

Sound Check: Auditing Recent Audio Dataset Practices

William Agnew¹, Julia Barnett², Annie Chu², Rachel Hong³, Michael Feffer¹, Robin Netzorg⁴, Harry H. Jiang¹, Ezra Awumey¹, Sauvik Das¹

¹Carnegie Mellon University, ²Northwestern University, ³University of Washington, ⁴University of California Berkeley
wagnew@andrew.cmu.edu

Abstract

Audio AI models are increasingly used for a broad range of applications including music and sound generation, text-to-speech (TTS), voice cloning, emotion analysis, transcription, and audio classification. However, we have little understanding of the datasets used to create audio AI models, a gap that leaves the field without a powerful tool for understanding potential biases, toxicity, copyright violations, and other ethical and performance issues of downstream models. To bridge this gap, we conduct a mapping literature review of hundreds of audio datasets used in recent music, sound, and speech AI papers. We first assess sourcing, size, and usage of these datasets, finding that while there are hundreds of audio datasets, few are widely used. Next, we identify nine representative datasets, and conduct several analyses to understand biases, toxicity, representation, and quality. We find that these datasets are often biased against women, have stereotypes about marginalized communities, and contain significant amounts of copyrighted work. We also find that audio datasets often come with scant documentation. To address this gap, we extend Gebru's datasheets for datasets (Gebru et al. 2021) to audio data, providing domain-specific documentation guidance. Finally, to facilitate public exploration of dataset contents and accountability, we developed an audio datasets exploration web tool available at <https://audio-audit.vercel.app/>. All code for this project is available at <https://anonymous.4open.science/r/gen-audio-ethics-0F57/README.md>.

Content warning: this paper contains discussions of offensive language.

1 Introduction

Deep learning and ML-based techniques achieve state-of-the-art performance on a broad range of audio processing tasks, including speech transcription (Radford et al. 2023), pitch estimation (Kim et al. 2018; Liu et al. 2023b), and acoustic event classification (Wang et al. 2022). Beyond solving foundational problems in audio processing, these technologies support an increasing number of higher-level human-AI interactions. For example, virtual avatars, assistive technologies for the visually impaired, and novel UI paradigms use audio AI tools to improve user

experiences (Danielescu et al. 2023; Upadhyay et al. 2023; Yu, Parde, and Chattopadhyay 2023). More recently, generative AI has led to the development of audio-based technologies for tasks including making reservations and generating music from text (Jacques et al. 2019; Wu et al. 2024).

Audio AI has been implicated and theorized in a range of ethical and societal concerns. Emotion analysis technologies have been deployed to make automated hiring decisions, with outcomes potentially worsened by bias (Akselrod and Venzke 2023). Voice actors and musicians have alleged that their intellectual property (IP) is being improperly used in audio datasets, in particular by the music generation startups Suno and Udio (Newton-Rex 2024a,b), the voice generation startup ElevenLabs (Soni 2024), and OpenAI, which allegedly copied Scarlet Johansson's voice without permission (Allyn 2024). Relatedly, voice AI, created using large and broad speech datasets, are able to accurately clone voices, and have been used in financial scams (Granda 2024) and misinformation campaigns (Bond 2024; ANI 2024).

Our paper contributes a broad analysis of audio data currently used to train AI models spanning the domains of music, speech, and sound. This work is motivated by the fact that while the downstream ethical risks and harms of generative text and vision models have been the subject of significant prior work (Birhane et al. 2024a; Birhane, Prabhu, and Kahembwe 2021; Bianchi et al. 2023; Hong et al. 2024), there has been comparatively little focus and understanding of these issues in the context of generative audio, leading to a “documentation debt” (Bender et al. 2021) where widely used audio datasets are often poorly documented and understood. Inspired by audits of vision and text datasets that helped crystallize discussion of the ethical harms present in those generative AI modalities, in this paper we make the following contributions:

- We chart the recent landscape of audio datasets, systematizing their creation, size, usage, and submodality through a mapping literature review of hundreds of audio AI papers.
- We uncover issues with bias, representation, improper use of intellectual property, and toxicity in nine widely used and representative audio datasets.
- We extend datasheets for datasets (Gebru et al. 2021) to

audio datasets to help address their longstanding documentation debt.

- We create a web-based searchable index of these datasets to enable further exploration by the public.

2 Background

Before detailing the methods and results of our audio dataset and model audits, we first briefly describe prior work in audio AI, research regarding ethical considerations of such developments, and work in dataset audits more generally.

2.1 Audio AI

To derive models from audio datasets, a number of methods have been introduced in the audio domain, mirroring developments in deep learning for text and images. This includes RNNs (Sturm et al. 2019), CNNs (Team 2019; Huang et al. 2017), combinations of the two (Donahue, Lipton, and McAuley 2017), and a recent turn to techniques leveraging transformer architectures (Agostinelli et al. 2023; Donahue et al. 2023; Garcia et al. 2023) and diffusion (Forsgren and Martiros 2022; Wang et al. 2023), especially as generative AI (GenAI) has captured public interest. The main architectures used in audio modeling are similar to those utilized in image and text, with adjustments made to handle the specifics of audio or speech data, such as the Audio Spectrogram Transformer (Gong, Chung, and Glass 2021), which utilize frequency representations of audio (e.g., spectrograms) to better model the different audio modalities. Large audio language models often will combine pre-trained large language models (LLMs) with audio-specific encoders to extend LLM capabilities to audio modalities (Tang et al. 2024; Gong et al. 2023a; Chu et al. 2023). Generative audio approaches often borrow from LLMs (Borsos et al. 2023a) or diffusion models that act on spectrograms or waveforms directly (Liu et al. 2023a).

While there are many semantically distinct audio modalities, they largely fall into three categories: music, speech, and (environmental) sound—largely referred to as audio by the community. While historically modeling in each of these modalities has primarily fallen under separate fields, recent advances in large audio language models have seen these distinct modalities being united under single frameworks (Ghosh et al. 2024; Gong et al. 2023a). These models aim to perform classical tasks, such as Automatic Speech Recognition and Note Identification under a single modeling paradigm (Tang et al. 2024) and support applications spanning accessibility technologies to music generation. For further information on AI for audio, Civit et al. (2022) provide a review of music generation, Mehrish et al. (2023) review AI for speech processing, and Nogueira et al. (2022); Kelley and Dickerson (2020), and Palaniappan, Sundaraj, and Sundaraj (2014) review AI for sound processing.

2.2 Ethics and Societal Implications of Audio AI

In light of the recent turn to GenAI and improvements of other ML-based audio technologies, some researchers have started to grapple with corresponding ethical concerns and implications. However, as detailed by the literature review conducted by Barnett (2023) surveying generative audio research papers, few papers consider the potential negative impacts of their work. Even further, Morreale, Sharma, and Wei (2023) find that audio datasets are often created without permission of audio owners and creators. Some of these harms have started to be addressed, especially recently, such as training data attribution of generative audio models (Barnett, Garcia, and Pardo 2024; Bralios et al. 2024).

Shelby et al. (2023) highlight the sociotechnical nature of AI harms and emphasize that harms cannot exist independent of societal norms and structures—they have to exist in a set of systems. Ruha Benjamin (Benjamin 2019) sheds light on how technological harms (not unlike most societal harms) have a disproportionate effect on people of color; technologies were built for the people in power and “often adopt the default norms and power structures of society.” Audio harms are no exception; they can even be magnified when we do not understand the contents of the data. Modern voice cloning models boast of being able to capture diverse voices, but both being able and being unable to model voice archetypes like a “gay voice” comes with a host of both safety and representational harms (Sigurgeirsson and Ungless 2024). Being unable to model a particular voice archetype could result in harms of not representing all types of voices, while being able to could lead to harmful stereotyping and or require data collection that could put participants at risk. Audio AI has also raised concerns about representation, culture, and data rights. For example, the Māori people have accused OpenAI’s Whisper transcription AI of training on recordings of their language, *te reo* without permission, and then incorrectly transcribing *te reo*, potentially damaging its integrity (Mahelona et al. 2023a).

Audio deepfakes present a whole new set of harms separate from those already realized by visual and even video deepfakes. In 2023 alone, music featuring deepfake voices of popular artists went viral on social media (Coscarelli 2023; Feffer, Lipton, and Donahue 2023), prompting the music industry to start grappling with intellectual property concerns entailed by generative audio models (Hoover 2023; Johnson 2023; Patel 2023; Sisario 2024) and even take down online communities where deepfake audio was proliferating (Hook 2023). Audio deepfakes are particularly dangerous in the case of phishing and fraud, where bad actors can impersonate voices with high believability and deceive people or even bypass voice security systems (Habib et al. 2019; Sisman et al. 2020; Kim, Kim, and Yoon 2022). Many audio papers, especially text-to-speech papers (TTS), note the potential for misuse in form of audio deepfake (Wang et al. 2020; Kim et al. 2020; Kim, Kim, and Yoon 2022) and some even noted they had no plans to release their models due to the strong potential for misuse via deepfake (Kim, Kim, and Yoon

2022). Hutiri, Papakyriakopoulos, and Xiang (2024) detail the specific harms inherent to speech generators such as voice clones of voice actors, “bringing back the dead,” and audio deepfakes of public figures. Battle-Roca et al. (2023) focus on the specific aspect of transparency within generative music and highlight the link between transparency and creativity, originality, and ownership of AI-generated music, suggesting that we should move towards more transparent AI-based music generation. Within TTS there are questions with regards to liability of harmful speech (Henderson, Hashimoto, and Lemley 2023) as well as harms of the reverse—a high potential for hallucination in speech-to-text (Koenecke et al. 2024) with an estimate of about 1% of audio transcriptions being entirely hallucinated.

Other ethical quandaries remain regarding the contribution of these models to climate change (Douwes, Esling, and Briot 2021; Holzapfel, Kaila, and Jääskeläinen 2024), speaker privacy and security (O’Reilly et al. 2024; Champion 2024), creativity (Khosrowi, Finn, and Clark 2023), GenAI’s effect on music creators as a whole (Barnett 2023; Lee et al. 2022) and the ethics of using voice synthesis on deceased people (Feffer, Lipton, and Donahue 2023; Lee et al. 2022). Beyond risks for misinformation and economic harms to artists, recent high-profile instances of fraud (e.g., the transfer of millions of dollars to scammers leveraging GenAI to deceive targets (Lo 2024; Milmo 2024)), physiognomy (e.g., gender and sexual orientation classification (Lee et al. 2024)), and surveillance (e.g., gunshot detection for predictive policing (Crocco et al. 2016)) illustrate the real-world privacy, security, and ethical risks of these technologies.

2.3 Dataset Audits

Audits of datasets have proven vital for understanding the behavior and forecasting biases, toxicity, and other harms of downstream models. Prabhu and Birhane (2021) found that the 80 Million Tiny Images dataset contained racist and non-consensual intimate imagery (NCII), (Johnson 2020), and Birhane, Prabhu, and Kahembwe (2021) and Thiel (2023) uncovered evidence of child sexual abuse material (CSAM) in the LAION5B text-image dataset (Schuhmann et al. 2022), leading to removal of these datasets (Johnson 2020; Cole 2023). While dataset audits incorporate a variety of methods and aims—representation, toxicity, privacy, or copyright concerns—they all help determine how the targeted dataset’s contents align with expectations in efforts to achieve accountability (Birhane et al. 2024b). Paullada et al. (2021) surveyed dataset audits and found they reveal representational harms and the presence of problematic content overlooked during data curation. Despite the impact of audits, Bender et al. (2021) argue that machine learning faces a dataset “documentation debt,” with popular datasets having little if any documentation. Audio suffers acutely from documentation debt, with very few analyses of audio datasets outside of their suitability for increasing technical performance, with the notable exceptions of bias and representation audits of Mozilla Common Voice (Shuyo 2014) and Leschanowsky

et al. (2024)’s audit of speaker recognition datasets used between 2012 and 2021. Our paper builds on these analyses to include recently used datasets across three major subgenres of audio—music, speech, and sounds—and concerns including bias, representation, stereotypes, data quality and quantity, data sourcing, and copyright.

3 Mapping Review of Current Audio Datasets and Models

To understand how many audio datasets exist, the distribution of their usage in the research community, and how these datasets were sourced, we conducted a mapping review enhanced with systematic elements (Ferrari 2015) utilizing the STAMP sampling method (Rogge et al. 2024) on audio modeling papers submitted to arXiv, a preprint platform previous studies have found to be an effective source for current and important audio AI papers (Barnett 2023). We searched for papers uploaded between May 1 2023 and May 1 2024 to capture one year of data and annotated the datasets included in these papers. We chose this time frame to capture the a recent usage of datasets in a field that has undergone rapid changes in recent years given fast growing academic and commercial interest in generative AI. We analyzed audio modeling papers until we approached dataset saturation (i.e., further analysis yielded few new datasets).

Our final corpus for the mapping review included 66 papers about music, 59 papers about speech, 19 papers about general audio (either environmental non-music, non-speech sounds, or general purpose audio), and five papers about music and speech (typically singing voice synthesis)—these categories are mutually exclusive. The authors then went through these papers and identified any audio datasets used for training or evaluation, yielding 175 unique datasets. Figure 5 in Appendix 2 provides a visual representation of this process.

3.1 Analysis of Current Audio Datasets

We analyzed the 175 audio datasets found through our mapping review in order to understand practices, uses, and creation methods. For each dataset, we noted the number of times it was used by papers in our corpus, the number of times it had been cited overall (beyond our sample), its size in hours, the categorization of its contents (music/speech/general sound), how its corresponding data was collected, and noted concerns related to potential copyright infringement (see Table 1). For further detail on how we annotated and calculated these features, see the Appendix 2.

Distribution and Usage of Datasets Of the 175 datasets, the vast majority of them were only used in one ($n = 99$; 57%) or two papers ($n = 45$; 26%). The full distribution can be found in Figure 6. Only a handful of datasets were used more than 5 times. Speech datasets had the largest skew: most datasets were only used by one paper, while VCTK (Yamagishi et al. 2019) was used by 14. We find a wide variation in length even among the most popular datasets: the Mozilla Common Voice dataset is nearly two orders

Overview of Datasets									
Category	Overall Datasets	Creation Method Scraped	Size (Hours)			Citations			Copyright Infringing
			Sum	Median	Mean	Sum	Median	Mean	
Music	61	36%	74,139	19	1,236	14,346	98	267	33%
Speech	80	24%	573,522	59	7,546	39,511	202	590	20%
Sounds	31	32%	37,178	64	1,377	11,998	159	461	35%
Music+Speech	3	0%	44	15	1	30	15	15	0%

Table 1: Descriptive statistics from 175 datasets identified in review. Split by music, speech and sound, we list the count of datasets, percent that were scraped, size in hours, total number of citations (beyond our corpus), and conservative estimate of the percent likely copyright infringing. All of this information was determined by two authors independently evaluating each dataset by reading the original papers proposing the datasets (when present), investigating all possible information provided online about the datasets, and lacking both of those downloading the dataset and personally assessing this information. For further detail on how we calculated hours and decided which datasets had potential copyright issues, see Appendix 2.

of magnitude larger than VCTK dataset, despite both being speech datasets. Speech datasets were also the largest by number of hours (see Appendix 2). We documented 573,522 hours of speech data (median = 59 hours), the vast majority of which came from VoxPopuli (Wang et al. 2021), a 400,000 hour dataset consisting of European parliament event recordings, and the Spotify Podcast Dataset (Clifton et al. 2020), 100,000 hours of Spotify Podcasts that has been removed from public access. Music datasets totaled 74,139 hours, with a similar skew (median of 19 hours) driven by The Million Song Dataset (Bertin-Mahieux et al. 2011) (50,000 hours), Irish Massive ABC Notation Dataset (Wu et al. 2023) (7,200 hours), and Free Music Archive (Defferrard et al. 2017) (5,920 hours). These findings stand in contrast to text and image modalities, where there exist a smaller number of very large, widely used datasets that have significant source overlaps (Schuhmann et al. 2022; Gadre et al. 2024; Raffel et al. 2020; Soldaini et al. 2024).

