

Fugatto 1

Foundational Generative Audio Transformer Opus 1

Anonymous authors

Paper under double-blind review

ABSTRACT

Fugatto is a versatile audio synthesis and transformation model capable of following free-form text instructions with optional audio inputs. While large language models (LLMs) trained with text on a simple next-token prediction objective can learn to infer instructions directly from the data, models trained solely on audio data lack this capacity. This is because audio data does not inherently contain the instructions that were used to generate it. To overcome this challenge, we introduce a specialized dataset generation approach optimized for producing a wide range of audio generation and transformation tasks, ensuring the data reveals meaningful relationships between audio and language. Another challenge lies in achieving compositional abilities – such as combining, interpolating between, or negating instructions – using data alone. To address it, we propose *ComposableART*, an inference-time technique that extends classifier-free guidance to compositional guidance. It enables the seamless and flexible composition of instructions, leading to highly customizable audio outputs outside the training distribution. Our evaluations across a diverse set of tasks demonstrate that *Fugatto* performs competitively with specialized models, while *ComposableART* enhances its sonic palette and control over synthesis. Most notably, we highlight our framework’s ability to synthesize emergent sounds – sonic phenomena that transcend conventional audio generation – unlocking new creative possibilities. Demo Website.

1 INTRODUCTION

Recent research has sparked a strong debate between **specialist** and **generalist** models. While specialist models excel at specific tasks, they tend to be brittle, often struggling with changes in data distribution or task requirements. In contrast, generalist models eliminate the need for task-specific designs, can process diverse data, and scale effectively with increased compute and data. They also demonstrate emergent capabilities, enabling unsupervised task learning by leveraging broader datasets (Radford et al., 2019) and demonstrations (Brown et al., 2020). In this paper, we propose a strategy for developing a generalist audio synthesis and transformation model, called *Fugatto*, and an inference method for composing instructions through latent space manipulation, including from different models, called Composable Audio Representation Transformation (*ComposableART*).

Table 1: Comparison of our proposed model *Fugatto* with other models.

| | <i>Fugatto</i> | AudioBox | NExT-GPT | UniAudio | AUDIT | VoiceLDM |
|-------------------------|----------------|----------|----------|----------|-------|----------|
| Emergent properties | ✓ | ? | ? | ? | ? | ? |
| Large-scale data | ✓ | ✗ | ✓ | ✗ | ✗ | ✗ |
| Supports numerous tasks | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ |
| Free-Form Instructions | ✓ | ✗ | ? | ? | ✓ | ✗ |
| Open-ended generation | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ |
| Compositionality | ✓ | ✗ | ✗ | ✗ | ✗ | ✓ |
| Multi-Modal Inputs | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ |

Large language models (LLMs) have demonstrated impressive unsupervised multitask learning capabilities in the text domain (Radford et al., 2019), where instructions can be inferred from the

054 data itself. However, such instructions are typically absent in the audio domain, making it difficult to
055 generalize to unseen tasks without explicit guidance. Although models such as (Wang et al., 2024;
056 Yang et al., 2023; Vyas et al., 2023) exist, they have several limitations enumerated in Table 1. In
057 this panorama, dataset and instruction generation is necessary.

058 We employ a **multifaceted data and instruction generation strategy** that considerably expands the
059 range of tasks of audio generation models. First, we use LLMs to generate and augment instructions
060 and captions, providing *Fugatto* with instructions that are closer to free-form instructions (Goel
061 et al., 2024; Doh et al., 2023). Second, we develop instructions that can be either absolute (e.g.,
062 ”synthesize a happy voice”) or relative (e.g., ”increase the happiness of this voice”), enabling *Fugatto*
063 to handle a wide array of dynamic tasks (OpenAI, 2024). Third, we leverage audio understanding
064 models (Kong et al., 2024a; Gong et al., 2023) to create descriptions and synthetic captions for
065 audio clips, enriching the dataset where annotations are sparse, allowing for better generalization
066 and more accurate performance (Kong et al., 2024b). Fourth, we transmute existing datasets to
067 uncover new relationships, enabling the creation of entirely new tasks without requiring additional
068 raw data. Finally, we use audio processing tools to create new connections between text, audio, and
069 their corresponding transformations.

