

## 424 B Seed artists

425 We use the following artists as a seed list for the traversal of the artist graph on Spotify.

```
426 starting_artists = [  
427     '2ye2Wgw4gimLv2eAKyk1NB', # Metallica  
428     '3WrFJ7ztbogyGnTHbHJF12', # The Beatles  
429     '74ASZWbe4lXaubB36ztrGX', # Bob Dylan  
430     '1Xyo4u8uXC1ZmMpatF05PJ', # The Weeknd  
431     '4y6J8jwRAw04dssiSmN91R', # Muddy Waters  
432     '776Uo845nYHJpNaStv1Ds4', # Jimi Hendrix  
433     '6kACVPfC0nqzgfEF5ryl0x', # Johnny Cash  
434     '67ea9eGLXYMs02eYQRui3w', # The Who  
435     '7dGJo4pcD2V6oG8kP0tJRR', # Eminem  
436     '5pKCCKE2ajJHZ9KAiaK11H', # Rihanna  
437     '0dmPX6ovclg0y8WWJaFEUU', # Kraftwerk  
438     '1ZwdS5xdxEREPySFridCfh', # 2pac  
439     '0kbYTNQb4Pb1rPbbaF0pT4', # Miles Davis  
440     '6tbjWDEIzxoDsBA1FuhfPW', # Madonna  
441     '7guDJrEfX3qb6FEbdPA5qi', # Stevie Wonder  
442     '2QsynagSdAqZj3U9HgDzjD', # Bob Marley  
443     '5aIqB5nVVvmFsvSdExz408', # Johann Sebastian Bach  
444     '0Kekt6CKSo0m5mivKcoH51', # Sergei Rachmaninoff  
445 ]
```

## 446 C Columns in DISCO-10M Dataset

```
447 dataset_columns = [  
448     'video_url_youtube',  
449     'video_title_youtube',  
450     'track_name_spotify',  
451     'video_duration_youtube_sec',  
452     'preview_url_spotify',  
453     'video_view_count_youtube',  
454     'video_thumbnail_url_youtube',  
455     'search_query_youtube',  
456     'video_description_youtube',  
457     'track_id_spotify',  
458     'album_id_spotify',  
459     'artist_id_spotify',  
460     'track_duration_spotify_ms',  
461     'primary_artist_name_spotify',  
462     'track_release_date_spotify',  
463     'explicit_content_spotify',  
464     'similarity_duration',  
465     'similarity_query_video_title',  
466     'similarity_query_description',  
467     'similarity_audio',  
468     'audio_embedding_spotify',  
469     'audio_embedding_youtube',  
470 ]
```

471 **D Additional Examples: Audio Similarity Comparison**

472 We demonstrate the results of our audio similarity approach on five additional samples (see Figure 6).  
473 Similarly to the spectrograms presented in Section 3.2, we also observe an overlap for these samples  
474 when the audio similarity is above  $\delta_a > 0.4$ . Even when two music snippets have the same frequency  
475 characteristics, there might still be small differences. This can be explained in part due to the audio  
476 quality of a YouTube video, which is dependent on the quality selected by the person uploading the  
477 video, and can therefore vary greatly, unlike the audio quality of Spotify. This difference can be seen  
478 best in the high-frequency content of the spectrogram, which tends to be weaker and less pronounced  
479 in the YouTube audio samples. We notice a strong dissimilarity in the first example between the  
480 Spotify preview audio spectrogram and the YouTube audio spectrogram. This is reflected by the low  
481 similarity score of  $\delta_a = 0.1985$ .