We also calculated the total citations these datasets had received¹ to gauge popularity relative to usage. As seen in Figure 1, both usage and citations are quite fragmented, but actual citations have a heavier focus on a few important datasets. Similar to usage and hours, speech dominated the total citation count receiving 39,511 cumulative citations (median = 202), with LibriSpeech (Panayotov et al. 2015) in the lead with 6,136 citations. Music was second with 14,346 cumulative citations (median = 98); GTZAN led music datasets in citations with 4,345 citations. Datasets including sounds were least cited with 11,998 cumulative citations (median = 159), with AudioSet being most popularly cited at 3,204 citations.

In contrast to citing public datasets, many papers in our corpus did not release the data they used. Out of the 157 papers in our sample, 65 papers used at least one proprietary dataset, and there were 77 proprietary datasets in total. Of these 77 datasets, 79% ($n = 61$) were not released, often without a rationale.

¹Citation data accessed from Google Scholar in Fall of 2024 and may not be exhaustive. Moreover, oftentimes when someone cites a dataset, they are doing so in acknowledgment of the field (e.g., in the related literature section) as opposed to actually using that data for training or evaluation.

Language Contents of Datasets When considering the linguistic diversity in speech datasets, the inclusion of underrepresented languages is frequently not prioritized, with many datasets predominantly featuring only one language. Out of the 77 speech datasets we examined, the majority (61) were monolingual, with 50 solely in English, followed by 7 in Mandarin. In contrast, 16 datasets encompassed between 2-30 languages, while only two datasets included more than 50 languages.

Sources of Datasets The two most salient audio data sources were YouTube ($n = 25$ datasets) and LibriVox ($n = 13$). Other standouts were freesound.org ($n = 6$), Spotify ($n = 4$), and VCTK (Yamagishi et al. 2019) ($n = 3$). Other popular sources for audio content included podcasts ($n = 6$), marketplace websites ($n = 4$), TED talks ($n = 4$), TV shows ($n = 3$), and parliament/public speeches ($n = 3$). We found that LibriVox and its derivatives were referenced in 35 out of the 59 speech papers included in our literature review. LibriSpeech (Panayotov et al. 2015), a dataset comprised of public domain audiobooks read by volunteers across various languages, serves as a foundation for 13 derivative datasets, including both direct derivatives like the LibriSpeech dataset (Panayotov et al. 2015) and Musan (Snyder, Chen, and Povey 2015), as well as derivatives of derivatives like LibriLight (Kahn et al. 2020).

The most cited derivative datasets include LibriSpeech ($n = 12$), LJSpeech ($n = 9$), LibriTTS ($n = 9$), and LibriLight ($n = 6$). This trend aligns with the overall popularity of these datasets, gauged by citation numbers. It is crucial to emphasize that LibriVox predominantly comprises century-old texts, encompassing outdated and potentially problematic language, cultural perspectives, and social norms, which we corroborate in our audit in Section 4. Researchers utilizing LibriVox should be mindful of the potential introduction of bias and toxicity inherent in this dataset.

Takeaways from the Mapping Review The composition of audio datasets used by the research and commercial audio community is vastly different than that of text and vision. It is extremely fragmented—beyond a few notable datasets (e.g., VCTK (Veaux, Yamagishi, and Mac-

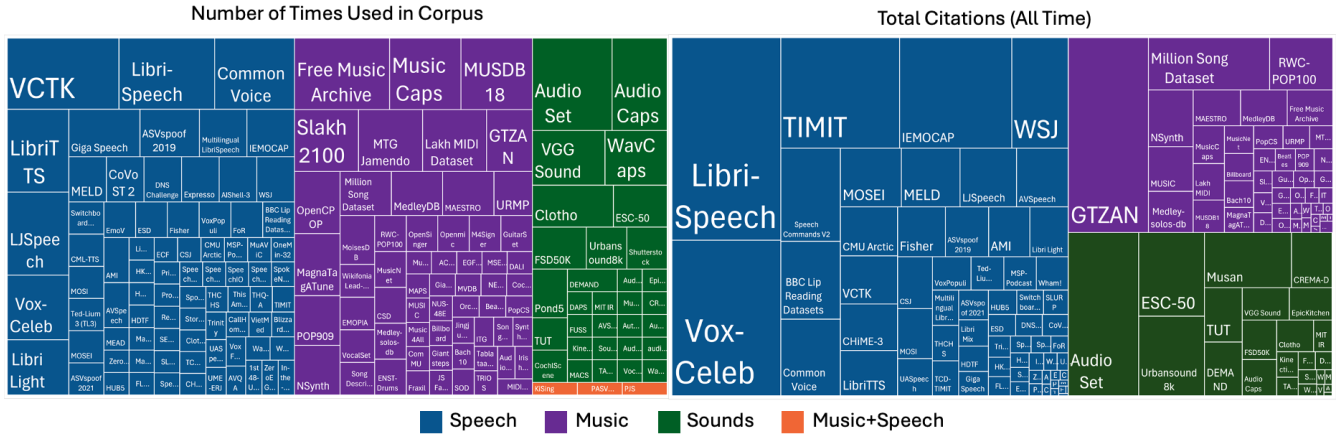


Figure 1: Area charts displaying the proportion of all 175 audited datasets by (1) number of times used in our paper corpus, (2) the cumulative total of citations received to date. Split into 4 categories: Speech, Music, Non-Music/Non-Speech sounds, Music and Speech combined.

Donald 2017) and AudioSet (Gemmeke et al. 2017)), researchers tended to use one-off, idiosyncratic datasets. These datasets also come with their own suite of problems; much of the contents in these datasets were likely sourced without creators’ knowledge or consent, potentially infringing copyright or distributing content at scales beyond the original understanding of data providers. There is no standardized way to prepare, create, release, or even discuss audio datasets used in research. Much of the work done in this literature review was brute-force data diving to understand the composition of datasets—future datasets should be prepared and released in a more mindful and documented manner. We discuss a suggested approach at the end of our paper, adapting datasheets for datasets (Gebu et al. 2021) to the audio domain.

4 Audit of Audio Datasets

Given the lack of documentation and sheer number of audio datasets in recent use, we next conduct deeper inspections of nine audio datasets that we identified as representative. We choose these datasets by identifying the five most frequently used and the five largest datasets in each submodality: music, speech, and sound ($n = 24$ total after removing duplicates). We removed datasets that are not freely accessible ($n = 3$: Audiosparx, Million Song Dataset, and Spotify Podcast), or that were drawn from the same sources as larger datasets in our audit ($n = 10$: AudioCaps, Auto-ACD, Clotho, Kinnectics 70, Librilight, LibriTS, LJSpeech, MusicCaps, SLakh 2100, and VGG Sound). $n = 5$ datasets appeared multiple times in the largest and most used rankings. We exclude Voxpopulli ($n = 1$), as while it is very large, only a small fraction of it is transcribed, limiting its usefulness for audio AI. We exclude musdb18 ($n = 1$) due to its small size (150 tracks), and exclude the Irish ABC ($n = 1$) due to its specificity and lack of live recordings limiting its usefulness for training general purpose music or audio models. We are left with nine remaining datasets ($n = 9$), representative of the largest and most popular au-

dio datasets in recent use across submodalities. We assess for bias, toxicity, representation, presence of potentially copyrighted content, and other important features for predicting behavior of downstream models.

1. *AudioSet* (Gemmeke et al. 2017), a dataset of 2 million 10-second YouTube clips, comprising a range of music, speech, and sounds;
2. *Mozilla Common Voice 17* (Ardila et al. 2020), a corpus of crowd-sourced sentences read by volunteers;
3. *VCTK* (Veaux, Yamagishi, and MacDonald 2017), a dataset of sentences read from the Herald, a Scottish newspaper, and several shorter accent elicitation texts;
4. *LibriVox* (LibriVox 2025), a dataset of volunteer recordings of public domain books;
5. *Free Music Archive* (Defferrard et al. 2017), a music-specific dataset scraped from an online repository;
6. *Jamendo* (Bogdanov et al. 2019), another scraped music-specific dataset,
7. *Wav-Caps* (Mei et al. 2024), a dataset of sounds sourced from AudioSet (Youtube), The BBC Sounds Effects library, FreeSound, and SoundBible;
8. *GigaSpeech* (Chen et al. 2021) a speech dataset sourced from YouTube, audiobooks, and podcasts; and lastly
9. the *Lakh MIDI Dataset* (Raffel 2016), a MIDI subset of the Million Song Dataset (Bertin-Mahieux et al. 2011), a currently unavailable dataset of music taken by The Echo Nest, a now defunct music analytics company.

4.1 Overview of Audio Datasets

Creation Dates Text-based audio datasets have two relevant creation dates: audio creation date and text creation date. VCTK was recorded in 2013 and features newspaper articles up to 2013 (Veaux, Yamagishi, and MacDonald 2017). LibriVox recordings were made between 2005 and present day (LibriVox 2025), but rely on public domain texts that typically enter the public domain 70 years

after the death of their last living author (of California 2024), and are thus often over a century old. The Lakh MIDI dataset is derived from the Million Song Dataset, composed of songs released before 2012 (Raffel 2016). Similarly, AudioSet is composed of YouTube videos released before 2016 (Gemmeke et al. 2017). Of the audited datasets, only Mozilla Common Voice features both contemporary audio recordings and text, with creation dates between 2017 and present. Training cutoff date, or the date of the most recent training data, has emerged as a crucial concern of LLMs, determining which events they may possess knowledge (Cheng et al. 2024). The relative age of most audio datasets could lead to downstream models performing worse on new, or newly popular, words, sounds, and genres, reflecting a bias towards the past (Birhane et al. 2022).

Identity Representation Who is represented in the datasets is a central question of dataset audits with significant downstream impacts on bias. In Table 2 we summarize available information about audio creator demographics. One notable finding is that audio datasets often have limited demographic information, and several important demographic categories—including sexual orientation, disability, race, and ethnicity—are not documented in any of the datasets we audited. To assess representation and bias of different demographics, we instead investigated identity keywords in audio transcripts. In Figure 2, we present counts of identity keywords, using keywords from prior work (Dodge et al. 2021), for each dataset’s transcripts, with plurals and common alternative spellings merged. With the exception of Mozilla Common Voice, we find “man” is used 2-10x more often than “woman” in these datasets. We also find that many identity keywords have a very low count in many datasets, with “Muslim” only appearing a thousand times, 3.5x less than “Christian” and “Nonbinary” appearing only 24 times. GigaSpeech and Mozilla Common Voice are the only datasets that contains at least one instance of every identity keyword studied. While this list of identity words is not comprehensive, and some words have alternate non-identity meanings, our results evidence low representation of marginalized groups in audio datasets.

Language Representation In Figure 3, we show the approximate number of hours of audio in each dataset for English and non-English languages, and in Table 5 we break down by language for non-English languages. We find that most audio datasets contain between 2x and 10x more data in English than non-English, with the exception of Mozilla Common Voice, which contains 24,424 hours of non-English data and 3,507 hours of English data.

Sociodemographic bias We assess binary gender bias by comparing words with the highest pointwise mutual information (PMI) difference between “woman” and “man” keywords across all datasets (Appendix Figure 10). A higher PMI indicates that words are more strongly associated with each other. We find “woman”-related words are more associated with terms about families and childcare

than “man”-related words, while “man”-related words are not correlated with typically gendered terms. In addition, we find that “woman”-related words have stronger associations with “baby,” “beauty,” and “b*tch.” while “man”-related words have stronger associations with “dead” and “power.” Overall, these associations provide evidence that women are commonly depicted in relation to families, childcare, and as subjects of the male gaze (Bloom 2017; Mulvey 2013). Nonbinary genders were not represented well enough in these datasets to assess bias. However, given the prevalence of biased and toxic content towards queer people online (Queerina et al. 2023), it is likely larger audio datasets will contain biased and toxic content towards these groups.

4.2 Toxicity

In this section we assess toxicity in the representative audio datasets. In Figure 8, we show the most common profane words in the datasets. We note that profanity is not the same as toxicity, and in some contexts these words are neither profane nor toxic. We find that, while all datasets contain at least some profane words, FMA, GigaSpeech and LibriVox contain by far the most, with many thousands of occurrences of racist and queerphobic terms. This finding may be due to LibriVox’s sourcing from public domain texts, which are at a minimum 70 years old and generally much older, and represent times where where toxic dialogues about marginalized populations were more overt.

In Figure 4, we present the number of hours of content in each dataset classified as toxic by the `pysentimento` toxicity classifier. This classifier considers not just profanity but additional textual cues to assess whether a text is toxic. Examples of and further discussion of sentences classified as toxic are available in Appendix A.4. While we find that each dataset has only approximately 1-3% of content flagged as toxic, this still amounts to hundreds of hours of toxic content. While levels of toxicity are low relative to the size of each dataset, LLMs have displayed an ability to recall text encountered only a few times during training (Carlini et al. 2022), raising the possibility that large audio models will exhibit similar behavior with even small amounts of profane or toxic content.

4.3 Audio Datasets Licensing

In Table 6, we summarize licenses of audio datasets. Mozilla Common Voice, VCTK, and LibriVox have permissive licenses that allow any use with minimal restrictions. The other six datasets all have licenses that potentially impact the ability of these datasets to be used for training or commercial applications. Lakh is derived from the Million Song Dataset, itself derived from Echo Nest, a music data service, which is subject to the Echo Nest License (Nest 2015), which prohibits commercial use. Lakh also contains many copyright tags, and we present the most frequent tags in Appendix Figure 11. AudioSet is derived from YouTube videos, which are licensed either under Creative Commons Licenses (YouTube 2024a) with different levels of permissiveness, or the YouTube License (YouTube

Dataset	Age	Gender	Sexual Orientation	Language	Locale/Country	Accent	Race/Ethnicity	Disability
Mozilla Common Voice	yes	yes	no	yes	yes	no	no	no
VCTK	yes	yes (binary)	no	yes	yes	yes	no	no
LibriVox	no	yes (binary)	no	yes	no	yes	no	no
Lakh	no	no	no	no	no	no	no	no
AudioSet	no	no	no	yes	yes	no	no	no
Free Music Archive	no	no	no	yes	yes	no	no	no
Jamendo	no	no	no	no	no	no	no	no
Wav-Caps	no	no	no	no	no	no	no	no
GigaSpeech	no	no	no	yes	no	no	no	no

Table 2: Documentation of demographics in audited datasets.