070 By combining these approaches, we ensure that *Fugatto* has access to diverse and enriched datasets,
071 allowing it to learn from various audio domains and contexts. This strategy enhances the model’s task
072 diversity and lays the groundwork for unsupervised multitask learning at scale, unlocking emergent
073 abilities such as synthesizing entirely new sounds, such as “a saxophone barking”.

074 While data is important, the ability to **compose, interpolate between, and negate instructions** is
075 generally difficult to obtain through data alone. Negation data, for instance, is typically unavailable,
076 and generating outputs that represent the composition or interpolation of instructions is equally chal-
077 lenging. Though this has been explored in the vision domain using Energy Based Models (EBMs) (Du
078 et al., 2020) and EBMs in latent space (Nie et al., 2021), EBMs require training a new model per
079 attribute, which can be cumbersome and impractical for a large number of attributes.

080 To overcome this limitation, we propose an inference-time technique called *ComposableART*, which
081 is based on classifier-free guidance (CFG) (Ho & Salimans, 2022) and further expands on the dual
082 classifier-free guidance in (Yang et al., 2024; Lee et al., 2024). We propose a generalized framework
083 that leverages the weighted combination of vector fields between instructions, frame indices and
084 models. This approach enables *Fugatto* to handle complex instruction-based operations, such as
085 smoothly interpolating between instructions or negating specific instructions to exclude undesired
086 effects. In contrast, models like (Vyas et al., 2023; Yang et al., 2023) rely on more rigid methods,
087 often requiring external classifiers or manual intervention to achieve similar results.

088 In this paper, we present a detailed exploration of *Fugatto*’s dataset and instruction generation,
089 training strategies, implementation improvements, and performance across a wide range of tasks.
090 Through rigorous evaluations and comparisons with specialist and generalist models, we establish
091 *Fugatto* as a new benchmark for foundation models for audio synthesis and transformation. Similarly,
092 we establish *ComposableART* as a highly desirable framework for compositional guidance that
093 unlocks the full potential of score-based generative models. Our contributions include:

- 094 • We offer a comprehensive strategy for building a foundation model for audio generation and
095 transformation given text and audio inputs, delivering strong performance across a wide
096 range of tasks and providing a robust framework for both research and practical applications.
- 097 • We demonstrate how to enhance and create contextually rich audio and text datasets while
098 generating flexible instructions with LLMs, enabling our community to replicate and adapt
099 these techniques for their own models.
- 100 • We demonstrate how to perform composition, interpolation, and negation of instructions by
101 extending CFG to compositional guidance, enabling better control over the model’s outputs.

102 2 APPROACH

103
104 Our approach to audio generation given text and optional audio is similar to recent approaches in
105 LLMs, focusing on large scale compute and datasets, followed by pre-training and fine-tuning stages.
106 However, our approach differs in two aspects. First, our dataset generation mechanism requires
107

108 going beyond unsupervised next token prediction (Section 2.1). Second, we propose inference time
109 techniques to control audio generation (Section 2.4). In Appendix A.1 we provide a detailed list of
110 tasks and instructions our model supports.

111 112 113 114 2.1 DATASET GENERATION

115
116 We focus on building as large and diverse a dataset in order to collect demonstrations of as many
117 audio tasks and contexts as possible. We emphasize that our ultimate goal is not to just excel on such
118 tasks, but to drive, as a community, towards a future where unsupervised multitask learning emerges
119 from data and model scale. Towards this goal, we propose a dataset generation strategy built on five
120 pillars below, and provide, in Appendix A.1.1, a thorough description of each pillar.

121 **I – Generating Free-Form Instructions with LLMs:** Our strategy consists of prompting an
122 LLM to create programmatic task-specific instruction generators, similar to Kocielnik et al. (2023).
123 We prompt an LLM to create python methods that generate instructions of different lengths and
124 for different personas (standard, young-crowd, thirty-somethings, professional), given task-specific
125 inputs including audio description, language, and others. Each persona has a set of verbs, adverbs,
126 connectors, and other assets to create the instructions.

127 **II – Generating Absolute and Relative Instructions:** Following GPT4-o’s (OpenAI, 2024) ability
128 to perform a relative change in speech, we aim to obtain an audio generation model that is able to
129 follow instructions that are absolute or relative, such as ”synthesize a happy voice” or ”increase the
130 happiness in this voice.”. Given that such data is normally not available, we later describe a strategy
131 to create such data by transmuting existing datasets and leveraging audio processing tools. Once
132 such datasets are created, we can select one sample and create an absolute instruction for it, or select
133 two samples and create a relative instruction. Similarly to absolute instructions, we generate relative
134 instructions with a python method, produced by an LLM, that creates an instruction given the task,
135 the attribute being modified and how it should be modified.

