

# ARIA-MIDI: A DATASET OF PIANO MIDI FILES FOR SYMBOLIC MUSIC MODELING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

We introduce an extensive new dataset of MIDI files, created by transcribing audio recordings of piano performances into their constituent notes. The data pipeline we use is multi-stage, employing a language model to autonomously crawl and score audio recordings from the internet based on their metadata, followed by a stage of pruning and segmentation using an audio classifier. The resulting dataset contains over one million distinct MIDI files, comprising roughly 100,000 hours of transcribed audio. We provide an in-depth analysis of our techniques, offering statistical insights, and investigate the content by extracting metadata tags, which we also provide.

## 1 INTRODUCTION

Central to the success of deep learning as a paradigm has been the datasets used to train neural networks. With the rapid technical advancements and ever increasing availability of computational power, music has become a popular target for deep learning research, and deep learning in turn has had a notable impact on the study and creation of musical works (Briot et al., 2019). These advancements in music-oriented deep learning rely heavily on diverse and well-structured datasets. Music is inherently structured and can be represented computationally in a variety of forms (Wiggins, 2016). Of particular relevance to this work are symbolic representations of music, such as MIDI (Musical Instrument Digital Interface), which are used by people and machines alike in both the generative process and the analysis of compositions.

In fields outside of computational music, the significance of comprehensive datasets is ubiquitous. For example, in computer vision, the ImageNet dataset (Deng et al., 2009) catalyzed research for almost a decade, providing both high-quality training data and robust benchmarks. Similarly, datasets such as Common Crawl (2024), C4 (Raffel et al., 2020), and the Pile (Gao et al., 2020) have been instrumental in advancing natural language processing. These resources have enabled the study of scaling-based approaches towards language modelling, enhanced available techniques, and provided foundation models for researchers with restricted resources.

The situation for music-oriented deep learning research is mixed. While numerous publicly available audio-based datasets exist (Gemmeke et al., 2017; Hawthorne et al., 2018; Thickstun et al., 2016; Bertin-Mahieux et al., 2011), symbolic datasets, which represent music in formats like MIDI, are comparatively lacking in both quality and quantity. The Lakh dataset (Raffel, 2016), comprising 176,581 MIDI files scraped from the internet, has been widely adopted in model training (Thickstun et al., 2023; Zeng et al., 2021) due to its scale. However, its files, often created through software sequencing or digital score conversion, often lack the nuances of expressive human performances, and vary significantly in quality. In contrast, the MAESTRO dataset (Hawthorne et al., 2018) offers high-quality Disklavier MIDI recordings from professional pianists, capturing the subtleties of human interpretation. However, its size and focus on virtuosic classical piano performances limits its applicability across diverse musical genres and compositional styles.

Underlying this is a common limitation: the manual transcription process creates bottlenecks for both the scale and quality. In recent years, researchers have turned to Automatic Music Transcription (AMT) (Benetos et al., 2019) to address these limitations, creating various large-scale symbolic datasets (Kong et al., 2020; Zhang et al., 2022; Edwards et al., 2023). Leading AMT techniques leverage neural networks to extract symbolic note-level information from audio (Sigtia et al., 2016;

Table 1: Comparison of publicly available datasets of symbolic music.

Dataset	# Files	# Hours	Genre	Source	Multi-track
MAESTRO	1,276	199	Classical	Piano Competitions	No
Mutopia <sup>1</sup>	1,862	69	Mixed	Lilypond	Yes
PiJAMA	2,777	223	Jazz	AMT	No
GiantMIDI	10,855	1,237	Classical	AMT	No
ATEPP	11,742	1,009	Classical	AMT	No
Lakh <sup>2</sup>	176,581	9,567	Mixed	Web-scrape	Yes
Aria-MIDI <sup>3</sup>	1,186,253	100,209	Mixed	AMT	No

<sup>1</sup> <https://www.mutopiaproject.org/><sup>2</sup> Size of the full dataset, including corrupted files. The commonly used *clean* subset contains 45,129 files.<sup>3</sup> Reduces to 745,109 files and 62,484 hours after compositional deduplication described in Section 4.

Hawthorne et al., 2017; Kong et al., 2021), theoretically enabling symbolic datasets to match the scale found in other modalities. Nevertheless, several challenges persist:

**Transcription Quality.** Some music forms, like solo-piano recordings, translate more accurately to MIDI than others. Additionally, training neural-AMT models relies on a small number of specific high-quality datasets of aligned audio-MIDI, e.g., MAESTRO for solo-piano. This limitation restricts use cases outside the distribution of the training data, leading to degraded transcriptions of recordings in different genres or with audio artifacts (Marták et al., 2024; Edwards et al., 2024).

**Pre-processing.** A dichotomy between quality and scale still presents itself. Current methods employed for audio curation and pre-processing (e.g., selection, pruning, and segmentation) are insufficient when applied to noisy and diverse audio corpora without human oversight. This is partially due to a lack of training data precisely labelled for the nuances of this application. Datasets maintaining high-quality standards have utilized a stage of machine-guided manual human verification to remove falsely identified audio (Zhang et al., 2022; Edwards et al., 2023), an approach that does not scale well.

In this work we address these challenges, focusing on creating a comprehensive dataset of piano transcriptions. We demonstrate that with strategic modifications to the data pipeline, AMT-based approaches can be scaled effectively. Our approach leverages a robust piano-AMT model [REDACTED], capable of accurately transcribing recordings across diverse timbres and recording qualities. While this model is crucial to our process, in this work we concentrate on developing techniques for precise curation, pre-processing, and metadata attribution for publicly available audio files. These techniques enable the creation of a large-scale, high-quality dataset suitable for various music information retrieval tasks and generative applications.

## 1.1 CONTRIBUTIONS OF THIS PAPER

More specifically, our contributions are as follows:

1. We introduce a scalable, language model-guided method for crawling and extracting metadata from specific types of videos. We analyze the effectiveness of this approach in the context of publicly available piano recordings.
2. We outline a process for distilling an audio source-separation model to train a classifier capable of accurately identifying and segmenting realistic piano recordings, which we open-source. This enabled an 8-fold improvement in identification of non-piano audio without human supervision when compared to previous work.
3. Using an existing piano-transcription model, we provide a new MIDI dataset of piano transcriptions, one of the largest and cleanest to date.

We hope the dataset released alongside this work has a positive impact on the MIR research community. We foresee several potential areas where it may accelerate research. Firstly, pretrained

generative models have had a large impact on the textual and visual domains (Zhou et al., 2023). These models rely on datasets typically in terabytes. Comprising approximately 20 gigabytes of MIDI files, Aria-MIDI isn't on this scale; however, it may still be useful for research into pretrained music models. Secondly, we are releasing accurate compositional metadata for each file, as well as piano audio-classifier scores, which due to our training methods can act as a proxy for recording quality. This information is valuable for many MIR tasks (Choi et al., 2017), as well as for making clean and *compositionally deduplicated* subsets.

## 1.2 RELATED WORK

The use of neural networks for automatic music transcription has its roots in the seminal work of Sigtia et al. (2015; 2016). This was followed by various works experimenting with different approaches and neural architectures (Hawthorne et al., 2017; 2021; Yan et al., 2021; Toyama et al., 2023). The high-resolution piano transcription model introduced in Kong et al. (2021), trained using the MAESTRO (Hawthorne et al., 2018) and MAPS (Emiya et al., 2010) datasets, became the de facto benchmark for accuracy. More recently, AMT research has extended to other instruments (Riley et al., 2024) and multi-track transcription (Gardner et al., 2021; Chang et al., 2024), where it has seen success.