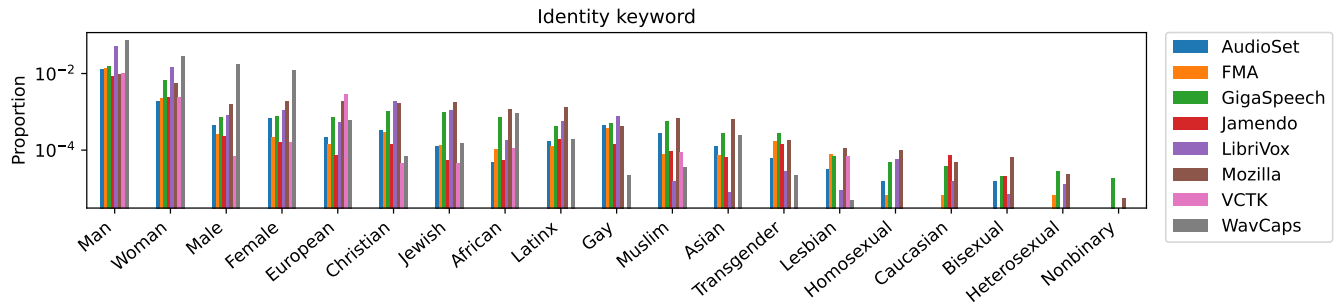


Figure 2: Proportion of identity keyword mentions for each dataset. Y-axis is in log-scale.

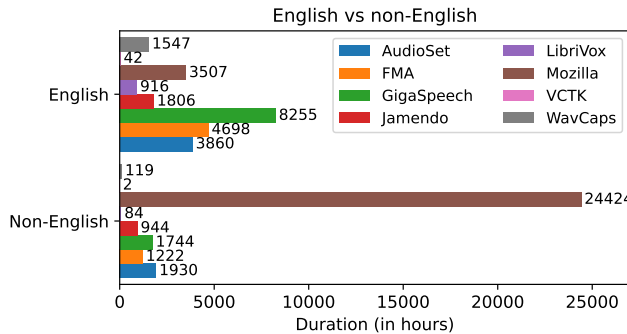


Figure 3: Estimated duration of each dataset by hours in English and not English.

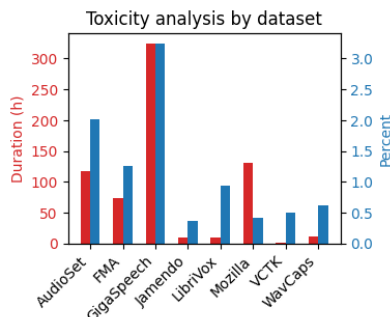


Figure 4: Estimated duration and percent of toxic-predicted English sentences split by dataset.

2024b), where the creator retains all ownership. However, YouTube can use or modify videos in connection with YouTube’s business, and users of YouTube can use or modify videos only as a feature of YouTube. We present the most common YouTube channel names in Figure 11 in our Appendix. Free Music Archive and Jamendo data are covered under several Creative Commons and other licenses, including many that prohibit commercial use and derivatives, require artist attribution, mandate that derivative works carry the same license, and restrict use to personal use only. Wav-Caps includes audio from YouTube under Creative Commons and Youtube licenses, and also content from the BBC Sound Effects Library under a BBC license that does not permit commercial use (BBC 2025). GigaSpeech includes content from YouTube under Creative Commons and Youtube licenses, and also a variety of different podcasts, each with different licenses (including the 99% Invisible license (Visible 2025) and Australian Broadcasting Company license (ABC 2025)), many of which do not permit commercial use. In short, our analysis uncovers extensive presence of copyrighted material in these datasets from a broad range of artists and creators.

4.4 Audio Content: Listening Audit

Audio data is meant to be heard—so we conducted a “listening audit” of the selected datasets. Our descriptive review offers a qualitative assessment of audio quality and content. Three authors independently reviewed samples from the seven primary datasets analyzed in this paper. All three evaluated samples from AudioSet (Gemmeke et al. 2017), while at least two authors reviewed samples from each of the remaining six datasets. After listening, authors discussed their observations to reach consensus.

Quality Among the speech datasets considered, VCTK (Veaux, Yamagishi, and MacDonald 2017) stood out as very high quality, though exhibits minor issues, as some samples contain reverberation, muffling, or instances of speakers stumbling over words. In contrast, the majority of other datasets exhibit relatively low quality, often characterized by significant background noise and, in some cases, recordings that are nearly unintelligible. Mozilla Common Voice (Ardila et al. 2020) exhibited notably lower audio quality compared to VCTK and LibriVox. Frequent background noise often overshadowed the speech, and environmental or recording artifacts—such as audible keyboard clicks used to stop recordings—were common, increasing the likelihood of these artifacts being learned by models. Wav-Caps and GigaSpeech had varying levels of quality, with some clips apparently created in recording studio, and others with low quality microphones in noisy conditions.

In contrast, music datasets demonstrated considerably higher audio quality, featuring stereo channels and a sampling rate of 44.1 kHz. Both FMA and MTG-Jamendo datasets included a mix of vocal and non-vocal tracks, with a strong representation of samples having digital audio production elements such as sampling, synthesizers, and applied audio effects. Among these, MTG-Jamendo (Bogdanov et al. 2019) offered exceptionally high-quality audio, while FMA (Defferrard et al. 2017) was of noticeably lower quality, often exhibiting environmental noise and a grainy or distorted sound. Another notable observation was the lack of normalization, leading to significant variation in loudness levels across samples. The Lakh MIDI (Raffel 2016) dataset also demonstrated high quality, though it consists of MIDI files rather than raw audio, limiting its direct comparison with other datasets.

Content The source material for LibriVox—public domain books—primarily consists of older English texts, resulting in a vocabulary set that is not representative of contemporary English usage. Moreover, utterance lengths vary wildly—ranging from single-word utterances to fragmented sentences. This abrupt splicing of sentences before their natural conclusion led to orators delivering them as though they were complete, even when they ended mid-clause. These recordings may misrepresent the natural flow of human prosodic speech, reducing their suitability for tasks that rely on realistic speech patterns. In addition, we found that the Free Music Archive dataset (developed for music) contains many environmental sounds.

AudioSet (Gemmeke et al. 2017) differs from the other datasets in that it is a collection of URLs and 10 second time stamps links to Youtube videos, each with associated content tags. Many clips are over a decade old, have under 1,000 views, and would seldom surface to contemporary Youtube users. This practical privacy is disrupted by their inclusion in AudioSet. This concern is especially pertinent given that some of these videos had children or minors talking to camera (e.g., vlogs), raising ethical questions surrounding consent for using these data to train models. Categorizations (tags) were frequently incor-

rect or inconsistent—for example, a video tagged “bicycle” might include a bicycle in the footage but no accompanying “bicycle”-like sound. Inconsistent tagging is also evident in instances where two audio samples with similar content were assigned tags with different levels of detail, such as a produced song with vocals being tagged as “Music” in one case and “Music, Singing, female singing, musical instrument” in another. We found some categorizations to be potentially offensive, in particular the category “funny music”² which we found applied to non-Western music. While “Speech” was a common tag, there was no distinction for language nor differentiation between background and foreground speech (e.g., speech meant to be intelligible vs. background conversation). Similarly, music is frequently tagged without differentiation between foreground music and incidental background music.

4.5 Privacy

Multiple datasets, including all datasets using YouTube, Free Music Archive, and Jamendo to source data, often contain metadata enabling determination of the names of speakers. Additional personally identifiable information is often available through metadata, audio, or—particularly on YouTube—video, including faces and locations. We found many instances of YouTube videos of children being included in datasets. Most critically, voiceprints are foundational to many of these datasets, and they contain voiceprints of tens of thousands of people. Voiceprints present with other PII, particularly names, are especially concerning, raising the possibility audio AI models could inadvertently learn to reproduce the voices of people in their training datasets, or malicious actors could use these datasets to impersonate the people within them.

5 Discussion

In this paper, we identified 175 audio datasets that have been recently used for generative audio AI. The distribution of their use is long-tailed, with a small number being used frequently, and a majority being used only once or twice. Many audio datasets were created specifically for research purposes, but approximately one third were scraped from online sources. These scraped datasets carry various licenses—often Creative Commons versions that bar commercial use, remixing, or require attribution—complicating their use for AI training. The corpora also contain many instances of profanity and toxicity. Sparse documentation and metadata obscure who is represented in audio datasets, highlighting a need for better documentation practices. Through a content analysis, however, we found that marginalized communities were less likely to be explicitly mentioned in audio datasets.

5.1 Impacts of Bias and Toxicity in Audio Datasets on Downstream Models

Our analysis indicates that the datasets we analyze in depth are biased both statistically (i.e., skewed inclusion

²https://research.google.com/audioset/dataset/funny_music.html

of certain types of data) and socially (i.e., certain keywords and terms correlated with identity), in addition to containing non-trivial amounts of profane and potentially toxic content. Given that models in other modalities have displayed an ability to recall data encountered only a few times during training (Carlini et al. 2022), this raises the possibility that large audio models will exhibit similar behavior with even small amounts of profane or toxic content. These issues can also lead to disparate performance across socially salient categories. For instance, we observe heavy emphasis on the English language across datasets. Downstream of these datasets, Radford et al. (2023) note that while their Whisper transcription model obtains state-of-the-art performance on several tasks, they are limited in that “Whisper’s speech recognition performance is still quite poor on many languages” due to their “pre-training dataset [being] currently very English-heavy due to biases of [their] data collection pipeline.” Disparate performance of transcription via audio models across languages, and by extension cultures, is thus already documented as an impact of representation bias in audio datasets (Fuckner et al. 2023; Nacimientogarcía, Díaz-Kaas-Nielsen, and González-González 2024).

Problems extend beyond transcription: translation, voice recognition, and audio generation face similar risks. Our PMI, identity keyword, and profanity analyses reveal gender associations with certain words, virtual omission of keywords pertaining to sexual orientation, and a high incidence of profanity related to sexual anatomy and racial slurs. If not properly addressed, these features of audio datasets may yield downstream audio generation models that perpetuate or even accelerate instances of social bias and toxicity that echo harmful stereotypes found in the wild (e.g., models that describe or portray women in a stereotyped way, models that are unable to create outputs relating to the LGBTQ experience, models that degenerate into profanity based on sexual anatomy or racial slurs without provocation, etc.).

Projects like the Common Voice corpus (Ardila et al. 2020), the Corpus of Regional African American Language (Kendall and Farrington 2023), and the Mid-Atlantic Gender Expansive Speech Corpus (Hope, Ward, and Lilley 2023) aim to improve the diversity of speech datasets. However, ethical dataset creation and use goes far beyond mere inclusion, especially in the context of AI and big tech—it can often be predatory, harming included communities by exposing them to data- and AI-intensified surveillance, poorly performing AI, and loss of control over data and culture (Mahelona et al. 2023b). Auditing datasets to uncover biases and a lack of representation is the start of the conversation, but effective solutions must be led by those who own the data (Mahelona et al. 2023b).

5.2 Scale of Audio Datasets

Audio datasets require considerably more storage per sentence than text datasets. AudioSet totals nearly 2.5TB in size and contains approximately 52 million spoken words. In contrast, C4 (Raffel et al. 2020), a text dataset approximately 30% of the (file) size of AudioSet, contains 153

billion words (Soldaini et al. 2024). Including the same breadth and depth of content in audio datasets as text datasets will require significantly more storage and compute which can both exacerbate known harms for large generative models, and create new ones. In particular, the requirement for more storage and compute limits which actors can train on such datasets. While most of the datasets of corporations training the largest audio models remain closed to external auditors, OpenAI Whisper (Radford et al. 2023) was trained on 680,000 hours of audio, and the BigSSL speech recognition model (Zhang et al. 2022) from Google was trained on one million hours of YouTube audio. Both of these datasets are over an order of magnitude larger than the largest freely accessible dataset in our survey. Thus, generative audio models may impose significantly higher barriers for participatory approaches to AI development than generative text models.

Finally, while the high carbon and water impacts (Cridle and Bryan 2024) of LLM training (Luccioni, Viguier, and Ligozat 2023) and inference (Everman et al. 2023; Luccioni, Jernite, and Strubell 2024) are well known, the scale of audio datasets raises concerns that audio and multimodal models will have even higher pollution and resource costs compared to purely text-based models (Douwes, Esling, and Briot 2021).

5.3 Source of Audio Datasets

The existence of proprietary datasets vastly larger than open datasets raises questions about the sources of these datasets. Google has explicitly sourced massive audio datasets from YouTube (Zhang et al. 2022), and Spotify released (and then removed) a 100,000 hour transcribed dataset from hosted podcasts (Clifton et al. 2020). While OpenAI has not indicated how the Whisper dataset was sourced, nor are they publicly known to host massive audio datasets, the few indications of the source of massive audio datasets we have point to existing commercial hosting and streaming corporations (Davis 2024).

Corporations once freely released massive audio datasets. Google released AudioSet, over 5,000 hours of YouTube audio, in 2017. However, AudioSet was released as links to YouTube videos that could be independently downloaded to obtain audio. When we tried to download this dataset in May 2024, we found YouTube had rate limiters and crawler IP blocking that would make downloading this dataset take several months. Similarly, Spotify released 100,000 hours of transcribed podcast audio in 2020, but then removed this dataset in December 2023, citing “shifting priorities” (Johnson 2023). Most explicitly, Universal removed its music from Tiktok in January 2024, partially citing concerns over AI generated music and AI covers of its songs (Hoskins 2024). In our analysis of current audio datasets (Section 3.1), there were 77 proprietary datasets utilized, of which 79% ($n = 61$) were not released. We argue these events constitute a pattern of increasing restrictiveness to proprietary audio datasets, mirroring recent trends in text data sources (Longpre et al. 2024). These restrictions seem in part motivated by the new value of massive audio datasets for

training generative audio AI. As Ojewale et al. (2024) have noted, restrictions on dataset access is a major barrier auditors face, making the trend towards closed-source and proprietary datasets especially concerning.

5.4 Impacts to Rights and Intellectual Property

Any recording of a person’s voice is a biometric identifier that can be misused for imitation, deception, or right-of-publicity violations. We found major audio datasets contained a wide array of content licensed under terms restricting commercial use, including audio from YouTube, the BBC, and a wide array of podcasts. A generated voice in an OpenAI product released in 2024 closely resembled that of actress Scarlett Johansson; she has since threatened legal action before the product, Sky, was quietly dropped (Allyn 2024; De Vynck 2024). The ability to mimic voice actors, musicians, and other artists raises concerns of economic and labor harms, where AI is used to undercut artists with their own data. With the threat to right of publicity greatly increased through GenAI, individuals in the US may seek recourse in the precedent of *Midler v. Ford Motor Co.* which ruled that voice may form part of an individual’s identity, and that imitation of voice without approval is unlawful, though the specific boundaries of what consists of protected “identity” remains murky (Lapter 2007). Recently proposed US Senate legislation seeks to protect the public from generative AI; if passed, the bipartisan 2024 No AI Fakes Act would serve to protect individuals’ voices and likenesses from AI deep-fakes (Salazar 2024).