136 **III – Generating Audio Descriptions with Audio Understanding Models:** We use audio under-
137 standing and classification models (Kong et al., 2024a) to produce synthetic captions for audio in
138 the wild, following recent research (Leng et al., 2023; Kong et al., 2024b) showing that it is possible
139 to drastically expand or improve text-to-audio models with synthetic captions. As such, we expand
140 the strategies described in (Leng et al., 2023; Kong et al., 2024b) to generate high quality synthetic
141 captions. For speech data, we implemented a prompt generation pipeline that automates the creation
142 of natural language descriptions for voices. The pipeline converts speech attributes predicted by
143 models – such as ”gender”, emotion, and speech quality – into detailed natural language descriptions
144 using LLM-generated templates. These templates describe voices in various ways based on the
145 speaker attributes, enhancing diversity by generating descriptions in multiple formats.

146 **IV – Creating New Tasks and Datasets by Transmuting Datasets:** We leverage implicit relation-
147 ships between samples in a dataset to enable new tasks. Generally speaking, we look for datasets
148 where one factor is held constant while other factors change. For example, we leverage emotional
149 speech synthesis datasets with different renditions of the same text (Livingstone & Russo, 2018) by
150 the same speaker to define a speech transformation task. Similarly, we leverage instrument synthesis
151 datasets with different renditions of the same note (Engel et al., 2017) to define an instrument trans-
152 formation task. We also leverage datasets that provide the individual parts of a sounds mixture (Rafii
153 et al., 2017) to support tasks such as source separation, and audio generation conditioned on audio
154 context and captions, possibly synthetic.

155 **V – Creating New Tasks and Datasets by Leveraging Audio Processing Tools:** We create
156 synthetic paired data for speech and audio by using Praat (Boersma & Van Heuven, 2001) and
157 Pedalboard (Spotify, 2024) to manipulate several speech and audio factors. For each factor, we apply
158 controlled modifications, allowing us to generate speech and audio samples with specific alterations.
159 With this strategy we can create speech and audio pairs that describe speech and audio transformation
160 tasks such as changing of one or multiple speech factors, for example ”increase the F0 variance
161 slightly and decrease the speech rate.” or ”add moderate reverb to this audio file”. For each factor, we
determine a practical range of adjustments and define increments that correspond to varying degrees
of change, such as ”slightly”, ”moderate”, and ”significant”.

162 With these pillars established and leveraging open source datasets, we are able to build a large text
 163 and audio dataset with at least 20 million rows, not including on-the-fly modifications to captions,
 164 instructions and audio. Assuming each row refers to 10 seconds of audio, our dataset is comprised of
 165 at least 2.9 million hours of audio, approximately 330 years. We emphasize that this compilation is
 166 exclusively comprised of open source data, and provide a full list of datasets, tasks, and instructions
 167 in Appendices A.1.2, A.1.3 and A.1.4 respectively.

168 2.2 INSTRUCTION GENERATION

169 Our approach supports template-based and free-form instructions.

170
 171 **Template-based Instructions:** In these instructions, the task is explicitly provided, followed by
 172 task-specific attributes, always in the same order, and with each attribute wrapped between start and
 173 end of attribute markers. We dynamically construct template-based instructions based on the task
 174 and the set of factors. Each factor is specified by its name and corresponding value, using the format
 175 `given {name}:{value}` followed by a closing HTML-like tag `</{factor}>`. The instruction
 176 always starts with the task, followed by its factors, and ends with `output : .` Additionally, a `given`
 177 `context` clause is appended at different locations, determined by the task at hand, when audio
 178 contexts are present. The structure of a template-based instruction is:

```
179 input:{task} given {factor}:{value}</{factor}>given context:<audio>output:
```

180 where `{task}` represents the specific task, `{factor}` and `{value}` correspond to the different
 181 factors and their respective values, and `<audio>` refers to the audio context provided. We provide
 182 examples of task and dataset specific instructions in Appendix A.1.4.