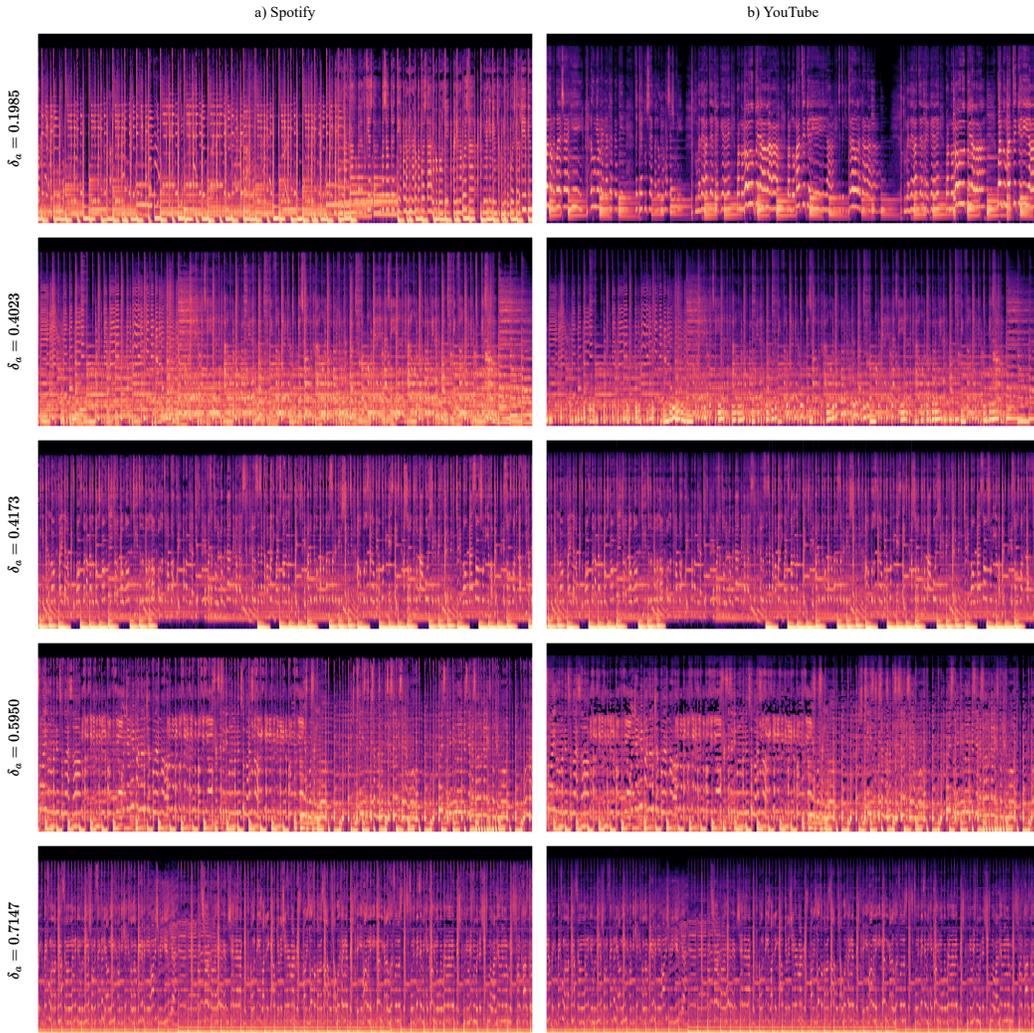


Figure 6: Comparison of audio similarity between Spotify preview audio and YouTube audio.  $\delta_a$  denotes the cosine similarity of the audio embedding. We observe that the similarity of our audio embeddings is related to the similarity of the Log-Mel Spectrograms, and that the similarity increases when the spectrograms are closer to each other.

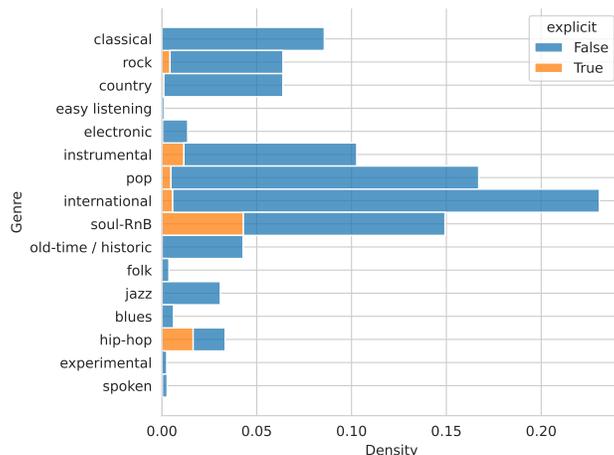


Figure 7: Genre distribution of FMA root genres in our dataset. The genre embeddings were computed using a CLAP encoder on the sentence *This audio is a <genre> song*. Each song is mapped to a genre according to the largest cosine similarity between the song embedding and the genre embedding.