6 Recommendations

As initial steps towards ensuring harms that have plagued other modalities, including bias, toxicity, and intellectual property concerns, are not also present in future audio AI, we recommend (1) audio dataset developers adopt improved documentation to enable better assessment of bias and representation and (2) audio dataset developers only use data that permits remixing and commercial use at a minimum, but ideally seek active and informed consent for usage in AI. Both of these aims may be achieved by continuing existing practices of creating datasets in-lab specifically for AI.

6.1 Datasheets for Audio Datasets

To enable effective audio AI dataset documentation, we adapt and extend Gebru et al.’s “datasheets for datasets” (Gebru et al. 2021) and Papakriakopoulos et al. (2023)’s augmented datasheets for speech datasets to the context of audio to guide ethical audio dataset development, documentation, and use. In particular, we add several questions that specifically speak to contents of and representation in audio data, in addition to questions assessing data provenance, consent, and copyright. This datasheet is intended to serve both as a reflexive practice for audio dataset creators, documentation standards for audio dataset publishers, and a guide for future audio dataset auditors. We provide the full datasheet in Appendix A.1.

6.2 Educating and Mobilizing Data Workers

Documenting existing datasets is insufficient to improve dataset practices. As Birhane et al. (2024b) note in their review of audits, “Most of the academic work we reviewed focused on the process of evaluating AI systems for bias, fairness, or disparate impacts. Conversely, these studies rarely focused on other stages of auditing crucial to accountability in non-academic work, such as discovering harms, communicating audit results, or organizing non-technical interventions and collective action.” Artists, creators, and other data workers learning about inclusion of their data in AI datasets have been a key catalyst for these communities discussing their wants and needs in regards to inclusion in AI datasets and subsequently organizing and advocating (Marx 2024). To support this collective action, we created <https://audio-audit.vercel.app/>, a website that enables anyone to search for their inclusion in prominent audio datasets. Modeled after dataset search tools used to assess inclusion of data in AI datasets (Willison 2023a, 2022b,a, 2023b), this website enables users to understand how their work is being used in audio AI.

7 Conclusion, Limitations, and Future Work

In this paper, we conducted a large-scale survey of audio datasets that are used in generative audio models: an audit broadly inclusive of speech, music, and sound datasets. We found that these datasets exhibit similar patterns of bias and toxicity as text and image datasets, raising the concern that audio models could, in turn, exhibit similar levels of bias and toxicity as LLMs and image models if these risks are not mitigated. We found that hundreds of audio datasets are in use, with several open datasets being significantly larger and more widely used than others. However, we found indications that datasets have recently started becoming more closed and commercial, with past sources of massive datasets, including YouTube and Spotify, taking down datasets or implementing new measures to block crawling of their audio repositories. The widespread presence of copyrighted material in many of these audio data sources, frequent use of proprietary audio datasets in research by corporations, and new commercial perceptions raise serious concerns for musicians, voice actors, and other audio workers releasing content online, especially on platforms like YouTube or Spotify.

Our audit has several limitations and is just the start of addressing the dataset documentation debt (Bender et al. 2021) in audio AI. We were unable to assess the demographics of the people in audio datasets, and each of music, speech, and sound require specific and unique audits and analyses. We were only able to conduct high-level analyses of the contents of a small number of datasets, and each submodality—speech, music, and sound—have important distinctions that merit specialized treatment. We hope future audits will provide a deeper understanding of acoustic qualities of audio datasets, and that audio dataset curators will take steps to diversify their datasets, mitigate biases and toxicity, and remove copyrighted and other non-consensually sourced material.

References

- ABC. 2025. Terms of Use.
- Agostinelli, A.; Denk, T. I.; Borsos, Z.; Engel, J.; Verzetti, M.; Caillon, A.; Huang, Q.; Jansen, A.; Roberts, A.; Tagliasacchi, M.; Sharifi, M.; Zeghidour, N.; and Frank, C. 2023. Musi-
cLM: Generating Music From Text. *arXiv:2301.11325*.
- AI, S. 2022. The obscenity list.
- Akselrod, O.; and Venzke, C. 2023. How Artificial Intelligence Might Prevent You From Getting Hired.
- Allyn, B. 2024. Scarlett Johansson says she is 'shocked, angered' over new ChatGPT voice.
- ANI. 2024. Ranveer Singh's deepfake endorsing political party goes viral; actor flags alert on Instagram.
- Ardila, R.; Branson, M.; Davis, K.; Kohler, M.; Meyer, J.; Henretty, M.; Morais, R.; Saunders, L.; Tyers, E.; and Weber, G. 2020. Common Voice: A Massively-Multilingual Speech Corpus. In *The Language Resources and Evaluation Conference (LREC)*, 4218–4222. Marseille, France: European Language Resources Association.
- arXiv. 2023. About arXiv.
- Barnett, J. 2023. The ethical implications of generative audio models: A systematic literature review. In *Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society*, 146–161. Montreal, Canada: Association for Computing Machinery (ACM).
- Barnett, J.; Garcia, H. F.; and Pardo, B. 2024. Exploring Musical Roots: Applying Audio Embeddings to Empower Influence Attribution for a Generative Music Model. *arXiv:2401.14542*.
- Battle-Roca, R.; Gómez, E.; Liao, W.; Serra, X.; and Mitsu-fuji, Y. 2023. Transparency in music-generative AI: A systematic literature review.
- BBC. 2025. Licensing.
- Bender, E. M.; Gebru, T.; McMillan-Major, A.; and Shmitchell, S. 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*, 610–623. Virtual Event, Canada: Association for Computing Machinery (ACM).
- Benjamin, R. 2019. *Race after technology: Abolitionist tools for the new Jim code*. John Wiley & Sons.
- Bertin-Mahieux, T.; Ellis, D. P.; Whitman, B.; and Lamere, P. 2011. The Million Song Dataset. In *Proceedings of the 12th International Conference on Music Information Retrieval (ISMIR 2011)*, 591–596. Miami, US: International Society for Music Information Retrieval.
- Bianchi, F.; Kalluri, P.; Durmus, E.; Ladhak, F.; Cheng, M.; Nozza, D.; Hashimoto, T.; Jurafsky, D.; Zou, J.; and Caliskan, A. 2023. Easily accessible text-to-image generation amplifies demographic stereotypes at large scale. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*, 1493–1504.
- Birhane, A.; Dehdashtian, S.; Prabhu, V.; and Boddeti, V. 2024a. The Dark Side of Dataset Scaling: Evaluating Racial Classification in Multimodal Models. In *The 2024 ACM Conference on Fairness, Accountability, and Transparency*, 1229–1244.
- Birhane, A.; Kalluri, P.; Card, D.; Agnew, W.; Dotan, R.; and Bao, M. 2022. The values encoded in machine learning research. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, 173–184.
- Birhane, A.; Prabhu, V.; Han, S.; Boddeti, V. N.; and Luc-cioni, A. S. 2023. Into the LAION's den: Investigating hate in multimodal datasets. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, 21268–21284. New Orleans, USA: Curran Associates, Inc.
- Birhane, A.; Prabhu, V. U.; and Kahembwe, E. 2021. Mul-timodal datasets: misogyny, pornography, and malignant stereotypes. *arXiv preprint arXiv:2110.01963*.
- Birhane, A.; Steed, R.; Ojewale, V.; Vecchione, B.; and Raji, I. D. 2024b. AI auditing: The broken bus on the road to AI accountability. In *2024 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML)*, 612–643. IEEE, Toronto, Canada: IEEE.
- Bloom, J. D. 2017. *Reading the male gaze in literature and culture: Studies in erotic epistemology*. Cham, Switzerland: Springer.
- Bogdanov, D.; Won, M.; Tovstogan, P.; Porter, A.; and Serra, X. 2019. The MTG-Jamendo Dataset for Automatic Mu-sic Tagging. In *Machine Learning for Music Discovery Workshop, International Conference on Machine Learning (ICML 2019)*. Long Beach, CA, United States: PMLR.
- Bond, S. 2024. A political consultant faces charges and fines for Biden deepfake robocalls.
- Borsos, Z.; Marinier, R.; Vincent, D.; Kharitonov, E.; Pietquin, O.; Sharifi, M.; Roblek, D.; Teboul, O.; Grang-ier, D.; Tagliasacchi, M.; et al. 2023a. Audioldm: a lan-guage modeling approach to audio generation. *IEEE/ACM transactions on audio, speech, and language processing*, 31: 2523–2533.
- Borsos, Z.; Sharifi, M.; Vincent, D.; Kharitonov, E.; Zeghi-dour, N.; and Tagliasacchi, M. 2023b. Soundstorm: Efficient parallel audio generation. *arXiv preprint arXiv:2305.09636*.
- Bralios, D.; Wichern, G.; Germain, F. G.; Pan, Z.; Khu-rana, S.; Hori, C.; and Le Roux, J. 2024. Generation or Replication: Auscultating Audio Latent Diffusion Models. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 1156–1160. IEEE, Seoul, South Korea: IEEE.
- Carlini, N.; Ippolito, D.; Jagielski, M.; Lee, K.; Tramer, F.; and Zhang, C. 2022. Quantifying memorization across neural language models. *arXiv preprint arXiv:2202.07646*.
- Champion, P. 2024. Anonymizing Speech: Evaluat-ing and Designing Speaker Anonymization Techniques. *arXiv:2308.04455*.
- Chen, G.; Chai, S.; Wang, G.; Du, J.; Zhang, W.-Q.; Weng, C.; Su, D.; Povey, D.; Trmal, J.; Zhang, J.; et al. 2021. Gigaspeech: An evolving, multi-domain asr corpus with 10,000 hours of transcribed audio. *arXiv preprint arXiv:2106.06909*.

- Cheng, J.; Marone, M.; Weller, O.; Lawrie, D.; Khashabi, D.; and Van Durme, B. 2024. Dated Data: Tracing Knowledge Cutoffs in Large Language Models. *arXiv preprint arXiv:2403.12958*.
- Chu, Y.; Xu, J.; Zhou, X.; Yang, Q.; Zhang, S.; Yan, Z.; Zhou, C.; and Zhou, J. 2023. Qwen-Audio: Advancing Universal Audio Understanding via Unified Large-Scale Audio-Language Models. *arXiv:2311.07919*.
- Civit, M.; Civit-Masot, J.; Cuadrado, F.; and Escalona, M. J. 2022. A systematic review of artificial intelligence-based music generation: Scope, applications, and future trends. *Expert Systems with Applications*, 209: 118190.
- Clifton, A.; Pappu, A.; Reddy, S.; Yu, Y.; Karlgren, J.; Carterette, B.; and Jones, R. 2020. The Spotify Podcast Dataset. *arXiv:2004.04270*.
- Cole, S. 2023. Largest Dataset Powering AI Images Removed After Discovery of Child Sexual Abuse Material.
- Commons, C. 2024a. CC BY 4.0.
- Commons, C. 2024b. CC0.
- Commons, O. D. 2024c. Open Data Commons Attribution License (ODC-By) v1.0.
- Coscarelli, J. 2023. An A.I. Hit of Fake ‘Drake’ and ‘The Weeknd’ Rattles the Music World. *The New York Times*.
- Criddle, C.; and Bryan, K. 2024. AI boom sparks concern over Big Tech’s water consumption. *Financial Times*.
- Crocco, M.; Cristani, M.; Trucco, A.; and Murino, V. 2016. Audio surveillance: A systematic review. *ACM Computing Surveys (CSUR)*, 48(4): 1–46.
- Danielescu, A.; Horowitz-Hendler, S. A.; Pabst, A.; Stewart, K. M.; Gallo, E. M.; and Aylett, M. P. 2023. Creating Inclusive Voices for the 21st Century: A Non-Binary Text-to-Speech for Conversational Assistants. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, CHI ’23. New York, NY, USA: Association for Computing Machinery. ISBN 9781450394215.
- Davis, W. 2024. OpenAI transcribed over a million hours of YouTube videos to train GPT-4.
- De Vynck, G. 2024. OpenAI rolls out voice mode after delaying it for safety reasons.
- Defferrard, M.; Benzi, K.; Vandergheynst, P.; and Bresson, X. 2017. FMA: A Dataset For Music Analysis. In *18th International Society for Music Information Retrieval Conference*, 316–323. Suzhou, China: International Society for Music Information Retrieval.
- Deshmukh, S.; Han, S.; Bukhari, H.; Elizalde, B.; Gamper, H.; Singh, R.; and Raj, B. 2024. Audio Entailment: Assessing Deductive Reasoning for Audio Understanding. *arXiv:2407.18062*.
- Dodge, J.; Sap, M.; Marasović, A.; Agnew, W.; Ilharco, G.; Groeneveld, D.; Mitchell, M.; and Gardner, M. 2021. Documenting Large Webtext Corpora: A Case Study on the Colossal Clean Crawled Corpus. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 1286–1305. Online and Punta Cana, Dominican Republic: Association for Computational Linguistics.
- Donahue, C.; Caillon, A.; Roberts, A.; Manilow, E.; Esling, P.; Agostinelli, A.; Verzetti, M.; Simon, I.; Pietquin, O.; Zeghidour, N.; and Engel, J. 2023. SingSong: Generating musical accompaniments from singing. *arXiv:2301.12662*.
- Donahue, C.; Lipton, Z. C.; and McAuley, J. 2017. Dance dance convolution. In *International conference on machine learning*, 1039–1048. PMLR, Sydney, Australia: PMLR.
- Douwes, C.; Esling, P.; and Briot, J.-P. 2021. Energy Consumption of Deep Generative Audio Models. *arXiv:2107.02621*.
- Everman, B.; Villwock, T.; Chen, D.; Soto, N.; Zhang, O.; and Zong, Z. 2023. Evaluating the carbon impact of large language models at the inference stage. In *2023 IEEE international performance, computing, and communications conference (IPCCC)*, 150–157. IEEE, Anaheim, USA: IEEE.
- Feffer, M.; Lipton, Z. C.; and Donahue, C. 2023. Deepdrake ft. bts-gan and taylorvc: an exploratory analysis of musical deepfakes and hosting platforms. In *Proceedings of the 2nd Workshop on Human-Centric Music Information Retrieval*. Milan, Italy: CEUR Workshop Proceedings.
- Ferrari, R. 2015. Writing narrative style literature reviews. *Medical Writing*, 24 (4), 230–235.
- Forsgren, S.; and Martiros, H. 2022. Riffusion: Stable diffusion for real-time music generation.
- Fuckner, M.; Horsman, S.; Wiggers, P.; and Janssen, I. 2023. Uncovering Bias in ASR Systems: Evaluating Wav2vec2 and Whisper for Dutch speakers. In *2023 International Conference on Speech Technology and Human-Computer Dialogue (SpeD)*, 146–151. IEEE.
- Gadre, S. Y.; Ilharco, G.; Fang, A.; Hayase, J.; Smyrnis, G.; Nguyen, T.; Marten, R.; Wortsman, M.; Ghosh, D.; Zhang, J.; et al. 2024. Datacomp: In search of the next generation of multimodal datasets. *Advances in Neural Information Processing Systems*, 36.
- Garcia, H. F. E.; Seetharaman, P.; Kumar, R.; and Pardo, B. 2023. VampNet: Music Generation via Masked Acoustic Token Modeling. In *Ismir 2023 Hybrid Conference*, 359–366. Milan, Italy: International Society for Music Information Retrieval.
- Garofolo, J.; Graff, D.; Paul, D.; and Pallett, D. 1993. CSR-I (WSJ0) complete ldc93s6a. *Web Download*. Philadelphia: Linguistic Data Consortium, 83.
- Gebru, T.; Morgenstern, J.; Vecchione, B.; Vaughan, J. W.; Wallach, H.; Iii, H. D.; and Crawford, K. 2021. Datasheets for datasets. *Communications of the ACM*, 64(12): 86–92.
- Gemmeke, J. F.; Ellis, D. P. W.; Freedman, D.; Jansen, A.; Lawrence, W.; Moore, R. C.; Plakal, M.; and Ritter, M. 2017. Audio Set: An ontology and human-labeled dataset for audio events. In *Proc. IEEE ICASSP 2017*, 776–780. New Orleans, LA: IEEE.
- Ghosh, S.; Kumar, S.; Seth, A.; Evuru, C. K. R.; Tyagi, U.; Sakshi, S.; Nieto, O.; Duraiswami, R.; and Manocha, D. 2024. GAMA: A Large Audio-Language Model with Advanced Audio Understanding and Complex Reasoning Abilities. *arXiv:2406.11768*.