183
 184 **Free-form Instructions:** We dynamically construct free-form instructions by using instruction
 185 generators introduced in Section 2.1. Unlike template-based instructions where we know before
 186 hand the reference in text to each audio context, in free-form it is not straight forward to determine
 187 which word in the text refers to the audio. As such, in free-form instructions we simply append to the
 188 instruction `given context_k:<audio>` for each audio context, resulting in this structure:

```
189 input:{instruction} [given context:<audio_k>]output:
```

190
 191 In order to promote simplicity and agility during development, we decided to use raw text instead
 192 of learnable tokens for `given factor` and `</factor>`. Following LLM practice, the full
 193 instruction and each audio are wrapped with learnable `<start of>` and `<end of>` tokens.

194 2.3 MODEL AND TRAINING

195
 196 In this section we describe the text and audio representations used in our model, as well as the
 197 training objective and architecture. We provide a graphical depiction and details for hyperparameters,
 198 objective function, training stages and oversampling in A.2.

199
 200 **Text and Audio representation:** The text representation is obtained by encoding the previously
 201 described instructions with a pre-trained language model held frozen during training. In this *Opus*,
 202 we use the *byT5* tokenizer free language model (Xue et al., 2022), which supports a large set of
 203 characters, including IPA. In this opus, the audio representation is a 22khz mel-spectrogram with 80
 204 bins, which is subsumed by a relatively shallow learnable transformer encoder. The mel spectrogram
 205 is scaled to have approximately 0 mean and 1 standard deviation.

206
 207 **Training Objective and Architecture:** We train our model with the Optimal Transport Conditional
 208 Flow Matching (OT-CFM) objective (Lipman et al., 2022; Tong et al., 2023), and use a T5-based (Raf-
 209 fel et al., 2020) Transformer (Vaswani, 2017) with Adaptive Layer Norm (Xu et al., 2019) as the
 210 parameterization of the vector field estimator. We replace the Transformer MLPs with kernel size
 211 3 convolutions. After independently projecting the encoded text and audio to a shared embedding
 212 space, we time-wise concatenate the embedded audio and text. The model cross-attends to this
 213 representation, applying Adaptive Layer Norm (Xu et al., 2019) to them on every layer.

214 We observed that certain implementation choices yielded better training curves. Specifically, adaptive
 215 layer norm is completely computed in FP32, GELU uses approximate *tanh*, and the final layers are
 initialized to outputs zeros, which is approximately the mean of our scaled mel-distribution.

Training Stages: *Fugatto* training follows curriculum learning¹ (Bengio et al., 2009). We start with template-based instructions and a subset of tasks. Eventually, once we informally establish that the model is able to follow template-based instructions by observing validation scores and listening to samples, we proceed with an equal mixture of template-based and free-form instructions. Given that some tasks are underrepresented in the data, we find that oversampling leads to better results. Empirically, we find that sampling from a multinomial distribution with upsampling parameter $\beta = 0.25$, similar to Le et al., 2024, is sufficient. As training progresses, we adjust each dataset’s weight according to validation scores on target tasks.

2.4 COMPOSABLE AUDIO REPRESENTATION TRANSFORMATION (*ComposableART*)

We extend Classifier Free Guidance(CFG) (Ho & Salimans, 2022) to support compositional guidance. Compositional guidance provides an OT-CFM model with the ability to independently control and generate (unseen) combinations of instructions and tasks, including with vector fields from different models (Karras et al., 2024). Compositional generation has been explored in Diffusion Models (Liu et al., 2023; Yang et al., 2024; Lee et al., 2024), for images and audio. To the best of our knowledge, we are the first to showcase novel ways of applying compositional guidance, expanding it to not just attributes, but also tasks, models and temporal composition of attributes.

Compositional Guidance Method Classifier Free Guidance(CFG), generally applied to diffusion models, combines the conditional and unconditional score estimates to obtain samples of higher quality and diversity pertaining to the condition. The following equation summarizes the application of CFG, where ϵ_θ represents the score estimate and γ represents the gradient scaling factor:

$$\epsilon_\theta(\mathbf{z}_\lambda, \mathbf{c}) = \epsilon_\theta(\mathbf{z}_\lambda, \mathbf{c}) + \gamma(\epsilon_\theta(\mathbf{z}_\lambda, \mathbf{c}) - \epsilon_\theta(\mathbf{z}_\lambda)) \quad (1)$$