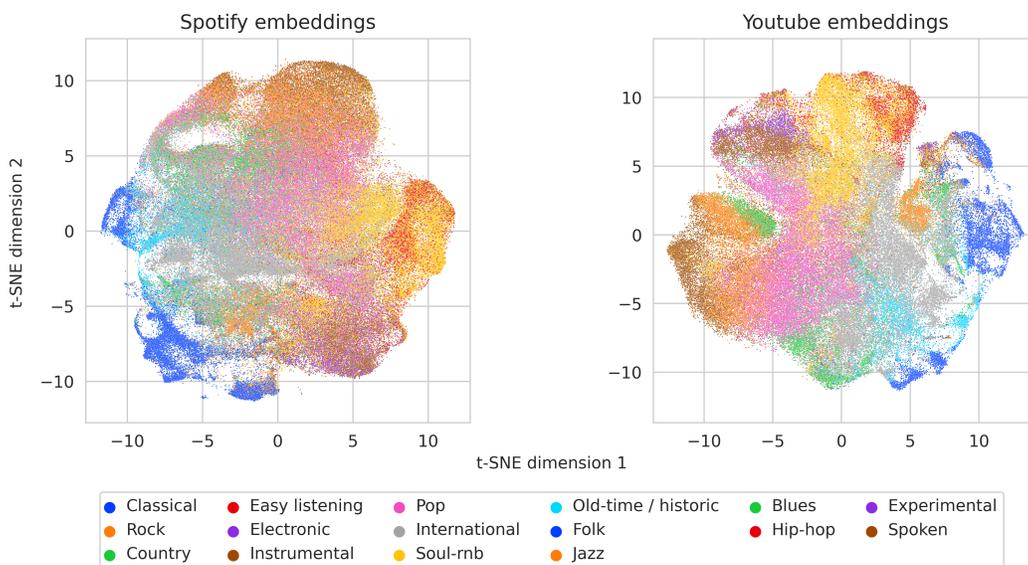


Figure 8: t-SNE plots for Spotify preview embeddings and YouTube audio embeddings computed with CLAP. Colors represent FMA root genres that were computed from the Spotify embeddings with zero-shot genre classification. The t-SNE plots show the relative positions of music samples, where samples with similar embeddings in the CLAP latent space are located closer to each other and dissimilar points farther apart. We observe that genres are well separated for both the YouTube embeddings and Spotify embeddings.

## 482 **E FMA Genre Analysis**

483 We repeat the zero-shot genre classification from Section 4 for the 16 FMA root genres. Figure 7  
484 shows the genre distribution for the FMA genres, while Figure 5 depicts the same t-SNE plot as  
485 shown in Figure 5 for these genres. We can observe that our results between overlapping genres are  
486 consistent and although there are more genres, we can observe meaningful relationships between  
487 them.

## 488 **F Subset Details**

489 As described in Section 3.3, we provide three different subsets of DISCO-10M.

490 DISCO-10K-random contains 10,000 random samples from DISCO-10M. We select the samples  
491 randomly from unique Spotify track IDs, meaning that the dataset will contain exactly one YouTube  
492 video link per Spotify song.

493 DISCO-200k-high-quality contains 200,000 high-quality samples filtered more strictly to improve the  
494 quality of the matches. We created this subset by filtering DISCO-10M with  $(\delta_a > 0.7) \wedge (\delta_{yt} > 0.8)$ .

## 495 **G Audio Characteristics**

496 The total playtime of YouTube videos in DISCO-10M is around 1,062,604 hours or 121 years.

497 We provide further insights regarding the audio attributes of the sample rate, MP3 file bitrate, and  
498 number of channels on the DISCO-10K-random subset. When downloading Spotify previews in MP3  
499 format, the audio quality and characteristics remain consistent. All audited samples share common  
500 attributes: a sample rate of 44.1 kHz, a bitrate of 96 kbps, and a 2-channel stereo setup.

501 In the case of audio from YouTube, there is a noticeable resemblance, albeit slightly less uniform.  
502 99.78% of the examined videos employ a standard stereo 2-channel configuration. 0.18% of videos  
503 are mono channel, while 0.04% utilize 6 channel surround for audio output. 99.96% of videos have a  
504 sample rate of 44.1 kHz, aligning with the settings of Spotify previews—the remaining 0.04% deviate  
505 by having a sample rate of 48 kHz.

## 506 **J Ethical Considerations**

507 Copyright and licensing agreements are a complex issue, particularly when it comes to big data  
508 collection for training large machine learning models. We acknowledge the concerns of artists  
509 regarding the potential negative impact on their artistic work. However, we believe that openly  
510 sharing such data helps democratize the research of music understanding and music creation.