- Gong, Y.; Chung, Y.-A.; and Glass, J. 2021. AST: Audio Spectrogram Transformer. *arXiv:2104.01778*.
- Gong, Y.; Liu, A. H.; Luo, H.; Karlinsky, L.; and Glass, J. 2023a. Joint Audio and Speech Understanding. In *2023 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*.
- Gong, Y.; Luo, H.; Liu, A. H.; Karlinsky, L.; and Glass, J. 2023b. Listen, Think, and Understand. *arXiv preprint arXiv:2305.10790*.
- Granda, C. 2024. Fraudsters use voice-cloning AI to scam man out of 25,000.
- Habib, R.; Mariooryad, S.; Shannon, M.; Battenberg, E.; Skerry-Ryan, R.; Stanton, D.; Kao, D.; and Bagby, T. 2019. Semi-supervised generative modeling for controllable speech synthesis. *arXiv preprint arXiv:1910.01709*.
- Henderson, P.; Hashimoto, T.; and Lemley, M. 2023. Where's the Liability in harmful AI Speech? *J. Free Speech L.*, 3: 589.
- Holzappel, A.; Kaila, A.-K.; and Jääskeläinen, P. 2024. Green MIR?: Investigating computational cost of recent music-AI research in ISMIR. In *International Society for Music Information Retrieval Conference (ISMIR)*.
- Hong, R.; Agnew, W.; Kohno, T.; and Morgenstern, J. 2024. Who's in and who's out? A case study of multimodal CLIP-filtering in DataComp. *arXiv:2405.08209*.
- Hook, S. 2023. AI hub discord server hit with subpoena by RIAA.
- Hoover, A. 2023. Spotify Has an AI Music Problem—but Bots Love It. *Wired*.
- Hope, M.; Ward, C.; and Lilley, J. 2023. Nonbinary American English speakers encode gender in vowel acoustics. In *Proceedings of Interspeech*, 4713–4717. Dublin, Ireland: ISCA.
- Hoskins, P. 2024. Universal Music to pull songs from TikTok. *BBC*.
- Huang, C.-Z. A.; Cooijmans, T.; Roberts, A.; Courville, A.; and Eck, D. 2017. Counterpoint by Convolution. In *International Society for Music Information Retrieval (ISMIR)*, 211–218. Suzhou, China: International Society for Music Information Retrieval.
- Huang, C.-Z. A.; Vaswani, A.; Uszkoreit, J.; Shazeer, N.; Simon, I.; Hawthorne, C.; Dai, A. M.; Hoffman, M. D.; Dinulescu, M.; and Eck, D. 2018. Music transformer. *arXiv preprint arXiv:1809.04281*.
- Hutiri, W.; Papakyriakopoulos, O.; and Xiang, A. 2024. Not My Voice! A Taxonomy of Ethical and Safety Harms of Speech Generators. In *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '24, 359–376. New York, NY, USA: Association for Computing Machinery. ISBN 9798400704505.
- Jacques, R.; Følstad, A.; Gerber, E.; Grudin, J.; Luger, E.; Monroy-Hernández, A.; and Wang, D. 2019. Conversational Agents: Acting on the Wave of Research and Development. In *Extended Abstracts of the 2019 CHI Conference on Human Factors in Computing Systems*, CHI EA '19, 1–8. New York, NY, USA: Association for Computing Machinery. ISBN 9781450359719.
- Johnson, A. 2023. Spotify Removes 'Tens Of Thousands' Of AI-Generated Songs: Here's Why. *Forbes*.
- Johnson, K. 2020. MIT takes down 80 Million Tiny Images data set due to racist and offensive content. *VentureBeat*.
- Kahn, J.; Rivière, M.; Zheng, W.; Kharitonov, E.; Xu, Q.; Mazaré, P. E.; Karadayi, J.; Liptchinsky, V.; Collobert, R.; Fuegen, C.; Likhomanenko, T.; Synnaeve, G.; Joulin, A.; Mohamed, A.; and Dupoux, E. 2020. Libri-Light: A Benchmark for ASR with Limited or No Supervision. In *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 7669–7673. Barcelona, Spain: IEEE. <https://github.com/facebookresearch/libri-light>.
- Kelley, T.; and Dickerson, K. 2020. A review of artificial intelligence (AI) algorithms for sound classification: Implications for human-robot interaction (hri). *Defense Technical Information Center (DTIC)*.
- Kendall, T.; and Farrington, C. 2023. The Corpus of Regional African American Language. Version 2023.06. Eugene, Ore.: The Online Resources for African American Language Project.
- Khosrowi, D.; Finn, E.; and Clark, E. 2023. Diffusing the Creator: Attributing Credit for Generative AI Outputs. In *Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society*, 890–900. Montreal, Canada: Association for Computing Machinery (ACM).
- Kim, J.; Kim, S.; Kong, J.; and Yoon, S. 2020. Glow-tts: A generative flow for text-to-speech via monotonic alignment search. *Advances in Neural Information Processing Systems*, 33: 8067–8077.
- Kim, J. W.; Salamon, J.; Li, P.; and Bello, J. P. 2018. Crepe: A convolutional representation for pitch estimation. In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 161–165. IEEE, Calgary, Canada: IEEE.
- Kim, S.; Kim, H.; and Yoon, S. 2022. Guided-tts 2: A diffusion model for high-quality adaptive text-to-speech with untranscribed data. *arXiv preprint arXiv:2205.15370*.
- Koenecke, A.; Choi, A. S. G.; Mei, K. X.; Schellmann, H.; and Sloane, M. 2024. Careless Whisper: Speech-to-Text Hallucination Harms. In *The 2024 ACM Conference on Fairness, Accountability, and Transparency*, 1672–1681.
- Lapter, A. J. 2007. How the Other Half Lives (Revisited): Twenty Years Since Midler v. Ford. *Texas Intellectual Property Law Journal*.
- Law, E.; West, K.; Mandel, M. I.; Bay, M.; and Downie, J. S. 2009. Evaluation of algorithms using games: The case of music tagging. In *ISMIR*, 387–392. Citeseer, Kobe, Japan: International Society for Music Information Retrieval.
- Lee, H. P.; Yang, Y. J.; Von Davier, T. S.; Forlizzi, J.; and Das, S. 2024. Deepfakes, Phrenology, Surveillance, and More! A Taxonomy of AI Privacy Risks. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*, 1–19. Honolulu, USA: Association for Computing Machinery (ACM).

- Lee, K.; Hitt, G.; Terada, E.; and Lee, J. H. 2022. Ethics of Singing Voice Synthesis: Perceptions of Users and Developers. In *Proceedings of the 23rd International Society for Music Information Retrieval Conference (ISMIR 2022)*, 733–740. Bengaluru, India: International Society for Music Information Retrieval.
- Leschanowsky, A.; Rusti, C.; Quinlan, C.; Pnacek, M.; Gorce, L.; and Hutiri, W. 2024. A Data Perspective on Ethical Challenges in Voice Biometrics Research. *IEEE Transactions on Biometrics, Behavior, and Identity Science*.
- LibriVox. 2025. <https://librivox.org/pages/about-librivox/>. Accessed: 2025.
- Liu, H.; Chen, Z.; Yuan, Y.; Mei, X.; Liu, X.; Mandic, D.; Wang, W.; and Plumbley, M. D. 2023a. AudioLDM: Text-to-Audio Generation with Latent Diffusion Models. arXiv:2301.12503.
- Liu, Y.; Wu, P.; Black, A. W.; and Anumanchipalli, G. K. 2023b. A Fast and Accurate Pitch Estimation Algorithm Based on the Pseudo Wigner-Ville Distribution. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 1–5. IEEE, Rhodes Island, Greece: IEEE.
- Lo, C. 2024. Hong Kong worker transfers HK\$4 million after call from deepfake UK firm ‘CFO’.
- Longpre, S.; Mahari, R.; Lee, A.; Lund, C.; Oderinwale, H.; Brannon, W.; Saxena, N.; Obeng-Marnu, N.; South, T.; Hunter, C.; et al. 2024. Consent in Crisis: The Rapid Decline of the AI Data Commons. *arXiv preprint arXiv:2407.14933*.
- Lu, K.; Mardziel, P.; Wu, F.; Amancharla, P.; and Datta, A. 2020. *Gender Bias in Neural Natural Language Processing*, 189–202. Cham: Springer International Publishing. ISBN 978-3-030-62077-6.
- Luccioni, A. S.; Viguier, S.; and Ligozat, A.-L. 2023. Estimating the carbon footprint of bloom, a 176b parameter language model. *Journal of Machine Learning Research*, 24(253): 1–15.
- Luccioni, S.; Jernite, Y.; and Strubell, E. 2024. Power hungry processing: Watts driving the cost of AI deployment? In *The 2024 ACM Conference on Fairness, Accountability, and Transparency*, 85–99. Rio de Janeiro, Brazil: Association for Computing Machinery (ACM).
- Lucy, L.; Gururangan, S.; Soldaini, L.; Strubell, E.; Bamman, D.; Klein, L. F.; and Dodge, J. 2024. AboutMe: Using Self-Descriptions in Webpages to Document the Effects of English Pretraining Data Filters. arXiv:2401.06408.
- Mahelona, K.; Leoni, G.; Duncan, S.; and Thompson, M. 2023a. OpenAI’s Whisper is another case study in Colonisation.
- Mahelona, K.; Leoni, G.; Duncan, S.; and Thompson, M. 2023b. OpenAI’s Whisper is another case study in Colonisation.
- Manco, I.; Weck, B.; Doh, S.; Won, M.; Zhang, Y.; Bogdanov, D.; Wu, Y.; Chen, K.; Tovstogan, P.; Benetos, E.; Quinton, E.; Fazekas, G.; and Nam, J. 2023. The Song Descriptor Dataset: a Corpus of Audio Captions for Music-and-Language Evaluation. In *Machine Learning for Audio Workshop at NeurIPS 2023*.
- Marx, P. 2024. How artists are fighting generative AI.
- Mehrish, A.; Majumder, N.; Bharadwaj, R.; Mihalcea, R.; and Poria, S. 2023. A review of deep learning techniques for speech processing. *Information Fusion*, 99: 101869.
- Mei, X.; Meng, C.; Liu, H.; Kong, Q.; Ko, T.; Zhao, C.; Plumbley, M. D.; Zou, Y.; and Wang, W. 2024. Wavcaps: A chatgpt-assisted weakly-labelled audio captioning dataset for audio-language multimodal research. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*.
- Milmo, D. 2024. UK engineering firm Arup falls victim to £20m deepfake scam. *The Guardian*.
- Morreale, F.; Sharma, M.; and Wei, I.-C. 2023. Data Collection in Music Generation Training Sets: A Critical Analysis. In *ISMIR 2023 Hybrid Conference*, 37–46. Milan, Italy: International Society for Music Information Retrieval.
- Mulvey, L. 2013. Visual pleasure and narrative cinema. In *Feminism and film theory*, 57–68. London, UK: Routledge.
- Nacimiento-García, E.; Díaz-Kaas-Nielsen, H. S.; and González-González, C. S. 2024. Gender and Accent Biases in AI-Based Tools for Spanish: A Comparative Study between Alexa and Whisper. *Applied Sciences*, 14(11): 4734.
- Nest, T. E. 2015. Terms of Service.
- Newton-Rex, E. 2024a. Suno is a music AI company aiming to generate \$120 billion per year. But is it trained on copyrighted recordings?
- Newton-Rex, E. 2024b. Yes... Udio’s output resembles copyrighted music, too.
- Nicolas, G.; and Skinner, A. L. 2012. “That’s So Gay!” Priming the General Negative Usage of the Word Gay Increases Implicit Anti-Gay Bias. *The Journal of social psychology*, 152(5): 654–658.
- Nogueira, A. F. R.; Oliveira, H. S.; Machado, J. J.; and Tavares, J. M. R. 2022. Sound classification and processing of urban environments: A systematic literature review. *Sensors*, 22(22): 8608.
- of California, T. U. 2024. The public domain.
- Ojewale, V.; Steed, R.; Vecchione, B.; Birhane, A.; and Raji, I. D. 2024. Towards AI Accountability Infrastructure: Gaps and Opportunities in AI Audit Tooling. *arXiv preprint arXiv:2402.17861*.
- O’Reilly, P.; Jin, Z.; Su, J.; and Pardo, B. 2024. Maskmark: Robust Neural watermarking for Real and Synthetic Speech. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 4650–4654. IEEE, Seoul, South Korea: IEEE.
- Palaniappan, R.; Sundaraj, K.; and Sundaraj, S. 2014. Artificial intelligence techniques used in respiratory sound analysis—a systematic review. *Biomedizinische Technik/Biomedical Engineering*, 59(1): 7–18.
- Panayotov, V.; Chen, G.; Povey, D.; and Khudanpur, S. 2015. Librispeech: an ASR corpus based on public domain audio