We extend the Classifier Free Guidance framework to support the combination of vector fields across multiple instructions (c_k), multiple mel-frame indices (f) and multiple models (θ_m). Let $v_{t,f}(c_k, \theta_m)$ be the vector field produced at flow-step t by model m , parameterized by θ_m , given condition c_k or \emptyset , for mel-frame f . Additionally, let $w_{k,f,m}$ be the flow-step invariant and user-determined scalar weight associated with $v_{t,f}(c_k, \theta_m)$. The equation for compositional guidance across instructions, frames, and models is defined as:

$$\tilde{v}_{t,f} = \sum_{k,m} w_{k,f,m} (v_{t,f}(c_k, \theta_m) + \gamma(v_{t,f}(c_k, \theta_m) - v_{t,f}(\emptyset, \theta_m))) \quad (2)$$

where γ refers to the gradient scale parameter from CFG, and $\tilde{v}_{t,f}$ is the resulting composed vector field for flow step t and frame f . This follows similar conditional independence assumptions in (Nie et al., 2021; Liu et al., 2023) to support compositional guidance. We apply this compositional guidance at every step of the flow-matching inference procedure.

Attribute/Event Composition: An attribute is a simple input prompt belonging to a particular task like speech synthesis. A task like audio event generation can have multiple prompts or attributes as input. With compositional guidance, we can support unique unseen combinations of attributes. This gives the users an ability to create artistic combinations like simulating a scene with multiple-audio events, e.g. by composing ‘thunder’, ‘rain’ and ‘wind’ a storm can be achieved.

Task Composition: There are many tasks the model has seen during training, but the model hasn’t seen a combination of tasks (like speech synthesis alongside audio event generation) during training. With compositional guidance, we can support synthesis of unique unseen combinations of tasks like speech synthesis with a particular audio event in the background.

Model Composition: We also showcase that the same technique can be extended to compose multiple different models together. This can be especially useful if we train multiple different domain specific versions of *Fugatto*, each having its own parameters and datasets. This enables users to synthesize a “mixture of experts” sample that uses multiple different models, each trained on independent domains like speech and audio events. Similar to the approach in (Karras et al.,

¹It simplifies development and supports incremental research, though not strictly necessary.

270 2024), we modify the above setup simply by using the velocities for each independent model with
 271 parameters θ_m .

272 **Temporal Composition:** Instead of broadcasting the same scalar for each frame f , we instead use
 273 a different weight $w_{k,f,m}$ per frame. This allows users to control the compositional output with any
 274 arbitrary curve (like sigmoid or linearly increasing or decreasing weight curve) while enjoying the
 275 benefits of composition of attributes, tasks and models.
 276

277 3 EXPERIMENTS

278 We present a comprehensive evaluation of *Fugatto* across multiple tasks to demonstrate its effective-
 279 ness and versatility. We begin with an ablation study, examining the impact of various design choices.
 280 Next, we evaluate *Fugatto*'s performance in audio synthesis and transformation tasks in speech, music
 281 and general sounds. Finally, we explore *Fugatto*'s emergent capabilities, and perform a thorough
 282 evaluation of our *ComposableART* method. Unless otherwise specified, we use template-based
 283 instructions, 2nd order Heun solver with 50 function evaluations, and task specific gradient scale γ .
 284
 285

286 3.1 ABLATIONS

287 We first analyze the impact of different t -sampling strategies in OT-CFM, comparing the traditional
 288 uniform sampling with others. Then, we examine the effect of model size on both loss metrics
 289 and emergent capabilities, including the ability to synthesize novel sounds not found in the training
 290 data—such as a “saxophone meowing” or “a person speaking while barking.”
 291

292 **t -sampling strategy:** In the OT-CFM framework, the timestep t is typically sampled from a uni-
 293 form distribution, $t \sim \mathcal{U}(0, 1)$. However, recent discussions within the community have introduced
 294 conflicting strategies for this sampling process. Notably, Stable Audio’s GitHub repository proposes
 295 sampling t from a sigmoid-transformed normal distribution, $t \sim \text{sigmoid}(\mathcal{N}(0, 1))$, thereby concen-
 296 trating samples around $t = 0.5$. On the other hand, (Lovelace et al., 2023) advocate for increased
 297 sampling from values of t closer to zero.
 298

299 In our experiments, we observe that although Stable Audio’s strategy provides a marginal improve-
 300 ment on TTA tasks, it renders the model unable to effectively attend to the transcript in text-to-speech
 301 (TTS) tasks – a critical requirement. We provide an explanation for this phenomenon in Appendix A.3.
 302

303 Training with $t \sim \mathcal{U}(0, 1)$ is effective across all tasks.
 304 Training with $t \sim \text{sigmoid}(\mathcal{N}(0, 1))$ significantly degrades TTS performance.