There are three predominant publicly available datasets of piano transcriptions, all of which utilized the transcription model introduced in Kong et al. (2021). GiantMIDI (Kong et al., 2020) was the first, comprising transcriptions of piano recordings matching names of musical works taken from the IMSLP (2006). From 143,701 initial recordings, 10,855 were identified by a model trained to detect solo-piano recordings. The ATEPP dataset (Zhang et al., 2022) took a different approach than GiantMIDI, focusing on repeat performances of standard classical piano repertoire, and using text-based techniques to determine opus and piece numbers. PiJAMA (Edwards et al., 2023), a dataset of jazz piano transcriptions, spans 120 different pianists across 244 recorded albums. Recordings were curated by matching tracks from albums performed by a manually curated list of pianists.

All three datasets utilized YouTube to match musical metadata with audio recordings and employed audio-based classifiers, trained using MAESTRO and AudioSet, to identify piano recordings. For ATEPP and PiJAMA, these classifiers were also used to remove applause and speech. The level of human intervention varied across datasets: GiantMIDI relied solely on automated processes, while ATEPP and PiJAMA incorporated manual checks. Table 1 presents a comparison of these datasets in context.

## 2 METHODOLOGY

In this section, we describe the methodology used to compile our dataset of MIDI files. Our approach consists of three distinct stages. First, we assemble a comparatively large corpus of candidate piano recordings using low-overhead, text-based methods. Next, we employ audio-based techniques to refine our initial corpus through pruning and pre-processing. Finally, we conduct a computationally intensive stage of transcription and metadata extraction. This multi-stage approach allows us to efficiently process a large volume of data while ensuring high-quality results in our final dataset.

### 2.1 CRAWLING

A common theme in previous work has been to compile candidate recordings by first obtaining a corpus of metadata (e.g., composers, performers, album titles) via various means, and then using the APIs provided by Spotify<sup>1</sup> and YouTube<sup>2</sup> to match piece titles to corresponding videos on YouTube. We take a different approach: Our method begins with a small collection of manually curated seed videos and uses YouTube's API to crawl related content. The crawling priority is determined by a language model, which performs two tasks: 1) parsing the title and description of each video, and 2) scoring the likelihood of the video containing solo-piano content on a scale of 1-5. More specifically, starting from fifty solo-piano seed videos which span a variety of genres and styles, we follow a two-step procedure which we cycle through repeatedly:

<sup>1</sup><https://developer.spotify.com/documentation/web-api>

<sup>2</sup><https://developers.google.com/youtube/v3>

1. For each unscored video, we prompt a language model with its YouTube title and description. The model is tasked with assigning a score from one to five, indicating the likelihood that the video contains a solo-piano performance.
2. In order of priority determined by the score, we use the YouTube API to fetch related video URLs, titles, and descriptions.

By taking advantage of the related videos endpoint of the YouTube API, we outsource the majority of the crawling process to YouTube’s own recommendation algorithm. We used the 70B parameter version of Llama 3.1 (Dubey et al., 2024) for the language model, observing that smaller models made obvious mistakes more frequently. The system prompt we used can be found in Appendix A.1. Overall, we found this process to be effective. Although initially this procedure tended to overrepresent recordings of well-known classical pieces, as these became less available later in the process, piano recordings representing diverse set of musical styles and genres were crawled thoroughly.

## 2.2 AUDIO CLASSIFICATION AND SEGMENTATION

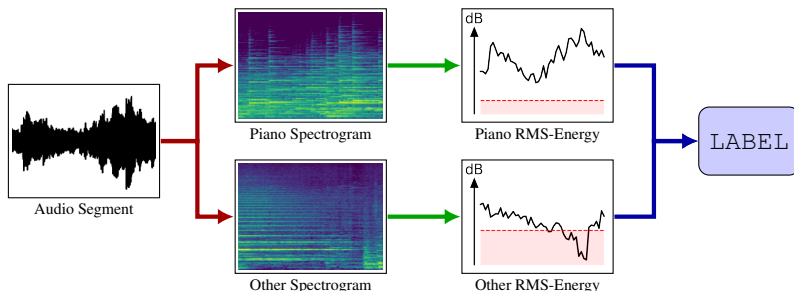


Figure 1: A visual representation of the pseudo-labeling process applied to a five-second excerpt from a piano concerto. As the non-piano component has a contiguous region with sufficient energy, this example is labelled as non-piano.

As we analyze in-depth in Section 3, relying solely on the score attributed to each recording during the crawling process results in an unacceptably high rate of misclassifications. Following other work (Kong et al., 2020; Zhang et al., 2022), we address this by using an audio classification model (Dieleman and Schrauwen, 2014) in the next stage of our pipeline. We identified the following problematic situations which we aimed to address with an audio classifier:

- Misclassifications due to logical mistakes by the language model or misleading/ambiguous YouTube data: The classifier identifies and removes such recordings while allowing us to retain those with some ambiguity, which would otherwise have to be pruned.
- Undesirable acoustic qualities in positively classified recordings: Despite positive classification by the language model, recordings can be inappropriate for transcription for a variety of reasons including incorrect instrumentation (e.g., harpsichord, organ, electric piano), low audio quality, or the presence of additional instruments. To mitigate this, we include representative examples from these categories in the training data for our classifier.
- Non-piano content in high-quality piano recordings: Many otherwise high-quality piano recordings contain segments of non-piano content, such as applause, commentary, or extended periods of silence. Using the algorithm we describe in Section 2.2.2, we adapt the classifier to segment recordings into contiguous regions of solo-piano performance, removing unwanted content.

A primary concern when building an audio classifier is the quality, diversity, and accuracy of the labels used for training. Initial investigations revealed that relying on well-known datasets such as MAESTRO (Hawthorne et al., 2018) and AudioSet (Gemmeke et al., 2017) was insufficient, as we display in Table 4. In an effort to achieve classification accuracy approximating human labels, we curated a mixed training dataset, representative of our corpus of crawled recordings. We used a

Table 2: Overview of supervised and unsupervised audio corpora used to train the piano audio classifier. Prop. denotes the proportion of solo-piano labels in the corresponding component. Notably, the raw GiantMIDI files, pre-processed with source-separation, only contain 87.35% solo-piano labels, and our collection of piano concertos contained 14.98% solo-piano labels.

Component	Length (h)	Weight (%)	Pseudo lab.	Prop. (%)	dB <sub>min</sub>	l <sub>min</sub> (s)
GiantMIDI	1040	43.24	True	87.35	-25dB	1.5
Score-4	676	28.11	True	44.81	-22dB	1.5
MAESTRO	198	8.23	False	99.62	N/A	N/A
Synthetic Data	143	5.95	False	98.34	N/A	N/A
Jazz Trio Database	139	5.78	True	3.61	-28dB	1.0
Piano Other	75	3.12	False	99.22	N/A	N/A
Non-Piano Other	71	2.95	False	0.00	N/A	N/A
Symphonies	40	1.66	False	0.00	N/A	N/A
Piano Concertos	23	0.96	True	14.98	-28dB	1.0
Total	2405	100.0				

novel approach, leveraging an audio source-separation model to accurately generate pseudo-labels for the unsupervised (unlabelled) parts of the training corpus.