- books. In *2015 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, 5206–5210. IEEE, Brisbane, Australia: IEEE.
- Papakyriakopoulos, O.; Choi, A. S. G.; Thong, W.; Zhao, D.; Andrews, J.; Bourke, R.; Xiang, A.; and Koenecke, A. 2023. Augmented datasheets for speech datasets and ethical decision-making. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*, 881–904.
- Patel, N. 2023. Google and YouTube are trying to have it both ways with AI and copyright - The Verge.
- Paullada, A.; Raji, I. D.; Bender, E. M.; Denton, E.; and Hanna, A. 2021. Data and its (dis) contents: A survey of dataset development and use in machine learning research. *Patterns*, 2(11).
- Peracha, O. 2022. JS Fake Chorales: a Synthetic Dataset of Polyphonic Music with Human Annotation. arXiv:2107.10388.
- Prabhu, V. U.; and Birhane, A. 2021. Large datasets: A pyrrhic win for computer vision. In *Institute of Electrical and Electronics Engineers/Computer Vision Foundation Conference on Applications of Computer Vision*. Virtual: IEEE.
- Queerina, O. O.; Ovalle, A.; Subramonian, A.; Singh, A.; Voelcker, C.; Sutherland, D. J.; Locatelli, D.; Breznik, E.; Klubicka, E.; Yuan, H.; et al. 2023. Queer in AI: a case study in community-led participatory AI. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*, 1882–1895. Chicago, USA: Association for Computing Machinery (ACM).
- Radford, A.; Kim, J. W.; Xu, T.; Brockman, G.; McLeavey, C.; and Sutskever, I. 2023. Robust speech recognition via large-scale weak supervision. In *International Conference on Machine Learning*, 28492–28518. PMLR, Honolulu, USA: PMLR.
- Raffel, C. 2016. *Learning-based methods for comparing sequences, with applications to audio-to-midi alignment and matching*. Columbia University.
- Raffel, C.; Shazeer, N.; Roberts, A.; Lee, K.; Narang, S.; Matena, M.; Zhou, Y.; Li, W.; and Liu, P. J. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140): 1–67.
- Rogge, A.; Anter, L.; Kunze, D.; Pomsel, K.; and Willenbrock, G. 2024. Standardized sampling for systematic literature reviews (STAMP method): ensuring reproducibility and replicability. *Media and Communication*, 12.
- Rudinger, R.; May, C.; and Van Durme, B. 2017. Social bias in elicited natural language inferences. In *Proceedings of the First ACL Workshop on Ethics in Natural Language Processing*, 74–79. Valencia, Spain: Association for Computational Linguistics (ACL).
- Salazar, M. E. 2024. Nurture Originals, Foster Art, and Keep Entertainment Safe Act of 2024. Draft bill, 118th Congress, 2^d Session, U.S. House of Representatives. Introduced 4 September 2024.
- Sap, M.; Card, D.; Gabriel, S.; Choi, Y.; and Smith, N. A. 2019. The risk of racial bias in hate speech detection. In *Proceedings of the 57th annual meeting of the association for computational linguistics*, 1668–1678. Florence, Italy: Association for Computational Linguistics (ACL).
- Schuhmann, C.; Beaumont, R.; Vencu, R.; Gordon, C.; Wightman, R.; Cherti, M.; Coombes, T.; Katta, A.; Mullis, C.; Wortsman, M.; et al. 2022. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in Neural Information Processing Systems*, 35: 25278–25294.
- Shelby, R.; Rismani, S.; Henne, K.; Moon, A.; Rostamzadeh, N.; Nicholas, P.; Yilla-Akbari, N.; Gallegos, J.; Smart, A.; Garcia, E.; 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*, 723–741.
- Shuyo, N. 2014. langdetect.
- Sigurgeirsson, A.; and Ungless, E. L. 2024. Just Because We Camp, Doesn't Mean We Should: The Ethics of Modelling Queer Voices. *arXiv preprint arXiv:2406.07504*.
- Sisario, B. 2024. Universal Music Group Pulls Songs From TikTok. *The New York Times*.
- Sisman, B.; Yamagishi, J.; King, S.; and Li, H. 2020. An overview of voice conversion and its challenges: From statistical modeling to deep learning. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 29: 132–157.
- Snyder, D.; Chen, G.; and Povey, D. 2015. MUSAN: A Music, Speech, and Noise Corpus. ArXiv:1510.08484v1, arXiv:1510.08484.
- Soldaini, L.; Kinney, R.; Bhagia, A.; Schwenk, D.; Atkinson, D.; Authur, R.; Bogin, B.; Chandu, K.; Dumas, J.; Elazar, Y.; Hofmann, V.; Jha, A. H.; Kumar, S.; Lucy, L.; Lyu, X.; Lambert, N.; Magnusson, I.; Morrison, J.; Muennighoff, N.; Naik, A.; Nam, C.; Peters, M. E.; Ravichander, A.; Richardson, K.; Shen, Z.; Strubell, E.; Subramani, N.; Tafjord, O.; Walsh, P.; Zettlemoyer, L.; Smith, N. A.; Hajishirzi, H.; Beltagy, I.; Groeneveld, D.; Dodge, J.; and Lo, K. 2024. Dolma: an Open Corpus of Three Trillion Tokens for Language Model Pretraining Research. arXiv:2402.00159.
- Soni, A. 2024. Actors Say AI Voice-Over Generator Eleven-Labs Cloned Likenesses.
- Sturm, B. L.; Ben-Tal, O.; Monaghan, Ú.; Collins, N.; Herremans, D.; Chew, E.; Hadjeres, G.; Deruty, E.; and Pachet, F. 2019. Machine learning research that matters for music creation: A case study. *Journal of New Music Research*, 48(1): 36–55.
- Tang, C.; Yu, W.; Sun, G.; Chen, X.; Tan, T.; Li, W.; Lu, L.; MA, Z.; and Zhang, C. 2024. SALMONN: Towards Generic Hearing Abilities for Large Language Models. In *The Twelfth International Conference on Learning Representations*.
- Team, G. D. 2019. Celebrating Johann Sebastian Bach Doodle - Google Doodles.
- Thiel, D. 2023. Identifying and eliminating csam in generative ml training data and models. Technical report, Technical report, Stanford University, Palo Alto, CA, 2023. URL <https://purl...>

- Upadhyay, P.; Heung, S.; Azenkot, S.; and Brewer, R. N. 2023. Studying Exploration & Long-Term Use of Voice Assistants by Older Adults. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, CHI '23. New York, NY, USA: Association for Computing Machinery. ISBN 9781450394215.
- Veaux, C.; Yamagishi, J.; and MacDonald, K. 2017. CSTR VCTK Corpus: English Multi-speaker Corpus for CSTR Voice Cloning Toolkit.
- Visible, . P. 2025. Terms of Service.
- Wang, C.; Riviere, M.; Lee, A.; Wu, A.; Talnikar, C.; Haziza, D.; Williamson, M.; Pino, J.; and Dupoux, E. 2021. VoxPopuli: A Large-Scale Multilingual Speech Corpus for Representation Learning, Semi-Supervised Learning and Interpretation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, 993–1003. Online: Association for Computational Linguistics.
- Wang, R.; Juefei-Xu, F.; Huang, Y.; Guo, Q.; Xie, X.; Ma, L.; and Liu, Y. 2020. Deepsonar: Towards effective and robust detection of ai-synthesized fake voices. In *Proceedings of the 28th ACM international conference on multimedia*, 1207–1216.
- Wang, Y.; Ju, Z.; Tan, X.; He, L.; Wu, Z.; Bian, J.; and Zhao, S. 2023. AUDIT: Audio Editing by Following Instructions with Latent Diffusion Models. arXiv:2304.00830.
- Wang, Y.; Wei-Kocsis, J.; Springer, J. A.; and Matson, E. T. 2022. Deep learning in audio classification. In *International Conference on Information and Software Technologies*, 64–77. Springer, Kaunas, Lithuania: Springer.
- Willison, S. 2022a. Exploring 10m scraped Shutterstock videos used to train Meta's Make-A-Video text-to-video model.
- Willison, S. 2022b. Exploring the training data behind Stable Diffusion.
- Willison, S. 2023a. Exploring MusicCaps, the evaluation data released to accompany Google's MusicLM text-to-music model.
- Willison, S. 2023b. What's in the RedPajama-Data-1T LLM training set.
- Wu, S.; Li, X.; Yu, F.; and Sun, M. 2023. TunesFormer: Forming Irish Tunes with Control Codes by Bar Patching. arXiv:2301.02884.
- Wu, S.-L.; Donahue, C.; Watanabe, S.; and Bryan, N. J. 2024. Music ControlNet: Multiple Time-Varying Controls for Music Generation. *IEEE/ACM Trans. Audio, Speech and Lang. Proc.*, 32: 2692–2703.
- Yamagishi, J.; Veaux, C.; MacDonald, K.; et al. 2019. Cstr vctk corpus: English multi-speaker corpus for cstr voice cloning toolkit (version 0.92). *University of Edinburgh. The Centre for Speech Technology Research (CSTR)*, 271–350.
- YouTube. 2024a. Creative Commons.
- YouTube. 2024b. Terms of Service.
- Yu, J. E.; Parde, N.; and Chattopadhyay, D. 2023. “Where is history”: Toward Designing a Voice Assistant to help Older Adults locate Interface Features quickly. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, CHI '23. New York, NY, USA: Association for Computing Machinery. ISBN 9781450394215.
- Zhang, Y.; Park, D. S.; Han, W.; Qin, J.; Gulati, A.; Shor, J.; Jansen, A.; Xu, Y.; Huang, Y.; Wang, S.; et al. 2022. Bigssl: Exploring the frontier of large-scale semi-supervised learning for automatic speech recognition. *IEEE Journal of Selected Topics in Signal Processing*, 16(6): 1519–1532.
- Zhuo, L.; Yuan, R.; Pan, J.; Ma, Y.; Li, Y.; Zhang, G.; Liu, S.; Dannenberg, R.; Fu, J.; Lin, C.; Benetos, E.; Xue, W.; and Guo, Y. 2024. LyricWhiz: Robust Multilingual Zero-shot Lyrics Transcription by Whispering to ChatGPT. arXiv:2306.17103.

A Appendix

A.1 1: Datasheets for audio datasets

Our audit uncovered a lack of documentation of audio datasets. This lack of documentation complicates assessing issues of bias and representation in these datasets, in addition to other considerations. Gebru et al. (2021) propose datasheets for datasets. Below we propose datasheets for audio datasets. We reproduce the documentation guide in (Gebru et al. 2021), and we add several questions and modifications specific to audio datasets (our modifications and additions in italics). While our paper answers some of these questions for some popular audio datasets, in general audio datasets have a high documentation debt (Bender et al. 2021). We hope future work will answer more of these questions for more audio datasets.

(1) Motivation

- (a) What is the purpose of this dataset?
- (b) Who is the funding source of this dataset?
- (c) Who is the creator of this dataset?
- (d) *What type of task was this dataset primarily composed for—e.g., generative modelling? Classification? Audio separation? Style transfer? TTS? Something else?*
- (e) *What is the intended domain of the output of this dataset—speech? Music? General audio non-discriminatory sounds? Intentional generalization of both speech and music? Non-music, non-speech sounds such as sound effects?*
- (f) *Was this dataset intended to be exclusively in the audio domain, or multi-modal?*
- (g) *Was this dataset created for a specific purpose, or is it intended to be general use?*
- (h) *Was this dataset recorded with the intent of being “pure” sound as in separated or clear audio streams, or was it intentionally created with some environmental or background noise present?*
- (i) *Was this dataset created to fill a vacancy in our existing data (such as a low-resource language), or was it to add to an existing pool of similar datasets?*
- (j) *Was this a participatory effort by the community to build this dataset, or was this created by someone who is tangential to or otherwise does not belong to the community of which the dataset comprises?*
- (k) *If separate, what efforts were made to connect with the community for whom this dataset was comprised?*
- (l) Any other comments?

(2) Composition

- (a) What do the instances that comprise this dataset represent (e.g., full songs, 10-second song clips, speech excerpts, singular utterances, 1-5 second sound effects)? Are there multiple types of instances?
- (b) How many instances are there total?
- (c) *What are the descriptive statistics of the time length of these instances in seconds/minutes/hours (e.g., mean, median, sum, mode)?*

- (d) *What data does each instance consist of—raw audio files? Metadata about the audio? Midi files? Other forms of symbolic representations?*
- (e) *If the dataset contains music:*
 - i. *What are the genres present? Is it an equal distribution?*
 - ii. *Are they pre-existing/pre-released songs, or were they recorded for the purpose of this dataset?*
 - iii. *Is the data created from computational generations of symbolic music?*
 - iv. *Is the data composed of original recordings of music?*
- (f) *If the dataset contains speech:*
 - i. *What languages are present? Is it an equal distribution?*
 - ii. *How many speakers are present in the dataset?*
 - iii. *Do the speakers speak multiple language or code-switch, and how is this identified?*
 - iv. *What topics are present in the dataset?*
 - v. *Are there a variety of emotions present in the dataset? Is this identified in any way?*
- (g) Is there a label or target associated with each instance?
- (h) Is any desired metadata missing from individual instances?
 - (i) Are there recommended data splits present (e.g., training/validation/test)?
 - (j) What sources of noise are present in this dataset, and are they intentional or unintentional?
 - (k) Are there any duplications present in this dataset (e.g., songs repeated multiple times), and if so, why?
 - (l) Is the dataset self-contained, or does it rely on/link to an external source (e.g., links to song recordings)? If so, why?
- (m) *Does the dataset contain data that may be copyrighted?*
- (n) Does the dataset contain anything that may be considered offensive?
- (o) *Has all of the data been listened to by an author or member of the data-construction team?*
- (p) *If the dataset contains sounds created by humans:*
 - i. *What demographic information about the speakers can you provide? Ex. age, gender, sexual orientation, language, locale/country, and accent.*
 - ii. Does the dataset contain audio files that could feasibly be used to identify human beings?
 - iii. Does the dataset contain any speech that might be considered confidential or sensitive in any way (e.g., race or ethnic origins, sexual orientations, religious beliefs, political opinions, financial data, health data, biometrics data, or private and secure information such as criminal records or SSNs)?
 - iv. *Does the dataset contain speech from a population that does not belong to the group they are trying to emulate (e.g., a white American native attempting to mimic an accent from a country to which they have no genuine relation)?*

(q) Any other comments?

(3) Collection Process

- (a) How was the data associated with each instance acquired?
- (b) Was it recorded explicitly for the purpose of this dataset, or sourced elsewhere?
- (c) *Was the audio scraped, recorded, computationally created (e.g., from MIDI files), or created in another manner?*
- (d) *Was the dataset sourced from audio in the public domain?*
- (e) *What technical setting was the speech recorded from, including the recording environment? What tools were used to process the audio?*
- (f) *What is the sampling rate of the audio? How many channels?*
- (g) What mechanisms or procedures were used to collect the data?
- (h) *Was the data created from a computational algorithm?*
 - (i) *Who was involved in the audio recordings, and how were they compensated?*
 - (j) *Were there any copyright agreements struck with the contributors to the dataset?*
 - (k) *What were the contributors to the dataset told their contribution would entail? Were they thoroughly briefed on the potential extent of the use of their audio, or was it left to general terms?*
- (l) Over what timeframe was the audio collected?
- (m) *Was the audio collected from any multimodal sources (e.g., stripped from video)?*
- (n) Were any ethical review processes conducted (e.g., from an IRB)?
- (o) *Has an analysis of the potential impact of the dataset been conducted?*
- (p) If the dataset contains sounds created by humans:
 - i. Were the individuals present in the dataset notified about the data collection?
 - ii. Did the individuals present in the dataset consent to the data collection and use of their data?
 - iii. If consent was obtained, were the individuals provided with a mechanism to revoke their consent at a later date?
- (q) Any other comments?