305 **Model capacity:** We evaluate how the learnable parameter count influences loss curves and emergent
 306 capabilities. We consider models with 0.8 B, 1.4 B params, and 2.5 B parameters. Under a fixed
 307 data composition and sampling weights, we observe that increasing the parameter count from the
 308 smallest to the largest model not only improves validation losses but also delays overfitting. In
 309 Appendix A.2, we provide task-specific validation loss plots that showcase the expected decrease
 310 in validation loss as parameter count increases. Informally and consistent with findings in (Radford
 311 et al., 2019), we observe that the smaller model does not exhibit emergent abilities comparable to the
 312 larger models, particularly in their ability to synthesize novel sounds absent from the training data,
 313 such as “saxophone barking”. We invite readers to evaluate samples in our supplementary materials.
 314

315 Emerging capabilities surface with sufficient model capacity and training data.

316 3.2 AUDIO SYNTHESIS

317 **Text-To-Voice Synthesis:** We evaluate in-context text-to-speech synthesis (TTS) and singing voice
 318 synthesis (SVS). For TTS, we follow the evaluation in Wang et al. (2023), using the same transcripts
 319 as Eskimez et al. (2024), to evaluate our model’s ability to perform speech synthesis given a transcript
 320 and a speech sample from an unseen speaker. Following our training strategy, during evaluation we
 321 always provide the speaker’s previous sentence when possible, otherwise a random sample different
 322 from the target. For SVS, we evaluate *Fugatto*'s ability to generate singing voice samples from
 323 instructions describing the desired lyrics and musical style without a backing track, for example:

“Showcases a female singer with an interesting sound, conveys the message through american english lyrics, and infuses country influences throughout.” The full list is available in Appendix A.4

Our zero-shot TTS results in Table 2a show that *Fugatto* has word error rates similar to ground truth, and is competitive with expert and generalist (Omni) models in terms of speaker similarity. Our SVS results in Table 2b shows high cosine similarity between CLAP embeddings of the synthetic samples and the captions used to create it. We note that the higher WER in SVS likely comes from the higher difficulty for the generative model and the speech transcription model.

Table 2: *Fugatto* is comparable to generalist models and expert models on in-context TTS benchmarks. On SVS, it synthesizes samples with high CLAP-similarity and low WER relative to the task at hand.

| (a) TTS on the LibriSpeech Test Clean benchmark in Wang et al. (2023) | | | | | (b) SVS: WER and CLAP-Scores on a set of 10 music styles and 13 lyrics snippets from famous songs. | | |
|---|------|-------|---------|---------|--|-------|--------|
| Model | Omni | WER ↓ | SIM-o ↑ | SIM-r ↑ | Model | WER ↓ | CLAP ↑ |
| Ground Truth | | 2.20 | | | | | |
| Vall-E (24khz) | ✗ | 5.90 | | 0.58 | <i>Fugatto</i> $\gamma = 2$ | 71.90 | 0.49 |
| Natural Speech 3 (16khz) | ✗ | 1.81 | 0.67 | 0.76 | <i>Fugatto</i> $\gamma = 4$ | 19.54 | 0.45 |
| AudioBox (16khz) | ✗ | 4.80 | 0.73 | | | | |
| UniAudio (16khz) | ✓ | 2.00 | | 0.71 | | | |
| <i>Fugatto</i> $\gamma = 2$ | ✓ | 2.66 | 0.60 | 0.61 | | | |
| <i>Fugatto</i> $\gamma = 3$ | ✓ | 2.44 | 0.61 | 0.62 | | | |

Text-To-Audio (TTA): We showcase *Fugatto*’s performance on traditional TTA benchmarks that measure a model’s ability to synthesize general sounds (AudioCAPS) and music (MusicCAPS) that follow instructions provided in text. We use the metrics (FD, FAD, and IS) and data splits (train, test) used in Kong et al. (2024b). Results in Table 3a and Table 3b shows that our model achieves strictly better scores than existing generalist models, while occasionally outperforming expert models.