511 To address artists who disagree with our assessment, we offer two options for reconciliation. First,  
512 artists can contact us at *anonymized* to request the removal of links associated with their art from  
513 our dataset. Second, artists may choose to take down their YouTube video or Spotify song from the  
514 respective platform, rendering the link contained in DISCO-10M invalid. It is important to note that  
515 our guarantee applies solely to our dataset, while other entities who hold private audio datasets may  
516 not offer the same level of control.

517 Creating a dataset of this magnitude is achievable using publicly available tools and a reasonable  
518 timeframe. By making our dataset open-source, we also aim to raise awareness on the ease of creating  
519 big datasets and uncover the potential existence of similar datasets held by private institutions. Our  
520 goal is to provide an opportunity for anyone interested to explore ideas with this dataset, and to  
521 enhance our understanding of music creation and safety with large datasets. We emphasize that this  
522 dataset serves as a starting point, pushing the boundaries and fostering research of enhanced datasets  
523 for various tasks in machine learning for music. Access to such extensive datasets is crucial, not only  
524 in the visual domain as demonstrated by Laion-5B, but also in the domain of music.

525 In summary, our ethical framework emphasizes the importance of respecting artists' concerns,  
526 providing options for data exclusion, promoting transparency in dataset creation, and facilitating  
527 meaningful exploration of ML-assisted music creation while prioritizing safety considerations.

528 **K Datasheet for Datasets**

529 **K.1 Motivation**

530 1. **For what purpose was the dataset created?** Was there a specific task in mind? Was there  
531 a specific gap that needed to be filled? Please provide a description.

532 • We want to provide an open-source large-scale music dataset for the research com-  
533 munity. Such large datasets do not yet exist in this domain, and we believe they are  
534 needed to democratize innovation in music research and ML-assisted music creation.  
535 Working with large data also has inherent risks, which are better analyzed openly by a  
536 large community rather than by private institutions behind closed doors.

537 2. **Who created the dataset (e.g., which team, research group) and on behalf of which**  
538 **entity (e.g., company, institution, organization)?**

539 • *Anonymized*

540 3. **Who funded the creation of the dataset?** If there is an associated grant, please provide the  
541 name of the grantor and the grant name and number.

542 • *Anonymized*

543 4. **Any other comments?**

544 • No.

545 **K.2 Composition**

546 1. **What do the instances that comprise the dataset represent (e.g., documents, photos,**  
547 **people, countries)?** Are there multiple types of instances (e.g., movies, users, and ratings;  
548 people and interactions between them; nodes and edges)? Please provide a description.

549 • We share 11,018,816 YouTube links to music with metadata and associated Spotify  
550 music metadata. In addition we provided similarity measures between the YouTube  
551 video title, the YouTube description, the Song title, and name of the artist. Additionally,  
552 contribution includes providing audio embeddings for the YouTube video and the  
553 Spotify song preview computed with Laion-CLAP [35]. The metadata includes an  
554 explicit flag to allow users to filter for explicit or non-explicit music.

555 2. **How many instances are there in total (of each type, if appropriate)?**

556 • 11,018,816

557 3. **Does the dataset contain all possible instances or is it a sample (not necessarily random)**  
558 **of instances from a larger set?** If the dataset is a sample, then what is the larger set? Is the  
559 sample representative of the larger set (e.g., geographic coverage)? If so, please describe  
560 how this representativeness was validated/verified. If it is not representative of the larger set,  
561 please describe why not (e.g., to cover a more diverse range of instances, because instances  
562 were withheld or unavailable).

563 • No, DISCO-10M does not cover all artists on Spotify and only a selection of popular  
564 songs of those that we do consider. The YouTube search results only contain 20 matches  
565 we take into consideration. To improve the dataset quality, we filter out matches that  
566 do not meet the threshold described in Section 3.2.

567 4. **What data does each instance consist of?** “Raw” data (e.g., unprocessed text or images)  
568 or features? In either case, please provide a description.

569 • URLs to YouTube videos and Spotify song previews as well as song specific metadata,  
570 such as artist names, artist/song IDs, YouTube video title/description snippet, video  
571 views, duration, Spotify song duration and creation date

572 5. **Is there a label or target associated with each instance?** If so, please provide a description.

573 • No.