Given an audio file, we used the MVSep Piano source-separation model (Uhlich et al., 2024; Fabbro, 2024; Solovyev et al., 2023) to decompose it into its constituent parts, separating the piano component from the other audio content. For each five-second clip, we resample each component to 22,050 Hz and calculate a spectrogram with 2048 bins, using a frame and hop length of 2048 and 512 respectively. By calculating the RMS energy of each frame and converting to dBFS (Zölzer, 2022), we classify contiguous segments above an energy threshold as *non-silent*. Given parameters dB<sub>min</sub> and l<sub>min</sub>, we classify the five-second audio clip as *non-piano* if the component labeled *other* has a non-silent region longer than l<sub>min</sub>, according to the energy level dB<sub>min</sub>. Similarly, we classify the audio clip as non-piano if the component labeled *piano* has a contiguous silent region, according to the energy level -20dB, of length greater than four seconds. This process is visualised in Figure 1.

We applied this labelling procedure to various collections of publicly available audio files, displayed in Table 2, the main constituent being 10,000 YouTube videos from our corpus which were given a score of four by the language model. We also used the GiantMIDI audio files, the Jazz Trio Database (Cheston et al., 2024), and smaller collections of piano and non-piano recordings which we curated manually. Including the pseudo-labelled audio allows us to distill the source-separation model (Hinton et al., 2015), bypassing the high computational cost associated with applying source separation to our entire inference corpus, which we estimate at about 5,000 A100 hours<sup>3</sup>.

### 2.2.1 TRAINING

For our solo-piano classifier, we chose a CNN-based architecture (LeCun et al., 2015) with five convolutional layers followed by two dense layers and a single output neuron. The input to the classifier consists of mel-spectrograms calculated from five-second audio clips. We used a sample rate of 22,050 Hz, 2048 spectrogram frequency bins, 256 mel bins, and a hop length of 220 (corresponding to 10ms hops). We trained the model for ten epochs using the AdamW optimizer (Loshchilov and Hutter, 2019) with  $\beta_1, \beta_2 = 0.9, 0.95$ ,  $\epsilon = 1e-6$  and an L2 weight decay of 0.01. A linear learning rate scheduler was used, decaying to 10% of the initial learning rate after a warmup over the first 500 optimizer steps.

One consequence of training with pseudo-labels obtained using source separation was having to use relatively sensitive energy thresholds in order to correctly label training examples with a quiet but notable non-piano component. These thresholds occasionally result in incorrect training labels for solo-piano recordings with significant but acceptable background audio artifacts like noise,

<sup>3</sup>In comparison, classification of 100,000 hours of audio using our model only took 20 A100 hours, I/O being the main bottle-neck.

distortion, and reverb. To mitigate this, we trained with the corresponding audio augmentation in approximately 10% of batches, as well as randomly applying pitch shifting and bandpass filters. We also included labeled examples, representative of such piano recordings, as part of our training data.

### 2.2.2 INFERENCE

As well as per-file classification, we also use our classification model to segment audio recordings into their standalone components of contiguous piano performance. To do this, we employ a sliding-window based technique adapted from standard approaches (Keogh et al., 2004), aimed at accurately removing non-piano content whilst being robust to short-lived classification mistakes.

Given an audio recording, we score each five-second interval, calculated with a stride of one second, by passing the inputs through our model. We classify a region  $(n, m + 5)$ ,  $m \geq n + d$ , as non-piano if and only if all segments starting between  $n$  and  $m$  are scored below  $\lambda$ . The parameters  $d$  and  $\lambda$  control the sensitivity and minimum length of non-piano segments, which we set to 3 and 0.5 respectively. After excluding all non-piano segments, we classify the remaining segments as piano if they are longer than 45 seconds. Finally, we discard piano segments with an average score lower than 0.7. Our choice of algorithm and hyperparameters was motivated to reduce the chance of a solo-piano segment being prematurely interrupted due to instability in scoring. As our classifier designates intervals which are mostly silent as non-piano, this approach also segments piano performances separated by at least  $d + 5$  seconds of silence. In Section 3 we investigate the accuracy of both classification and segmentation of our proposed approach.

## 2.3 TRANSCRIPTION

We used the piano-AMT model [REDACTED], introduced in [REDACTED], to transcribe the segmented audio recordings. This choice was informed by the model’s robustness in transcribing audio from a diverse set of recording environments, compared with models used in previous work (Kong et al., 2020; Zhang et al., 2022) (See Appendix A.3). Transcription of the 100,209 hours of audio took 765 hours using an NVIDIA H100 GPU with a batch size of 128, representing an inference speed of roughly 131x real-time at a rate of approximately 2327 transcribed notes per second.

## 2.4 METADATA EXTRACTION

Access to per-file metadata labels provides a straightforward way to create different dataset splits, appropriate for various generative and MIR tasks. A central concern of ours was *entity resolution* (ER) (Christen, 2011), i.e., identifying the compositional source of each recording and addressing overrepresentation of popular pieces in the corpus<sup>4</sup>.

Inspired by our crawling methodology, we chose a similar strategy for extracting metadata. We used Llama 3.1 (70B) to process YouTube titles and descriptions for files that passed language model and audio detection filters. The prompt (see Appendix A.2) extracted composer, opus numbers (e.g., Op., BWV, K., D.) and piece identifiers, as well as performer, genre, and form labels. Accurate metadata labels provide a way to quantify and control *compositional duplication*, as well as supervised labels useful for MIR. We investigate the accuracy of these labels and their distributions in Sections 3 and 4 respectively.

## 3 METHODOLOGICAL ANALYSIS

In this section, we evaluate the effectiveness of the components in our data pipeline. Where applicable, we compare our methods to those used in previous work, in particular the GiantMIDI, ATEPP, and PiJAMA datasets. For baselines and to determine ground truth, we relied on human labels obtained from two musically trained pianists familiar with popular classical and jazz repertoire.

**Language Model Classification.** We first analyze the ability of a language model to correctly classify a video solo-piano according to its YouTube title and description. In our experiment we chose a random sample of 250 videos from those crawled, and calculated the accuracy of the labels provided by different language models, judged relative to the audio content. We also asked human

<sup>4</sup>For example, *moonlight* appears in 6,819 titles, likely referring to Beethoven’s Moonlight Sonata.

Table 3: Classification precision, recall, and F1-scores for various language models and score classification thresholds. Results indicate that Llama 3.1 (70B) with a score threshold of 4 offers a appropriate balance between inference cost and accuracy.

Model	Score $\geq 3$			Score $\geq 4$			Score $\geq 5$		
	P (%)	R (%)	F1 (%)	P (%)	R (%)	F1 (%)	P (%)	R (%)	F1 (%)
Llama 3.1 8B	73.08	81.43	77.03	82.26	72.86	77.27	86.36	27.14	41.30
Llama 3.1 70B	64.76	97.14	77.71	70.83	<b>97.14</b>	81.93	84.51	85.71	85.11
Llama 3.1 405B	77.01	95.71	85.35	80.49	94.29	<b>86.84</b>	<b>94.44</b>	72.86	82.26
Human labels	73.63	95.71	83.23	83.56	87.14	85.31	85.71	25.71	39.56

participants to label the videos according to the same prompt given to the language models. Classification precision, recall, and F1 scores can be seen in Table 3. In comparison to human-derived labels, language models perform well at this text-based classification task. We attribute this to the depth of knowledge of different composers, performances, and pieces, which the language models have access to. Despite this, there remains a discrepancy between the audio ground truth and labels obtained from titles and descriptions alone.

