from the original dataset, without checking or downloading image content. While simple to implement, we again caution that this method would not remove *all* NSFW content and would not impact copies of the dataset already downloaded nor models already trained on portions of CommonPool. The NSFW-categorized domains corresponding to the subset of 400,000 samples are available upon request.

C.2 Run hash-based detection to known CSAM lists

Dataset creators may follow perceptual hashing methods [48, 120] to compare data samples to known hash sets of CSAM provided by existing organizations like National Center for Missing and Exploited Children (NCMEC) and Canadian Centre for Child Protection (C3P) [27, 87]. Other APIs like Microsoft’s PhotoDNA or Thorn’s Safer tools can also be used to identify CSAM [81, 122].

A lack of positive match here indicates that the dataset does not contain images that correspond to *known* CSAM hashes, but does not reveal insight about *unknown* instances of CSAM, much less instances of NCII or NSFW materials. If a sample is a positive match, there may be complications on its removal [120], as the dataset has likely been downloaded and copied many times. Removing CSAM then re-uploading the dataset could allow actors with old copies of the dataset to quickly identify CSAM [34].

C.3 Remove people from dataset

Even if all instances of CSAM could be removed from datasets, the problematic ability of models to generate CSAM remains. Due to the capabilities of generative models to compose effectively [91], prior investigation shows that bad actors may distort benign images of children into sexualized images or use de-aging techniques on NSFW-generated content, which would also circumvent prompt checks [80, 123]. Removing all humans from training datasets would mitigate both CSAM in training data and generated imagery. Such an intervention would also mitigate some other deepfake harms, including deepfakes of political figures [22]. Automated implementation on a large-scale dataset, however, would not guarantee complete removal, as existing person and face detection models have their own biases and error rates [113, 132].

C.4 Restrict dataset license to research-only usage

Gadre et al. [49] release CommonPool as an index of image url-text pairs under a Creative Commons CC-BY-4.0 license [35], which allows for the commercial usage of CommonPool, including the deployment of text-to-image generative models trained on this dataset. To avoid potential harms of future deployed models [121], one avenue is to limit usage to research purposes only or adopting licenses with behavioral use conditions like the RAIL license [36]. This approach does not remedy the presence of harmful content, but at least prevents models trained on these datasets from being used at scale.

C.5 Rebuild dataset from ground up

In general, we find that the filtering approach, in which undesirable content is iteratively removed from a web dump, results in the potential for false negatives. As we show in Section 7.3.3, a filtering mechanism may miss harmful content and thus include it in the final dataset. Another perspective to dataset collection is to rebuild the dataset from the ground up by adding only data known to be safe and ethically obtained. We point to the ante-hoc frameworks and recommendations from ethical dataset curation research [6, 10, 110]. An example of this implementation is the inclusion of data from moderated sources like reputable news or stock photo sites. At the same time, however, sources may still include harmful content, and this method may not address copyright concerns [70].

C.6 Retract the dataset

Large-scale image datasets have been taken down in the past as a result of the discovery of CSAM or NCII – for instance, LAION-5B [34] due to analysis by Thiel [120] and Tiny Images [64] based on the audit by Birhane and Prabhu [16]. While this may prevent future distribution of problematic content contained in the dataset, prior work shows that several previously-deprecated datasets have still been circulated via derivative copies and used in peer-reviewed research [75, 95],

We also note that the removal of open-access datasets would still allow the existence of similar proprietary datasets, which may have the same issues, to remain in usage without the ability for researchers to audit. In this manner, researchers would no longer have access to these datasets and therefore would not be able to study machine learning practices or models that may correspond to commercial approaches. We raise the question of what the goals of research on open datasets are, and when, if ever, the growing body of evidence of the fundamental flaws of web-scraped datasets will be enough.

Received 15 April 2024