- 574 6. **Is any information missing from individual instances?** If so, please provide a description,  
575 explaining why this information is missing (e.g., because it was unavailable). This does not  
576 include intentionally removed information, but might include, e.g., redacted text.
- 577 • The artist names are not always known since we do not have this information for every  
578 artist ID in our dataset. This is the case in 3.46% of all datapoints. In addition we  
579 have 7.81% missing YouTube description snippets and 0.183% missing YouTube view  
580 counts.
- 581 7. **Are relationships between individual instances made explicit (e.g., users' movie ratings,  
582 social network links)?** If so, please describe how these relationships are made explicit.
- 583 • No.
- 584 8. **Are there recommended data splits (e.g., training, development/validation, testing)?** If  
585 so, please provide a description of these splits, explaining the rationale behind them.
- 586 • No.
- 587 9. **Are there any errors, sources of noise, or redundancies in the dataset?** If so, please  
588 provide a description.
- 589 • We acknowledge the existence of duplicate songs stemming from different YouTube  
590 videos corresponding to the same Spotify song. These duplicates can be removed by  
591 filtering with stricter thresholds (cf. Section 3.2).
- 592 10. **Is the dataset self-contained, or does it link to or otherwise rely on external resources  
593 (e.g., websites, tweets, other datasets)?** If it links to or relies on external resources, a) are  
594 there guarantees that they will exist, and remain constant, over time; b) are there official  
595 archival versions of the complete dataset (i.e., including the external resources as they  
596 existed at the time the dataset was created); c) are there any restrictions (e.g., licenses, fees)  
597 associated with any of the external resources that might apply to a dataset consumer? Please  
598 provide descriptions of all external resources and any restrictions associated with them, as  
599 well as links or other access points, as appropriate.
- 600 • DISCO-10M relies on the availability of the songs on YouTube and Spotify since we  
601 only link to those resources. The embeddings and other metadata are self-contained.
- 602 11. **Does the dataset contain data that might be considered confidential (e.g., data that is  
603 protected by legal privilege or by doctor– patient confidentiality, data that includes the  
604 content of individuals' non-public communications)?** If so, please provide a description.
- 605 • No, there are no confidential datapoints in DISCO-10M.
- 606 12. **Does the dataset contain data that, if viewed directly, might be offensive, insulting,  
607 threatening, or might otherwise cause anxiety?** If so, please describe why.
- 608 • DISCO-10M contains music with an explicit flag. We do not know in what ways the  
609 song is explicit (sexual, abusive or others) but the flag allows users to easily filter for  
610 such songs. Additionally, DISCO-10M does not contain any links to age-restricted  
611 YouTube video.
- 612 13. **Does the dataset identify any subpopulations (e.g., by age, gender)?** If so, please  
613 describe how these subpopulations are identified and provide a description of their respective  
614 distributions within the dataset.
- 615 • No.
- 616 14. **Is it possible to identify individuals (i.e., one or more natural persons), either directly or  
617 indirectly (i.e., in combination with other data) from the dataset?** If so, please describe  
618 how.
- 619 • Yes, the Spotify artist ID is directly related to one or multiple natural persons. Addi-  
620 tionally, the YouTube video URLs we provide in the dataset are uploaded by one or  
621 multiple natural persons.

- 622 15. **Does the dataset contain data that might be considered sensitive in any way (e.g.,**  
623 **data that reveals race or ethnic origins, sexual orientations, religious beliefs, political**  
624 **opinions or union member- ships, or locations; financial or health data; biometric**  
625 **or genetic data; forms of government identification, such as social security numbers;**  
626 **criminal history)?** If so, please provide a description.
- 627 • No.
- 628 16. **Any other comments?**
- 629 • We emphasize the focus of our dataset on music, and not on individuals. Additionally,  
630 we reiterate that this dataset is intended for research purposes only, as described in  
631 Section 5.