**Audio Segmentation.** We evaluate the performance of our audio classification model in identifying and segmenting solo-piano content within audio recordings. For comparison, we used MVSep directly to obtain binary labels, applying the same inference procedure as described in Section 2.2.2. For ablation, we trained a model without the pseudo-labeled training data listed in Table 2. To mimic classifiers used for segmentation in other work, notably for the GiantMIDI, ATEPP, and PiJAMA datasets, we include various noise, applause, and speech from AudioSet (Gemmeke et al., 2017) as negative training examples for our ablation model.

For this analysis, a random sample of 250 audio recordings with language model scores greater than or equal to 3 was selected, excluding those used during training. To establish ground truth, participants were tasked with segmenting recordings into regions of solo-piano content and assigning files into one of three categories: Not solo-piano, solo-piano with significant audio artifacts, or solo-piano with good to pristine recording quality. Human-labeled segments were post-processed in accordance with our inference algorithm: Non-piano segments shorter than eight seconds were ignored, and a minimum length of 45 seconds was imposed on piano segments. Segmentation accuracy results can be seen in Table 4. While the ablation model achieves accurate segmentation, being less likely to interrupt piano segments by misclassifying occasional noisy periods of extreme piano audio as non-piano content, it conversely mislabels non-piano audio as piano eight times more frequently than the proposed approach in absolute terms.

We next evaluated our model’s classification performance. To assess this on a per-file basis, we imposed minimum thresholds on the average score for predicted piano segments, negatively classifying files with no predicted piano segments after filtering. Table 5 reports the accuracy of this approach in identifying the non-piano recordings and the solo-piano recordings with significant audio artifacts, in our evaluation dataset. Additionally, we analyzed the audio files that constitute GiantMIDI, ATEPP, and PiJAMA. Our human participants manually categorized the files which fell below an empirically determined average score threshold of 0.7, which could indicate issues with recording quality or content. The resulting distributions of these categorizations are shown in Figure 2.

In both tasks, our approach performs well. For segmentation with  $\lambda=0.5$ , we achieve a 96.38% overlap with the ground truth for high-quality piano recordings, while removing 98.83% of non-piano audio over the full evaluation corpus. When additionally using a segment average-score threshold of  $T=0.7$ , this improves to 100%, while retaining 95.28% of the high-quality piano audio on a per-file basis. Furthermore, raising  $T$  to 0.9 removes the vast majority of low-quality audio, allowing us to curate a clean split of dataset, which we also provide.

**Metadata extraction.** We selected a random sample of 200 files and cross-referenced the metadata labels assigned by the language model with the YouTube titles and descriptions. For each file and metadata category, we manually checked for incorrect labels, e.g., misattributions, as well as labels which were missing despite information being present in the raw text. Results are displayed in

Table 4: Segmentation accuracy and overlap ratios for different techniques and hyperparameters. We consider a predicted segment correct if its beginning and end match the reference within tolerances of  $\pm 2$  seconds and  $\pm 5$  seconds, respectively. Each reference segment is matched to at most one predicted segment. Overlap ratios are calculated separately for piano and non-piano audio, each as the ratio of the duration of correctly identified audio to the total duration of the respective ground truth audio type.  $\text{dB}_{\min}$  and  $\lambda$  denote the sensitivity to non-piano content as described in Section 2.2.

Technique	Segmentation Accuracy			Segment Overlap	
	P (%)	R (%)	F1 (%)	Piano (%)	Non-Piano (%)
<i>Full corpus</i>					
MVsep, $\text{dB}_{\min}=-22\text{dB}$	65.62	70.47	67.96	91.96	97.22
MVsep, $\text{dB}_{\min}=-25\text{dB}$	58.28	63.76	60.90	88.74	98.18
MVsep, $\text{dB}_{\min}=-28\text{dB}$	49.38	53.69	51.45	82.23	98.66
Proposed, $\lambda=0.5$	<b>71.97</b>	<u>75.84</u>	<u>73.86</u>	<u>94.22</u>	98.83
Proposed, $\lambda=0.6$	70.70	74.50	72.55	92.67	<u>98.89</u>
Proposed, $\lambda=0.7$	68.39	71.14	69.74	91.04	<b>99.10</b>
Ablation, $\lambda=0.5$	<u>71.18</u>	<b>81.21</b>	<b>75.86</b>	<b>97.05</b>	91.10
<i>All solo-piano recordings</i>					
MVsep, $\text{dB}_{\min}=-22\text{dB}$	68.18	70.47	69.31	91.96	89.50
MVsep, $\text{dB}_{\min}=-25\text{dB}$	59.75	63.76	61.69	88.74	93.43
MVsep, $\text{dB}_{\min}=-28\text{dB}$	50.00	53.69	51.78	82.23	94.12
Proposed, $\lambda=0.5$	<u>72.44</u>	<u>75.84</u>	<u>74.10</u>	<u>94.22</u>	92.73
Proposed, $\lambda=0.6$	71.15	74.50	72.79	92.67	<u>93.18</u>
Proposed, $\lambda=0.7$	68.83	71.14	69.97	91.04	<b>94.66</b>
Ablation, $\lambda=0.5$	<b>77.56</b>	<b>81.21</b>	<b>79.34</b>	97.05	73.33
<i>Quality solo-piano recordings</i>					
MVsep, $\text{dB}_{\min}=-22\text{dB}$	70.00	75.97	72.86	94.71	85.13
MVsep, $\text{dB}_{\min}=-25\text{dB}$	61.64	69.77	65.45	92.06	87.04
MVsep, $\text{dB}_{\min}=-28\text{dB}$	50.33	58.91	54.29	87.41	<u>88.22</u>
Proposed, $\lambda=0.5$	<u>75.36</u>	<u>80.62</u>	<u>77.90</u>	<u>96.38</u>	84.80
Proposed, $\lambda=0.6$	73.91	79.07	76.40	95.33	86.13
Proposed, $\lambda=0.7$	71.53	75.97	73.68	93.92	<b>88.29</b>
Ablation, $\lambda=0.5$	<b>82.09</b>	<b>85.27</b>	<b>83.65</b>	<b>97.15</b>	83.13

Table 6. This manual verification confirms the accuracy of our metadata labeling process. Notably, in some instances, the language model provided accurate labels that were absent from the raw text.

Table 5: Classification performance for different segment average score thresholds calculated with on our human-labelled evaluation dataset. Segments were calculated using  $\lambda=0.5$ . *All solo-piano* measures the performance at identifying files with segments of solo-piano performance, regardless of audio artifacts, whereas *Quality solo-piano* measures the performance at identifying only recordings with good to pristine recording conditions. FP = False Positives.