(4) Preprocessing/Cleaning/Labeling

- (a) What preprocessing of the data was conducted in order to get it at the final stage the audio is in now?
- (b) *If there is metadata present, how was the metadata sourced?*
- (c) *Did users consent to release of this metadata?*
- (d) *If there is demographic information present about the individuals who recorded the audio, how was that obtained? Self-reported? Scraped? Sourced elsewhere?*

- (e) Was the original raw data saved in addition to the final version of the data (if different)?
- (f) *Was background noise (or otherwise present environmental noise) intentionally cleaned from this data?*
- (g) Is the code or other software used to prepare the data available to be published/released/acknowledged in line with the dataset?
- (h) *Are there any transcriptions available of the speech/lyrics? How was this processed? If human annotators, how were the annotators trained to create the transcriptions? If computationally annotated, what was this process? What is the accuracy rate?*
- (i) *Were there any content tagging such as hate-speech tags or swear word flagging?*
- (j) *Were any redactions of sensitive data performed on this dataset? How was this conducted?*
- (k) Any other comments?

(5) Uses

- (a) Has the dataset been used for any tasks already?
- (b) *If so, have the original voices/musicians in the dataset been made aware of the extended/full use of the dataset?*
- (c) *Is there any repository linking all papers or systems using (or with access to) the dataset?*
- (d) What tasks is the dataset designed for? Which are they suitable for? What could it potentially be used for?
- (e) Are there tasks for which the dataset should explicitly not be used?
- (f) Is there any element of the composition of the dataset that may impact how future uses of the dataset may be impacted?
- (g) *Is there any part of this dataset that is privately held but can be requested for research purposes?*
- (h) Any other comments?

(6) Distribution

- (a) Will the dataset be distributed to third parties outside of the entity (e.g., company, institutions, organization) on behalf of which the dataset was created?
- (b) *Are the creators of audio files aware of this?*
- (c) How will the dataset be distributed (e.g., open source, published on GitHub, by email request)?
- (d) When will the dataset be distributed?
- (e) Will the dataset be distributed under any copyright or other IP protection?

(7) Maintenance

- (a) Will the dataset be maintained or otherwise updated after the time of initial release? By whom?
- (b) How will the owner/curator of the dataset be contacted?
- (c) *If an artist discovers their work is present in this dataset, can they request it be removed? How?*
- (d) Will the dataset be updated (e.g., to correct any potential errors, adding new files)?

- (e) How long will the data be retained?
- (f) If others want to contribute to this data source, is there a mechanism for them to do so?

A.2 2: Literature Review Methodology

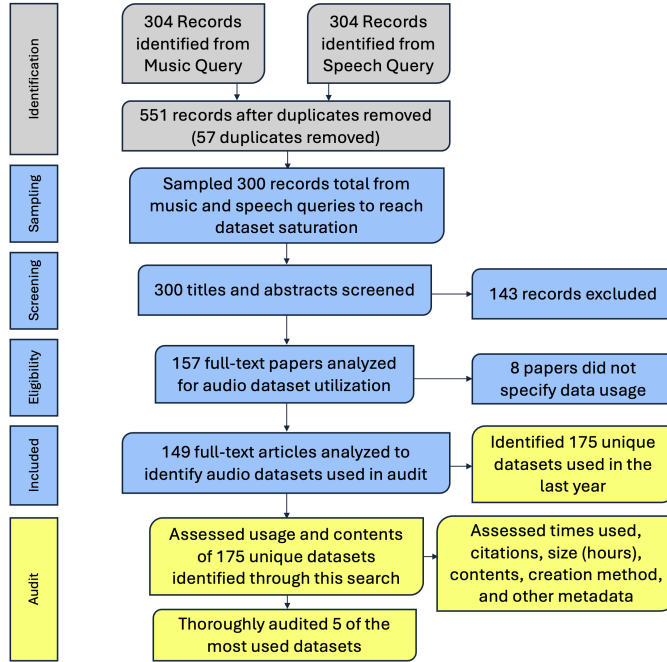


Figure 5: Flow diagram detailing the paper corpus used in the mapping review to produce the datasets audited in this paper. We started with 551 records from arXiv (grey) and analyzed 149 full-text articles (blue) to identify 175 unique datasets for analysis (yellow).

We were interested in identifying audio datasets currently used by researchers to both include in our audit and provide an overview of the data landscape. We chose to conduct a mapping literature review with systematic elements (Ferrari 2015) utilizing the STAMP sampling method (Rogge et al. 2024) of one year of arXiv (arXiv 2023) computer science works about audio models. We chose one year since there has been a paradigmatic shift in audio generative models, driven in part by the recent advances in large language models and both their translation to text-to-audio models and adopted language-model-style generation as seen in AudioLM (Borsos et al. 2023a), MusicLM (Agostinelli et al. 2023), SoundStorm (Borsos et al. 2023b), VampNet (Garcia et al. 2023), and more—researchers and commercial developers alike are largely abandoning the early approaches we saw in 2020-2022 such as the music transformer (Huang et al. 2018). Since 2023, approaches towards building Large Audio Language Models (LALMs) largely involve combining many datasets from a broad range of tasks in speech, audio, and music processing (Tang et al. 2024; Gong et al. 2023b; Chu et al. 2023), often using proprietary language models like the

GPT model family to supplement existing data with additional question-answer pairs (Ghosh et al. 2024; Deshmukh et al. 2024). One such dataset, OpenASQA, combines a total of 13 publicly available audio, music, and speech datasets to train their LALM, LTU-AS, and uses GPT-3.5-Turbo to generate QA pairs (Gong et al. 2023a). As researchers move towards curating giant meta-datasets of datasets, it becomes exceedingly vital to understand the origins, licensing, and limits of the many datasets that feed the creation of LALMs.

We chose to focus on arXiv submissions because Barnett’s (Barnett 2023) recent systematic literature review on the ethical implications of generative audio models resulted in a final corpus comprised of 91% arXiv works even after starting from a 50/50 split of arXiv and ACM works. Though it is difficult to quantify the most influential papers in anything other than citations, which is certainly not a perfect metric, we verified that all the major audio generation papers such as the ones already mentioned in this paper were present in arXiv as well to justify this focus of our exploratory systematic literature review.

Following Barnett’s literature review (Barnett 2023), we used the following query on arXiv (once for music, once for speech) for our mapping review. This query searches all fields on arXiv including title, abstract, and keywords of papers in this database, but does not search the full text of articles.

```

1 [
2   [
3     [All: "generative"] AND
4     [All: "$\langle$music/speech$\rangle$"] AND
5     [All: "model"]
6   ]
7   AND [date_range: from 2023-05-01 to
8     2024-05-01]
9 ]; classification: Computer Science (cs)

```

The music query produced 304 records and the speech query produced 1561. For context, the same date range from 2022-2023 included 114 records for music and 421 for speech, further confirming the rapid growth of this field. In order to have parity, we took the most recent 304 queries from speech. We then removed duplicates and sampled 300 records (150 each) for further analysis.

Within these 300 records, we then screened the titles and abstracts to assess eligibility as follows. We restricted papers to full-length audio papers using audio data for analysis or training. Of the 143 records we excluded, most were entirely focused on another modality such as text or vision (72%), conducted literature reviews or meta-analyses (10%), researched dance or motion generation (10%), were incomplete—e.g., stopped mid paragraph and was not a finished paper (3%), or had another salient quality (5%) such as not being in English or withdrawn from arXiv (which typically only occurs when a co-author did not grant permission for uploading to arXiv or something was substantially incorrect with the paper, withdrawals are not easily granted by request of authors). Of the remaining 157 papers in the corpus, we further excluded eight papers that did not specify their data usage even

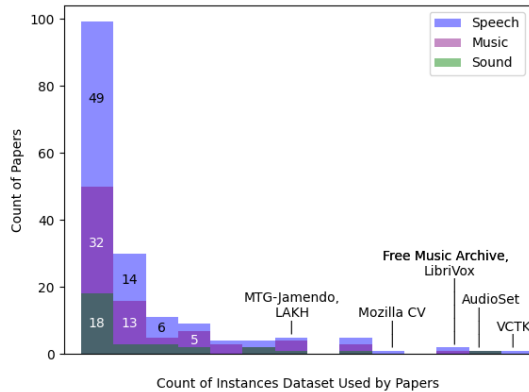


Figure 6: Stacked bar plot displaying count of times datasets were used by papers in the corpus. Split by Speech, Music, and Non-music/Non-Speech sounds. The vast majority of datasets were only used once, while a select few were used multiple times.

though they clearly used a dataset for their paper. This resulted in 149 full-text articles for inclusion. If the datasets were simply noted (e.g., in the related literature section), but not explicitly discussed as being part of the training or evaluation process described in the paper, they were not included in this analysis. This process yielded 175 unique audio datasets.

Dataset Labeling and Categorization

Calculating Dataset Content Duration (Hours) When explicitly stated or when we were provided with number of files and average file length, we listed the exact duration of the dataset. When we were provided with number of files and access to individual files, we got the average file length based on a small random sample and multiplied by number of files in the dataset. Otherwise, we made the following assumptions unless explicitly instructed not to do so: (1) average song length: 3 minutes; (2) average children’s song length: 1.5 minutes; (3) single sentence utterance: 12 words/5 seconds; (4) single word utterance: 1 second; (5) phonetically dense sentences: 15 seconds; (6) YouTube content: 10 or 30 seconds based on data examination.

Original Content (Yes/No): A dataset is considered original content if the creator(s) featured new recordings, synthesized data from an existing source, created new data from a model trained on an existing source (such as new chorales from JS Fake Chorales (Peracha 2022)), converted an existing source from one modality to another (e.g., WSJ0 (Garofolo et al. 1993), which includes speech recordings made from its sister WSJ text corpus), or added significant data derived from an existing source, such as crowd-sourced annotations to tag songs (e.g., MagnaTagATune (Law et al. 2009)). They are not regarded as original content if the dataset solely consists of scraped or crawled videos or their links to another source.

Potential for Copyright Infringement (Yes/No): We adopt a conservative and stringent approach to assess-

ing whether a dataset has the potential for copyright infringement. For example, any datasets scraped from YouTube are classified as having potential for copyright infringement as YouTube hosts a wide variety of content, much of which is protected by copyright. Scraping and distributing this content without proper authorization violates YouTube’s terms of service and copyright laws. For annotation-based datasets that provide links to audio recordings, we assess the original sources from which users can obtain the respective audio recordings. As such, we categorize MTG-Jamendo (Bogdanov et al. 2019) (and its derivative, Song Descriptor Dataset (Manco et al. 2023)) and Free Music Archive (Defferrard et al. 2017) as not infringing upon copyright because the sources, Jamendo.com and freemusicarchive.org, offer music within the public domain. Conversely, we categorize Million Song Dataset (Bertin-Mahieux et al. 2011) and MagnaTagATune (Law et al. 2009) as having potential for copyright infringement as their audio sources, 7digital.com and magnatune.com, necessitate purchasing licenses for access beyond private listening.

Method of Dataset Creation (Scraped, Created, or Augmented): If any portion of the dataset is created through scraping (even if others are not), we categorize it as ‘scraped’. If the dataset is a direct subset of another, we label it as ‘augmented’. We denote all other datasets as ‘created’.

We built a database of meta-data for each of these datasets, available at [redacted for blind review]. It contains more granular information such as download links, original purpose of the dataset, whether it was free to access, and if applicable, the language or genre of the contents.

A.3 3: Transcription and Textual Analysis

We first obtain high-quality transcripts. While Common Voice, VCTK, LibriVox, Wav-Caps, GigaSpeech, and the Lakh MIDI datasets contain transcripts, we found the transcripts included with AudioSet Youtube videos were of lower quality and not well-aligned. Therefore, we use Whisper-large (Radford et al. 2023) to transcribe over 50,000 randomly selected AudioSet Youtube clips. We also use Whisper-large with prompting improvements for music transcription from Zhuo et al. (2024) to transcribe Jamendo and the Free Music Archive. After transcription, we use a range of established text dataset audit techniques and tools. For toxicity analysis, we use the `pySentimento` library used to detect hate content (Birhane et al. 2023), as well as Surge AI’s profanity list (AI 2022). To detect language, we use the `langdetect` library (Shuyo 2014) along with Whisper language predictions on AudioSet, Jamendo, and Free Music Archive clips. To investigate the dataset content in relation to sociodemographic identities, we search for a set of keywords encompassing race, gender, religion, and sexual orientation, using lists established from prior work (Dodge et al. 2021; Hong et al. 2024; Lu et al. 2020). Finally, we track common words and calculate pointwise mutual information between these and identity keywords to determine

stereotypical associations (Rudinger, May, and Van Durme 2017).

We recognize the various limitations of these audit libraries and techniques, especially as hate speech models and profanity detection are more likely to wrongly flag data relating to LGBTQ+ or Black voices as toxic content (Dodge et al. 2021; Sap et al. 2019). In addition, language predictions are unreliable and rely on text content (Lucy et al. 2024). Despite these challenges, our goal is to present findings from a broad initial audit of these datasets to identify concerns that may warrant further in-depth investigation in future work; we argue that our findings are still useful indicators of broad trends of bias or toxicity.

A.4 4: AudioSet Contents

Video titles As AudioSet contains audio from music, speech, and general sounds, we assess its contents by analyzing titles of included texts. We first remove stop words from titles, and then we extract keywords directly from the words that remain as well as common associations with those keywords (Figure 7). While YouTube contains a broad range of videos, we find that many sound snippets in AudioSet cluster around specific topics. Video games (“video game”, “guitar hero”), music (“live music”, “singing”, “music festival”), reviews (“review video”), pets (“dog”, “cat”), and vehicles (“car”, “boat”) compose a significant fraction of AudioSet.

Sample transcripts In Table 3, we show examples of toxic-detected transcripts and include whether Whisper categorizes the audio as speech or music.

In Table 4, out of the 26 clips in AudioSet that mention *gay*-related keywords, we show several samples in which this term is used as a derogatory slur and depicts homophobic stereotypes. This problematic usage may propagate to speech generation models trained on AudioSet and amplify social bias (Nicolas and Skinner 2012).