### 632 K.3 Collection Process

- 633 1. **How was the data associated with each instance acquired?** Was the data directly ob-  
634 servable (e.g., raw text, movie ratings), reported by subjects (e.g., survey responses), or  
635 indirectly inferred/derived from other data (e.g., part-of-speech tags, model-based guesses  
636 for age or language)? If the data was reported by subjects or indirectly inferred/derived from  
637 other data, was the data validated/verified? If so, please describe how.
- 638 • The YouTube videos and metadata, and the Spotify tracks and metadata are observable  
639 and were collected by accessing the Spotify API as well as the YouTube API. The  
640 similarity scores and audio embeddings are computed by us.
- 641 2. **What mechanisms or procedures were used to collect the data (e.g., hardware appara-**  
642 **tuses or sensors, manual human curation, software programs, software APIs)?** How  
643 were these mechanisms or procedures validated?
- 644 • The Spotify API and the YouTube API. Our results were validated manually by assess-  
645 ing the quality of the retrieved information on random samples.
- 646 3. **If the dataset is a sample from a larger set, what was the sampling strategy (e.g.,**  
647 **deterministic, probabilistic with specific sampling probabilities)?**
- 648 • We started the Spotify artist scraping from the artist seed described in Appendix B.  
649 Additionally, we filter high-quality datapoints as described in Section 3.2.
- 650 4. **Who was involved in the data collection process (e.g., students, crowdworkers, contrac-**  
651 **tors) and how were they compensated (e.g., how much were crowdworkers paid)?**
- 652 • Only the authors of this paper were involved in the data collection process. *Author*  
653 *involvement and payment disclosed after acceptance.*
- 654 5. **Over what timeframe was the data collected?** Does this timeframe match the creation  
655 timeframe of the data associated with the instances (e.g., recent crawl of old news articles)?  
656 If not, please describe the time- frame in which the data associated with the instances was  
657 created.
- 658 • January 2023 to June 2023 was the timeframe of data collection. The creation time of  
659 the songs is diverse and can be seen in Figure 4.
- 660 6. **Were any ethical review processes conducted (e.g., by an institutional review board)?**  
661 If so, please provide a description of these review processes, including the outcomes, as well  
662 as a link or other access point to any supporting documentation.
- 663 • No.
- 664 7. **Did you collect the data from the individuals in question directly, or obtain it via third**  
665 **parties or other sources (e.g., websites)?**
- 666 • We collected data from Spotify and YouTube, not from artists directly.
- 667 8. **Were the individuals in question notified about the data collection?** If so, please describe  
668 (or show with screenshots or other information) how notice was provided, and provide a link  
669 or other access point to, or otherwise reproduce, the exact language of the notification itself.

- 670 • We did not notify any individuals about the data collection.
- 671 9. **Did the individuals in question consent to the collection and use of their data?** If so,
- 672 please describe (or show with screenshots or other information) how consent was requested
- 673 and provided, and provide a link or other access point to, or otherwise reproduce, the exact
- 674 language to which the individuals consented.
- 675 • We link to publicly available music on Spotify and YouTube. We allow every artist
- 676 contained in our dataset to have their entries removed upon request.
- 677 10. **If consent was obtained, were the consenting individuals provided with a mechanism to**
- 678 **revoke their consent in the future or for certain uses** If so, please provide a description,
- 679 as well as a link or other access point to the mechanism (if appropriate).
- 680 • Artists have the possibility to search our dataset for their YouTube video links, and their
- 681 Spotify artist ID and track IDs. If artists wish to remove their content from YouTube or
- 682 Spotify, they can contact those parties or remove it themselves, this would result in our
- 683 links becoming invalid. Additionally, we allow artists to contact us at *anonymized* to
- 684 request the removal of their datapoints.
- 685 11. **Has an analysis of the potential impact of the dataset and its use on data subjects (e.g.,**
- 686 **a data protection impact analysis) been conducted?** If so, please provide a description of
- 687 this analysis, including the outcomes, as well as a link or other access point to any supporting
- 688 documentation.
- 689 • Yes, we discuss the implications of our data collection pipeline and of our dataset in
- 690 Appendix J.
- 691 12. **Any other comments?**
- 692 • No.

#### 693 K.4 Preprocessing/cleaning/labeling

- 694 1. **Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or bucket-**
- 695 **ing, tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances,**
- 696 **processing of missing values)?**
- 697 • We performed preprocessing by filtering, as described in Section 3.2. We do not
- 698 process videos that are marked as age-restricted by YouTube, and we provide the
- 699 explicit content flag from Spotify.
- 700 2. **Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data (e.g., to**
- 701 **support unanticipated future uses)?** If so, please provide a link or other access point to
- 702 the “raw” data.
- 703 • No.
- 704 3. **Is the software that was used to preprocess/clean/label the data available?** If so, please
- 705 provide a link or other access point.
- 706 • We use `spotipy` to access the Spotify API, `youtubearchpython` to query the
- 707 YouTube search, and `pytube` to access the video on YouTube.
- 708 4. **Any other comments?**
- 709 • No.