Threshold	All solo-piano				Quality solo-piano			
	P (%)	R (%)	F1 (%)	FP	P (%)	R (%)	F1 (%)	FP
$T \geq 0.50$	99.28	93.24	96.17	1	88.49	96.85	92.48	16
$T \geq 0.60$	99.27	91.89	95.44	1	89.05	96.06	92.42	15
$T \geq 0.70$	100.00	89.86	94.66	0	90.98	95.28	93.08	12
$T \geq 0.80$	100.00	85.14	91.97	0	90.48	89.76	90.12	12
$T \geq 0.90$	100.00	75.68	86.15	0	95.54	84.25	89.54	5
$T \geq 0.95$	100.00	61.49	76.15	0	97.80	70.08	81.65	2



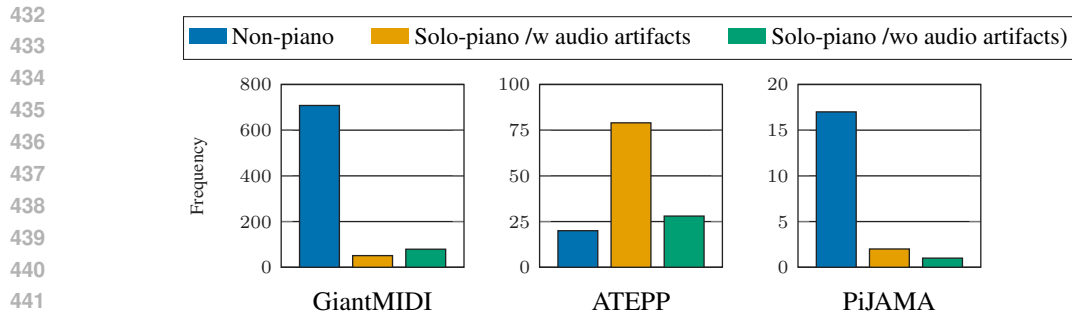


Figure 2: Distribution of files with average scores  $\leq 0.7$  across datasets. Files were manually categorized. Human-based pruning in ATEPP and PiJAMA explains distribution differences from GiantMIDI. See Table 1 for general dataset information.

Table 6: Analysis of metadata presence and accuracy across different attributes. For each attribute, presence indicates the percentage of files with assigned metadata, accuracy shows the percentage of correct labels among present metadata, and missed labels represents the percentage of files where metadata was omitted despite being inferrable from YouTube titles and descriptions. Accuracy was verified following the criteria specified in the system prompt (see Appendix A.2).

Attribute	Presence (%)	Accuracy (%)	Missed Labels (%)
Composer	71.0	99.3	2.7
Performer	62.0	99.2	0.8
Opus Number	32.0	100.0	1.5
Piece Number	22.0	93.2	4.3
Key Signature	23.0	97.8	0.0
Genre	86.5	94.2	0.6
Music Period	63.0	92.9	12.5

#### 4 DATASET STATISTICS

In this section, we analyze statistics about our methodology and the contents of the resulting dataset of MIDI files. Overall, when executing our data pipeline we collected YouTube data for 3,290,453 videos, from which we further processed 1,713,650 using our audio classifier. We then transcribed over 1 million audio segments into approximately 100,000 hours of transcribed solo-piano music. We present a breakdown of scores ascribed during each stage of processing in Figure 3. Taken together with the experiments in Section 3, we conclude that the techniques we have introduced

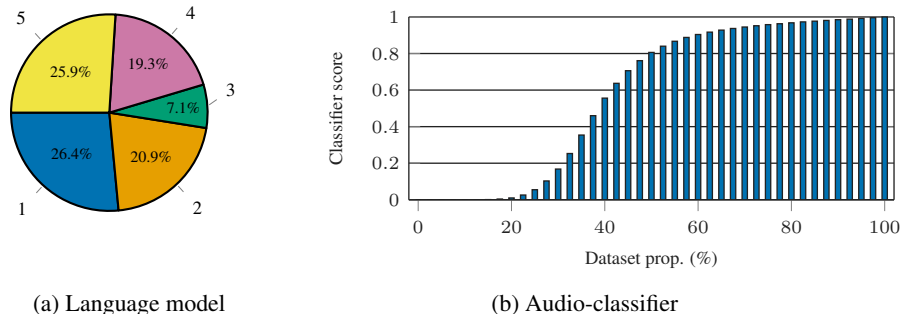


Figure 3: Score breakdowns for the different components of our data pipeline. Figure (a) displays the relative proportion of language model scores. Figure (b) shows the cumulative density of language model scores, calculated over recordings attributed a language model score of at least three.

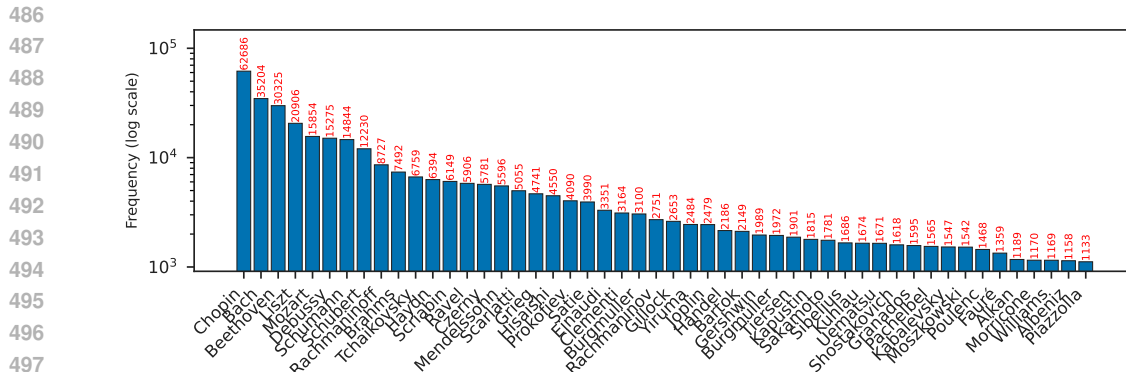


Figure 4: Number of transcriptions (log scale) in Aria-MIDI as attributed to the top 50 composers.

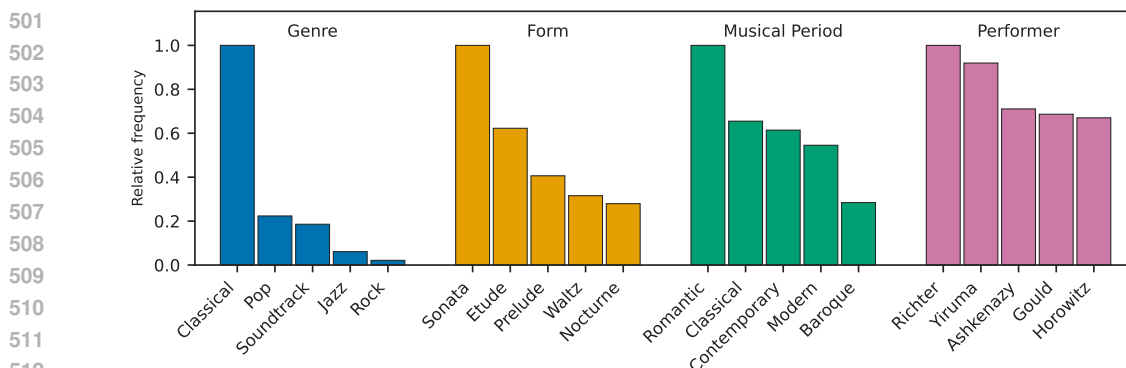


Figure 5: Relative frequency distribution of metadata across different categories, normalized to the most common category.

work well at scale. Moreover, extrapolating from the results in Table 5, we observe that the top-scoring 35,000 hours of MIDI files likely contain few transcriptions of non-solo-piano content.

To address compositional duplicates, we analyze metadata tags in three categories: composer, opus number, and piece number. To obtain a dataset split minimizing compositional duplicates within the text-based metadata constraints, we remove files that either match on all three tags (composer, opus number, and piece number) or match on both composer and opus number in cases where piece number tags are absent. For composers who appear more than 250 times across the dataset, we also prune all additional files that lack opus number and piece number tags. Overall, we identified 23,877 unique metadata triples, and after removing compositional duplicates using this procedure, 745,109 files remained. Figure 4 shows the frequency of performances by different composers in the complete dataset, illustrating their relative popularity.