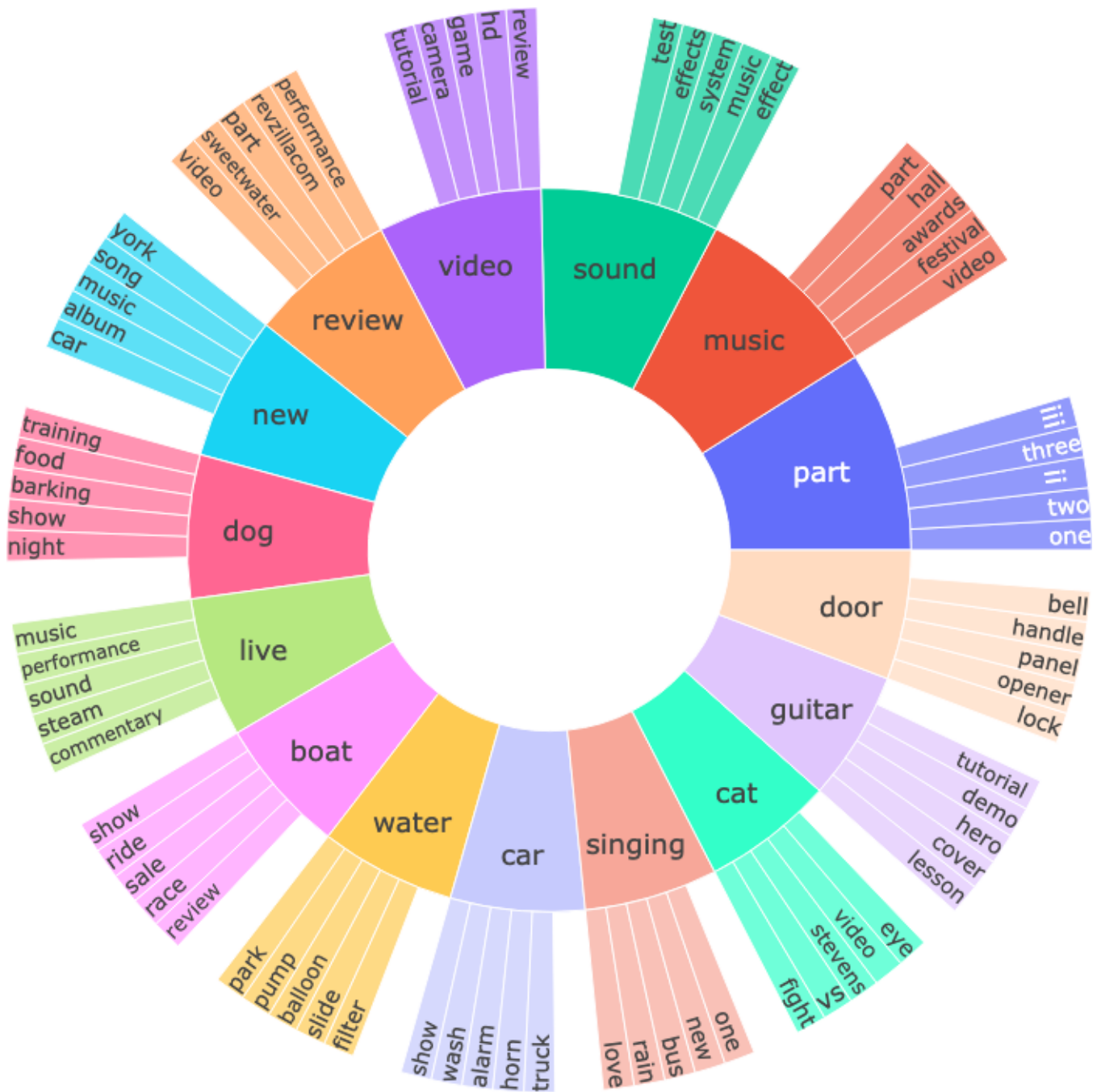


Figure 7: Common bigram breakdown of AudioSet video titles. The inner ring shows the top 25 most common words ignoring stopwords with sizes relatively to scale between words. The outer ring shows the top 5 words that follow each word in the inner ring with constant sizes for readability.

Table 3: Ten randomly sampled transcripts in AudioSet that are detected as hateful by `pysentiment0`. Profane words are redacted.

Category	Text
Speech	What? Alicia Fox pulling that hair. She still has.
Speech	and then shoot rocks up his *ss.
Speech	I don't know how you got here. I'm sorry. You're a real pain in the *ss. You're a real pain in the *ss. You're a real pain in the *ss.
Music	Pour some sh*t, matter of fact, go ahead and drink that Couple more shots, then we'll get freaky I peeped that, now I need that To the p*ssy like a record, go ahead and leak it It's real food, a dude up in the plate, party rockin'
Music	b*tch
Speech	do this without killing us.
Music	You're despicable.
Speech	Facebook it is. Hey look, I'm putting a video on Facebook. Me too. Look at her. Get the f*ck out of here, you d**chebag.
Speech	that woman had to the order and you you're an idiot no but i didn't come in here first because i didn't know the order
Speech	Stay with your p*ssy! Yes!

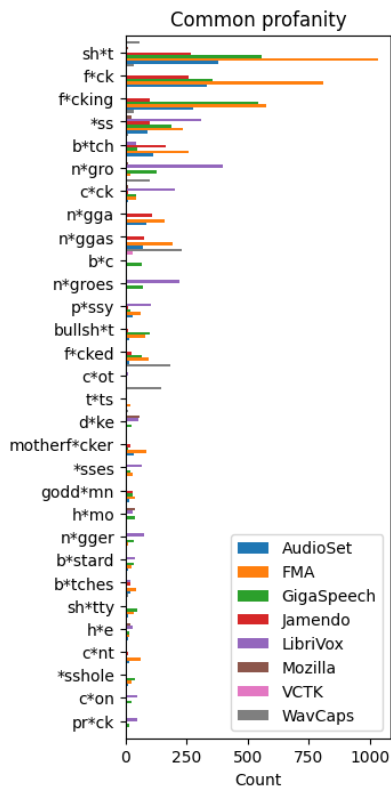


Figure 8: English profanity occurrences in all datasets.

Table 4: Sample transcripts in AudioSet that contain either `gay` or `gays` keywords. Profane words are redacted.

Text
You must be whipped, you must be gay. You must be whipped, whipped, you must be gay.
And how was the tea virus? F*ck that intro. It's gay.
Why do you suck your thumb? Whoa, whoa, f*ck you motherf*cker. I don't suck my f*cking thumb. He sucks his middle finger. Are you sure? That's gay. You're not gay for that.
Oh god, that's gay. Hold on. There we go, spooky! Isn't it f*cking cute? I think it's the cutest thing ever.
The zip code must have been like San Francisco as to gay people as to cats are to breeding. There was just ridiculous amounts of cats. And they would not shut up.
Why did you want this then? Lol, are you gay? Are you serious? Serious about what? What's f*pping?
Ha! Gay!

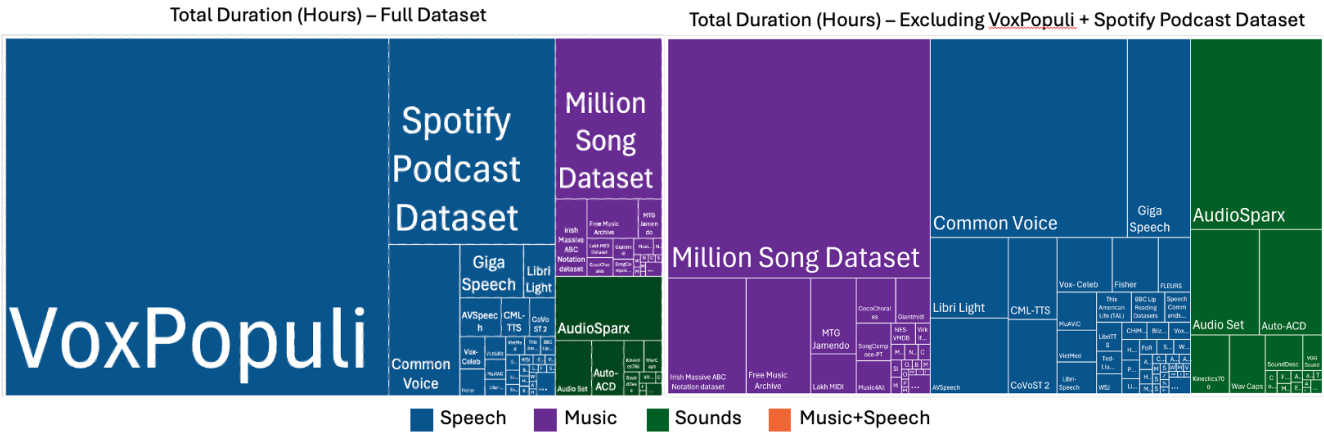


Figure 9: Area charts displaying the proportion of each dataset's estimated total duration relative to (1) all audited datasets and (2) all datasets excluding the two largest datasets, VoxPopuli and Spotify Podcast Dataset (right)

A.5 5: Binary Gender PMI

A.6 6: Additional Figures

A.7 7: Code and Data Availability

Code and annotation data are available at <https://anonymous.4open.science/r/gen-audio-ethics-F8CF/README.md>.

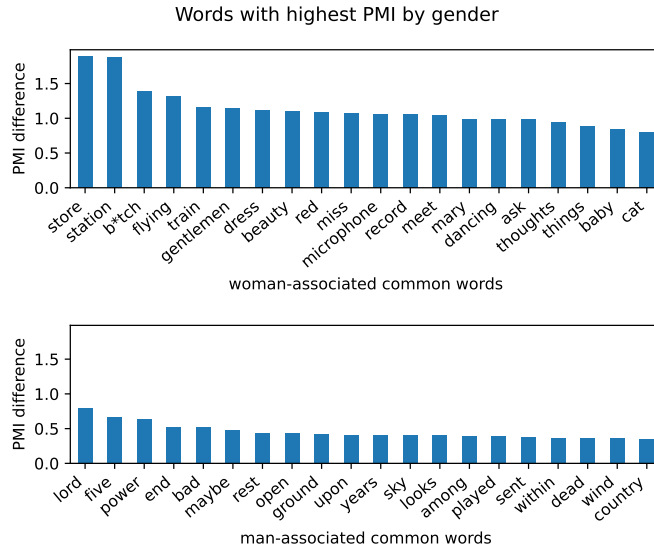


Figure 10: Common words with highest PMI to man-related versus woman-related keywords (from Lu et al. (2020)) across all datasets (randomly sampling 40 thousand sentences each). We consider words that appear at least 50 times in man-related and woman-related samples.

Language	Duration	Language	Duration
Catalan	3,587	German	923
Kinyarwanda	2,384	Norwegian	748
Spanish	2,219	Javanese	410
Esperanto	1,905	Russian	358
Belarusian	1,765	Vietnamese	336
German	1,424	Spanish	333
Bengali	1,273	Portuguese	322
French	1,147	Latin	312
Swahili	1,085	French	306
Chinese	1,061	Welsh	235

Table 5: The top 10 non-English languages by duration in hours for both the largest of these 9 datasets–Mozilla Common Voice, as well as for the other 8 datasets.

Dataset	License
Mozilla Common Voice	CC0 (Commons 2024b)
VCTK	Open Data Commons Attribution License (Commons 2024c)
LibriVox	CC BY 4.0 (Commons 2024a)
Lakh	Echo Nest License (Nest 2015)
AudioSet	Youtube License (YouTube 2024b) or Creative Common Licenses (YouTube 2024a)
Free Music Archive	Various Creative Commons
Jamendo	Various Creative Commons
Wav-Caps	Various Creative Commons, Youtube License (YouTube 2024b), BBC's Content Licence for RemArc (BBC 2025)
GigaSpeech	Many licenses, including 99% Invisible license (Visible 2025), Australian Broadcasting Company license(ABC 2025), YouTube License(YouTube 2024b)

Table 6: Licenses of audio datasets

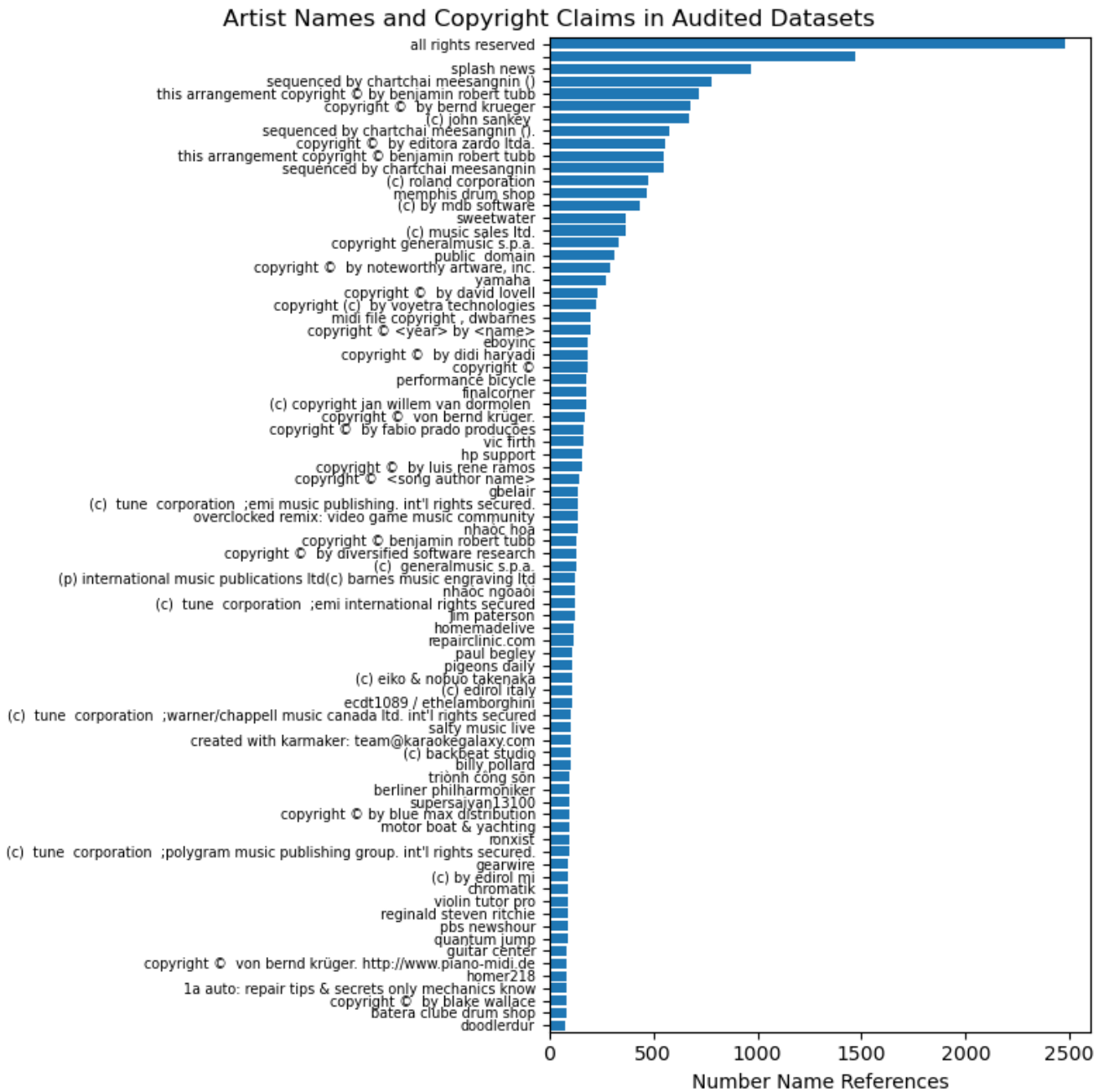


Figure 11: Count of appearances of artist names and copyright claims in the AudioSet and Lakh datasets.