#### 710 K.5 Uses

- 711 1. **Has the dataset been used for any tasks already?** If so, please provide a description.
- 712 • No.
- 713 2. **Is there a repository that links to any or all papers or systems that use the dataset?** If
- 714 so, please provide a link or other access point.

715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755  
756  
757  
758

- No.
3. **What (other) tasks could the dataset be used for?**
- We encourage the research community to use the dataset for music analysis, video analysis, music information retrieval, generative models for music, music genre recognition, as well as other possible down-stream tasks enabled by the provided embeddings.
4. **Is there anything about the composition of the dataset or the way it was collected and preprocessed/cleaned/labeled that might impact future uses?** For example, is there anything that a dataset consumer might need to know to avoid uses that could result in unfair treatment of individuals or groups (e.g., stereotyping, quality of service issues) or other risks or harms (e.g., legal risks, financial harms)? If so, please provide a description. Is there anything a dataset consumer could do to mitigate these risks or harms?
- Yes, as discussed in Section 5.
5. **Are there tasks for which the dataset should not be used?** If so, please provide a description.
- We strongly advise to use DISCO-10M only for research purposes and not for commercial applications.
6. **Any other comments?**
- No.

## K.6 Distribution

1. **Will the dataset be distributed to third parties outside of the entity (e.g., company, institution, organization) on behalf of which the dataset was created?** If so, please provide a description.
- Yes, the dataset will be open-source.
2. **How will the dataset will be distributed (e.g., tarball on website, API, GitHub)?** Does the dataset have a digital object identifier (DOI)?
- The dataset will be available on Hugging Face Datasets. DOI: 10.57967/hf/0754
3. **When will the dataset be distributed?**
- Starting from 14.06.2023.
4. **Will the dataset be distributed under a copyright or other intellectual property (IP) license, and/or under applicable terms of use (ToU)?** If so, please describe this license and/or ToU, and provide a link or other access point to, or otherwise reproduce, any relevant licensing terms or ToU, as well as any fees associated with these restrictions.
- CC-BY-4.0
5. **Have any third parties imposed IP-based or other restrictions on the data associated with the instances?** If so, please describe these restrictions, and provide a link or other access point to, or otherwise reproduce, any relevant licensing terms, as well as any fees associated with these restrictions.
- We do not own the copyright of the music accessible through the provided links.
6. **Do any export controls or other regulatory restrictions apply to the dataset or to individual instances?** If so, please describe these restrictions, and provide a link or other access point to, or otherwise reproduce, any supporting documentation.
- No.
7. **Any other comments?**
- No.

759 **K.7 Maintenance**

760 1. **Who will be supporting/hosting/maintaining the dataset?**

- 761 • Hugging Face Datasets will host the dataset and we will maintain the dataset.

762 2. **How can the owner/curator/manager of the dataset be contacted (e.g., email address)?**

- 763 • The authors can be contacted via *anonymized* email.

764 3. **Is there an erratum?** If so, please provide a link or other access point.

- 765 • Not initially, will be started when necessary, and will be documented with future  
766 releases.

767 4. **Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete in-**  
768 **stances)?** If so, please describe how often, by whom, and how updates will be communicated  
769 to dataset consumers (e.g., mailing list, GitHub)?

- 770 • No, except for updates due to removal of dataset entries. Updates will be communicated  
771 on Hugging Face Datasets.

772 5. **If the dataset relates to people, are there applicable limits on the retention of the data**  
773 **associated with the instances (e.g., were the individuals in question told that their data**  
774 **would be retained for a fixed period of time and then deleted)?** If so, please describe  
775 these limits and explain how they will be enforced.

- 776 • Artists may contact us to have entries excluded from our dataset.

777 6. **Will older versions of the dataset continue to be supported/hosted/maintained?** If so,  
778 please describe how. If not, please describe how its obsolescence will be communicated to  
779 dataset consumers.

- 780 • There are no older versions of DISCO-10M.

781 7. **If others want to extend/augment/build on/contribute to the dataset, is there a mech-**  
782 **anism for them to do so?** If so, please provide a description. Will these contributions  
783 be validated/verified? If so, please describe how. If not, why not? Is there a process for  
784 communicating/distributing these contributions to dataset consumers? If so, please provide  
785 a description.

- 786 • Updating and extending the dataset will be done on a case-by-case basis.

787 8. **Any other comments?**

- 788 • No.