Lastly, Figure 5 shows the distribution of metadata for other categories over the entire collection of MIDI files, without deduplication. Overwhelmingly, transcriptions of classical piano performances dominate; however, when accounting for the total size, many other genres are well represented.

## 5 CONCLUSION

We have introduced a new dataset of MIDI files, created by transcribing piano performances publicly accessible on the internet. In this paper, we provide an analysis of the components in our data pipeline and find them to be well-suited for our purposes. Going forward, we see several areas for future work: Primarily, extending our approach to other instruments such as guitar, as well as the multi-instrument case, could be approachable via variations of the source-separation-based approaches to audio pre-processing we have outlined. Secondly, further study into metadata attribution using language models, especially targeting improvements to compositional entity recognition.

## 540 COPYRIGHT DISCLAIMER

541  
542 The use of copyrighted works, and derivatives thereof such as MIDI transcriptions, for machine  
543 learning is a complex issue. To mitigate potential harms, we do not release audio files or raw  
544 metadata under direct copyright. We distribute this dataset under a CC-BY-NC-SA license ([Creative  
545 Commons](#)).

## 546 REFERENCES

- 547  
548 Emmanouil Benetos, Simon Dixon, Zhiyao Duan, and Sebastian Ewert. Automatic music transcrip-  
549 tion: An overview. *IEEE Signal Processing Magazine*, 36(1):20–30, 2019. doi: 10.1109/MSP.  
550 2018.2869928.
- 551  
552 Thierry Bertin-Mahieux, Daniel P.W. Ellis, Brian Whitman, and Paul Lamere. The million song  
553 dataset. In *Proceedings of the 12th International Conference on Music Information Retrieval*  
554 (*ISMIR 2011*), 2011.
- 555  
556 Jean-Pierre Briot, Gaëtan Hadjeres, and François-David Pachet. Deep learning techniques for music  
557 generation – a survey, 2019. URL <https://arxiv.org/abs/1709.01620>.
- 558  
559 Sungkyun Chang, Emmanouil Benetos, Holger Kirchhoff, and Simon Dixon. Yourmt3+: Multi-  
560 instrument music transcription with enhanced transformer architectures and cross-dataset stem  
561 augmentation. *arXiv preprint arXiv:2407.04822*, 2024.
- 562  
563 Huw Cheston, Joshua L Schlichting, Ian Cross, and Peter M C Harrison. Jazz trio database: Au-  
564 tomated annotation of jazz piano trio recordings processed using audio source separation. 2024.  
565 doi: 10.5334/tismir.186. URL <https://doi.org/10.5334/tismir.186>.
- 566  
567 Keunwoo Choi, György Fazekas, Kyunghyun Cho, and Mark Sandler. A tutorial on deep learning  
568 for music information retrieval. *arXiv preprint arXiv:1709.04396*, 2017.
- 569  
570 Peter Christen. A survey of indexing techniques for scalable record linkage and deduplication. *IEEE*  
571 *transactions on knowledge and data engineering*, 24(9):1537–1555, 2011.
- 572  
573 Common Crawl, 2024. URL <https://commoncrawl.org>.
- 574  
575 Creative Commons. Attribution-NonCommercial-ShareAlike 4.0 International (CC BY-NC-SA  
576 4.0). <https://creativecommons.org/licenses/by-nc-sa/4.0/>.
- 577  
578 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hi-  
579 erarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,  
580 pages 248–255. Ieee, 2009.
- 581  
582 Sander Dieleman and Benjamin Schrauwen. End-to-end learning for music audio. In *2014 IEEE*  
583 *International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6964–  
584 6968. IEEE, 2014.
- 585  
586 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha  
587 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.  
588 *arXiv preprint arXiv:2407.21783*, 2024.
- 589  
590 Drew Edwards, Simon Dixon, and Emmanouil Benetos. Pijama: Piano jazz with automatic midi  
591 annotations. *Transactions of the International Society for Music Information Retrieval*, 2023.
- 592  
593 Drew Edwards, Simon Dixon, Emmanouil Benetos, Akira Maezawa, and Yuta Kusaka. A data-  
594 driven analysis of robust automatic piano transcription. *IEEE Signal Processing Letters*, 2024.
- 595  
596 Valentin Emiya, Nancy Bertin, Bertrand David, and Roland Badeau. Maps-a piano database for  
597 multipitch estimation and automatic transcription of music. 2010.
- 598  
599 Uhlich S. Lai C.-H. Choi W. Martínez-Ramírez M. Liao W. Gadelha I. Ramos G. Hsu E. Rodrigues  
600 H. Stöter F.-R. Défossez A. Luo Y. Yu J. Chakraborty D. Mohanty S. Solovyev R. Stempkovskiy  
601 A. Habruseva T. Goswami N. Harada T. Kim M. Lee J. H. Dong Y. Zhang X. Liu J. Mitsufuji Y  
602 Fabbro, G. The sound demixing challenge 2023-music demixing track. *Trans. Int. Soc. Music.*  
603 *Inf. Retr.*, 7(1):63–84, 2024.

- 594 Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason  
595 Phang, Horace He, Anish Thite, Noa Nabeshima, et al. The pile: An 800gb dataset of diverse text  
596 for language modeling. *arXiv preprint arXiv:2101.00027*, 2020.
- 597 Josh Gardner, Ian Simon, Ethan Manilow, Curtis Hawthorne, and Jesse Engel. Mt3: Multi-task  
598 multitrack music transcription. *arXiv preprint arXiv:2111.03017*, 2021.
- 600 Jort F Gemmeke, Daniel PW Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R Channing  
601 Moore, Manoj Plakal, and Marvin Ritter. Audio set: An ontology and human-labeled dataset for  
602 audio events. In *2017 IEEE international conference on acoustics, speech and signal processing*  
603 *(ICASSP)*, pages 776–780. IEEE, 2017.
- 604 Curtis Hawthorne, Erich Elsen, Jialin Song, Adam Roberts, Ian Simon, Colin Raffel, Jesse Engel,  
605 Sageev Oore, and Douglas Eck. Onsets and frames: Dual-objective piano transcription. *arXiv*  
606 *preprint arXiv:1710.11153*, 2017.
- 608 Curtis Hawthorne, Andriy Stasyuk, Adam Roberts, Ian Simon, Cheng-Zhi Anna Huang, Sander  
609 Dieleman, Erich Elsen, Jesse Engel, and Douglas Eck. Enabling factorized piano music modeling  
610 and generation with the maestro dataset. *arXiv preprint arXiv:1810.12247*, 2018.
- 611 Curtis Hawthorne, Ian Simon, Rigel Swavely, Ethan Manilow, and Jesse Engel. Sequence-to-  
612 sequence piano transcription with transformers. *arXiv preprint arXiv:2107.09142*, 2021.
- 613 Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network, 2015.  
614 URL <https://arxiv.org/abs/1503.02531>.
- 615 IMSLP. IMSLP/Petrucci music library, 2006. URL <https://imslp.org>.
- 617 Eamonn Keogh, Selina Chu, David Hart, and Michael Pazzani. Segmenting time series: A survey  
618 and novel approach. In *Data mining in time series databases*, pages 1–21. World Scientific, 2004.
- 619 Qiuqiang Kong, Bochen Li, Jitong Chen, and Yuxuan Wang. Giantmidi-piano: A large-scale midi  
620 dataset for classical piano music. *arXiv preprint arXiv:2010.07061*, 2020.
- 622 Qiuqiang Kong, Bochen Li, Xuchen Song, Yuan Wan, and Yuxuan Wang. High-resolution piano  
623 transcription with pedals by regressing onset and offset times. *IEEE/ACM Transactions on Audio,*  
624 *Speech, and Language Processing*, 29:3707–3717, 2021.
- 625 Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436–444,  
626 2015.
- 627 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019. URL <https://arxiv.org/abs/1711.05101>.
- 628 Lukáš Samuel Marták, Patricia Hu, and Gerhard Widmer. Quantifying the corpus bias problem in  
629 automatic music transcription systems. *arXiv preprint arXiv:2408.04737*, 2024.
- 630 Colin Raffel. *Learning-Based Methods for Comparing Sequences, with Applications to Audio-to-*  
631 *MIDI Alignment and Matching*. PhD thesis, Columbia University, 2016.
- 632 Colin Raffel, Brian McFee, Eric J Humphrey, Justin Salamon, Oriol Nieto, Dawen Liang, Daniel PW  
633 Ellis, and C Colin Raffel. Mir\_eval: A transparent implementation of common mir metrics. In  
634 *ISMIR*, volume 10, page 2014, 2014.
- 635 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi  
636 Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text  
637 transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
- 638 Xavier Riley, Zixun Guo, Drew Edwards, and Simon Dixon. Gaps: A large and diverse classical  
639 guitar dataset and benchmark transcription model. *arXiv preprint arXiv:2408.08653*, 2024.
- 640 Siddharth Sigtia, Emmanouil Benetos, Nicolas Boulanger-Lewandowski, Tillman Weyde, Artur  
641 S d’Avila Garcez, and Simon Dixon. A hybrid recurrent neural network for music transcription. In  
642 *2015 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages  
643 2061–2065. IEEE, 2015.
- 644

- 648 Siddharth Sigtia, Emmanouil Benetos, and Simon Dixon. An end-to-end neural network for poly-  
649 phonic piano music transcription. *IEEE/ACM Transactions on Audio, Speech, and Language*  
650 *Processing*, 24(5):927–939, 2016.
- 651 Roman Solovyev, Alexander Stempkovskiy, and Tatiana Habruseva. Benchmarks and leaderboards  
652 for sound demixing tasks, 2023.
- 653
- 654 John Thickstun, Zaid Harchaoui, and Sham Kakade. Learning features of music from scratch. *arXiv*  
655 *preprint arXiv:1611.09827*, 2016.
- 656
- 657 John Thickstun, David Hall, Chris Donahue, and Percy Liang. Anticipatory music transformer.  
658 *arXiv preprint arXiv:2306.08620*, 2023.
- 659 Keisuke Toyama, Taketo Akama, Yukara Ikemiya, Yuhta Takida, Wei-Hsiang Liao, and Yuki Mitsu-  
660 fuji. Automatic piano transcription with hierarchical frequency-time transformer. *arXiv preprint*  
661 *arXiv:2307.04305*, 2023.
- 662
- 663 Stefan Uhlich, Giorgio Fabbro, Masato Hirano, Shusuke Takahashi, Gordon Wichern, Jonathan  
664 Le Roux, Dipam Chakraborty, Sharada Mohanty, Kai Li, Yi Luo, et al. The sound demixing  
665 challenge 2023-cinematic demixing track. *Trans. Int. Soc. Music. Inf. Retr.*, 7(1):44–62, 2024.
- 666 Geraint A Wiggins. Computer representation of music in the research environment. In *Modern*  
667 *Methods for Musicology*, pages 7–22. Routledge, 2016.
- 668
- 669 Yujia Yan, Frank Cwitkowitz, and Zhiyao Duan. Skipping the frame-level: Event-based piano  
670 transcription with neural semi-crfs. *Advances in Neural Information Processing Systems*, 34:  
671 20583–20595, 2021.
- 672 Mingliang Zeng, Xu Tan, Rui Wang, Zeqian Ju, Tao Qin, and Tie-Yan Liu. Musicbert: Symbolic  
673 music understanding with large-scale pre-training. *arXiv preprint arXiv:2106.05630*, 2021.
- 674
- 675 Huan Zhang, Jingjing Tang, Syed Rifat Mahmud Rafee, and Simon Dixon György Fazekas. Atepp:  
676 A dataset of automatically transcribed expressive piano performance. In *ISMIR 2022 Hybrid*  
677 *Conference, 2022*.
- 678 Ce Zhou, Qian Li, Chen Li, Jun Yu, Yixin Liu, Guangjing Wang, Kai Zhang, Cheng Ji, Qiben Yan,  
679 Lifang He, et al. A comprehensive survey on pretrained foundation models: A history from bert  
680 to chatgpt. *arXiv preprint arXiv:2302.09419*, 2023.
- 681 Udo Zölzer. *Digital audio signal processing*. John Wiley & Sons, 2022.
- 682
- 683
- 684
- 685
- 686
- 687
- 688
- 689
- 690
- 691
- 692
- 693
- 694
- 695
- 696
- 697
- 698
- 699
- 700
- 701

## 702 A APPENDIX

### 704 A.1 CRAWLING SYSTEM PROMPT

```

706
707 1
708 2 Analyze the YouTube video title and description to determine if it's
709 3     likely a solo piano performance. Consider the following:
710 4
711 5 1. Is the content music-related?
712 6 2. Are there explicit mentions of solo piano or pianist names?
713 7 3. Does it mention other instruments, vocalists, or non-musical elements
714 8     ?
715 9 4. Is it an educational video (tutorial, lesson) rather than a
716 10     performance?
717 11 5. If a piece name is provided, is it typically for solo piano?
718 12
719 13 Pay special attention to these factors, which suggest the content is NOT
720 14     a pure solo piano performance:
721 15
722 16 - Presence of other instruments or vocalists
723 17 - Educational content (lessons, tutorials)
724 18 - Non-piano keyboard instruments (e.g., organ, harpsichord)
725 19 - Significant narration or spoken content
726 20 - Orchestral accompaniment (e.g., piano concertos)
727 21 - Audio content beyond solo piano
728 22 - Repetitive tracks (e.g., loop videos)
729 23
730 24 The presence of any of these elements should generally result in a lower
731 25     rating.
732 26
733 27 Assign a score from 0-5 where:
734 28
735 29 5 = Certainly a solo piano performance only
736 30 - Clear indication of a solo pianist performing
737 31 - No signs of additional instruments, vocals, or non-performance
738 32     elements
739 33
740 34 4 = Very likely a solo piano performance, but not entirely certain
741 35 - Strong indications of solo piano, but some minor ambiguity
742 36 - No clear signs of additional elements, but not explicitly ruled out
743 37
744 38 3 = Possibly a solo piano performance, but with significant uncertainty
745 39 - Some indications of solo piano, but also hints of potential additional
746 40     elements
747 41 - Could be a piano-focused piece with minimal additional content
748 42
749 43 2 = Likely includes elements other than solo piano
750 44 - Clear indications of additional instruments, educational content, or
751 45     non-performance elements
752 46 - Still primarily piano-focused, but definitely not a pure solo
753 47     performance
754 48
755 49 1 = Mostly not a solo piano performance
756 50 - Significant presence of other instruments, vocals, or non-musical
757     content
758 - Piano may be present but is not the sole or main focus
759
760 0 = Definitely not a solo piano performance or not piano-related at all
761 - No indication of solo piano content
762 - Completely unrelated to piano performances
763
764 Examples:
765 "Chopin Nocturne Op. 9 No. 2 - Arthur Rubinstein" => 5
766 "The Art of Fugue - Glenn Gould (Piano)" => 5

```

```

756 51 "Bohemian Rhapsody - Piano Cover with Sheet Music" => 4
757 52 "Beethoven - Ode To Joy | VERY EASY Piano Tutorial" => 3
758 53 "Mozart Piano Concerto No. 21 - London Symphony Orchestra" => 1
759 54 "Top 10 Guitar Solos of All Time" => 0
760 55
761 56 Think step by step concisely, and then provide your score as a JSON
762    string: {"score": X}

```

## 764 A.2 METADATA EXTRACTION SYSTEM PROMPT

```

765 1 Analyze the YouTube video title and description provided within XML tags
766   . If it contains information about a solo-piano performance, extract
767   the following metadata and provide it as a JSON string:
768 2
769 3 - composer: Last name of the composer, if applicable (string, omit if
770   not present or uncertain)
771 4 - opus: Opus number (e.g., Op., BWV, K., D.), if applicable (integer,
772   omit if not present or uncertain)
773 5 - piece_number: Number or identifier within the opus, if applicable (
774   integer, omit if not present or uncertain)
775 6 - genre: Primary genre of the piece (string: "classical", "jazz", "pop",
776   "blues", "ragtime", "atonal", "rock", "soundtrack", "ambient", "
777   folk", omit if uncertain)
778 7 - form: Musical form (e.g., "sonata", "etude", "improvisation", "fantasy
779   ", etc.) (string, omit if unknown or not applicable)
780 8 - performer: Last name of the pianist or performer, if known (string,
781   omit if unknown or uncertain)
782 9 - key_signature: Key signature of the piece (string, omit if not
783   mentioned or uncertain)
784 10 - difficulty: Estimated difficulty level (string: "beginner", "
785   intermediate", "advanced", "virtuoso", omit if uncertain)
786 11 - music_period: Primary musical period (string: "classical", "romantic",
787   "baroque", "impressionist", "contemporary", "modern", omit if
788   uncertain)
789 12
790 13 Rules:
791 14 1. Omit keys and values entirely for fields not present, unknown, or
792   uncertain. Do not include empty strings or placeholder values.
793 15 2. Be cautious not to include fields unless you are reasonably certain
794   they are correct.
795 16 3. Provide opus and piece_number as integers only (e.g. don't include
796   BWV, K., S., or Op.). Omit if not clearly a number or if zero.
797 17 4. In the case of well-known pieces (e.g., Moonlight Sonata, Fantaisie-
798   Impromptu, etc.), add the opus and piece_number if you are certain,
799   even if it is not in the raw text.
800 18 5. Provide form as a single word each, using very general and well-known
801   terms.
802 19 6. Infer difficulty and period from context when possible, but omit if
803   uncertain.
804 20 7. For all strings only provide a single word in lowercase ASCII.
805 21 8. For composer and performer, use only the last name. If unsure which
806   name is the last name, omit the field.
807 22 9. Provide key_signatures using standard ASCII musical notation: Use 'b'
808   for flat, '#' for sharp, and 'm' for minor. Major keys should not
809   have a suffix. Examples: 'c', 'f#m', 'bb'.
810 23 10. Only include opus and piece_number if the video is a complete
811   performance of a single traditional opus number (Op., BWV, K., D.)
812   and its movements/variations. Omit both fields for compilations,
813   multiple works, ambiguous titles, or when using non-traditional/
814   modern catalog numbers.
815 24 11. Do not confuse piece_number with other identifiers like sonata
816   numbers (e.g., "Sonata No. 14") or separate opus numbers (e.g., in "
817   Op. 37-38", neither 37 nor 38 is a piece_number). Only use
818   piece_number when it s part of an opus and subordinate to it.

```

```

810 25
811 26 Examples:
812 27
813 28 1. Input:
814 29 <title>Chopin - Nocturne in E-flat major, Op. 9 No. 2 | Rousseau</title>
815 30 <description>Fr d ric Chopin's Nocturne in E-flat major, Op. 9, No. 2.
816 31 One of the most famous classical piano pieces from the Romantic era
817 32 . Performed by Rousseau.
818 33 #chopin #nocturne #classical #piano</description>
819 34 Output:
820 35 {
821 36   "composer": "chopin",
822 37   "opus": 9,
823 38   "piece_number": 2,
824 39   "genre": "classical",
825 40   "form": "nocturne",
826 41   "performer": "rousseau",
827 42   "key_signature": "eb",
828 43   "difficulty": "advanced",
829 44   "music_period": "romantic"
830 45 }
831 46 2. Input:
832 47 <title>Glenn Gould plays Bach Partita No.2 in C-minor (FULL)</title>
833 48 <description>1959 Studio recording DISCLAIMER: I do not own any material
834 49 shown in this video. This is for entertainment purposes ONLY.
835 50 Unlawful distribution of this material can result in bad stuff,
836 51 apparently, SO DON'T DO IT!</description>
837 52 Output:
838 53 {
839 54   "composer": "bach",
840 55   "genre": "classical",
841 56   "form": "partita",
842 57   "performer": "gould",
843 58   "key_signature": "cm",
844 59   "difficulty": "advanced",
845 60   "music_period": "baroque"
846 61 }
847 62 3. Input:
848 63 <title>Jazz Piano - Bill Evans - The Solo Sessions, Voll [ Full Album
849 64 ]</title>
850 65 <description></description>
851 66 Output:
852 67 {
853 68   "performer": "evans",
854 69   "genre": "jazz",
855 70   "music_period": "modern"
856 71 }
857 72 4. Input:
858 73 <title>Martha Argerich plays Beethoven Sonata No. 31, Op. 110</title>
859 74 <description>00:00 1. Moderato cantabile molto espressivo
860 75 06:12 2. Allegro molto
861 76 08:18 3. Adagio ma non troppo - Allegro ma non troppo
862 77 </description>
863 78 Output:
864 79 {
865 80   "composer": "beethoven",
866 81   "opus": "110",
867 82   "genre": "classical",
868 83

```



```

864 84  "form": "sonata",
865 85  "performer": "argerich",
866 86  "difficulty": "advanced",
867 87  "music_period": "classical"
868 88  }
869 89
870 90 Think step by step concisely, and then provide the metadata as a JSON
      string.

```

### A.3 TRANSCRIPTION ACCURACY

Table 7: Piano transcription note accuracy of the transcription model used for Aria-MIDI, evaluated on the MAESTRO (v3) and MAPS test sets. Results are calculated using the mir\_eval library (Raffel et al., 2014) with default settings. We compare to the model introduced in Kong et al. (2021), which was used for GiantMIDI, ATEPP, and PiJAMA.

	Note			Note /w offset			Note /w offset & vel		
	P (%)	R (%)	F1 (%)	P (%)	R (%)	F1 (%)	P (%)	R (%)	F1 (%)
<i>Chosen model</i>									
MAESTRO	98.86	96.45	97.63	91.63	89.42	90.50	90.56	88.39	89.45
MAPS	91.78	89.47	90.58	-	-	-	-	-	-
<i>Kong et al. (2021)</i>									
MAESTRO	98.82	95.53	96.82	86.51	84.21	85.33	84.97	82.72	83.82
MAPS	79.37	87.43	83.10	-	-	-	-	-	-