Table 3: *Fugatto* outperforms generalists models and occasionally outperforms specialist models in TTA benchmarks on AudioCaps and MusicCaps.

| (a) TTA on the AudioCaps benchmark in Kong et al., 2024b | | | | | (b) TTA on the MusicCaps benchmark in Kong et al., 2024b | | | | |
|--|------|-------|-------|-------|--|------|-------|-------|------|
| Model | Omni | FD ↓ | FAD ↓ | IS ↑ | Model | Omni | FD ↓ | FAD ↓ | IS ↑ |
| VoiceLDM-Maudio | ✗ | | 2.50 | | MusicGen (medium) | ✗ | 35.52 | 5.02 | 1.94 |
| AudioBox (16khz) | ✗ | 10.14 | 1.10 | 11.90 | AudioLDM-2-large | ✗ | 16.12 | 2.74 | 2.30 |
| NExT-GPT | ✓ | | 1.68 | | Tango-AF&AC-FT-MC | ✗ | 21.84 | 1.99 | 2.21 |
| UniAudio | ✓ | | 3.12 | | UniAudio* | ✓ | | 3.65 | |
| <i>Fugatto</i> $\gamma = 1$ | ✓ | 16.73 | 1.36 | 9.72 | <i>Fugatto</i> $\gamma = 1$ | ✓ | 11.52 | 1.43 | 2.73 |
| <i>Fugatto</i> $\gamma = 2$ | ✓ | 20.20 | 2.21 | 10.21 | <i>Fugatto</i> $\gamma = 2$ | ✓ | 13.18 | 1.93 | 2.97 |

3.3 AUDIO TRANSFORMATIONS

Speech Enhancement: We evaluate speech denoising and bandwidth extension tasks. Speech denoising evaluates the ability to extract speech from an additive mixture comprised of speech and noise. Bandwidth extension (sometimes referred to as “upsampling”) evaluates the ability to recreate missing content from audio that is low passed filtered and downsampled not to include frequencies above a certain threshold frequency². For speech denoising, we use the DNS-Noisy benchmark and traditional metrics PESQ and STOI described in Kong et al., 2023. For bandwidth extension, we use the VCTK benchmark and the traditional metric LSD described in Liu et al., 2024. In *Fugatto*, this task is interpreted as source separation, in contrast with the enhancement task that modifies the acoustic qualities of the target audio. Though *Fugatto* is comparable to specialist models in bandwidth extension, work remains to be done to close the gap in denoising.

Speech Modulation: In this task we evaluate our model’s ability to transform a person’s emotion in speech into another emotion, while preserving their speaker identity and the transcript. For this purpose, construct a train and test set based on the ESD dataset. We use open-source models to report emotion classification and correlation with Valence, Arousal and Dominance (VAD). We establish

²We use librosa after observing that torch.audio leaks frequencies above the threshold frequency.

Table 4: *Fugatto* is comparable to specialist models for speech denoising and upsampling.

| (a) Speech Denoising on the DNS benchmark in Kong et al., 2023 | | | | | (b) Upsampling on the VCTK benchmark in Liu et al., 2024 | | | |
|--|------|----------------------|----------------------|--------|--|------|------------|------------|
| Model | Omni | PESQ _{WB} ↑ | PESQ _{NB} ↑ | STOI ↑ | Model | Omni | LSD 4kHz ↓ | LSD 8kHz ↓ |
| Noisy dataset | | 1.59 | 2.16 | 91.60 | Unprocessed (to 22kHz) | | 2.74 | 1.84 |
| FullSubNet | ✗ | 2.90 | 3.37 | 96.40 | NuWave (to 22kHz) | ✗ | 1.37 | 0.88 |
| FAIR-Denoiser | ✗ | 2.66 | 3.23 | 96.60 | NVSR (to 22kHz) | ✗ | 1.49 | 1.37 |
| CleanUNet 2 | ✗ | 3.26 | 3.66 | 97.70 | AudioSR (to 22kHz) | ✗ | 1.25 | 1.08 |
| <i>Fugatto</i> $\gamma = 0.1$ | ✓ | 2.77 | 3.32 | 95.70 | <i>Fugatto</i> $\gamma = 0.1$ | ✓ | 1.29 | 1.25 |
| <i>Fugatto</i> $\gamma = 1$ | ✓ | 2.73 | 3.34 | 95.90 | <i>Fugatto</i> $\gamma = 1$ | ✓ | 1.38 | 1.34 |

upper bounds by also computing scores on ground truth data. Our results in Table 5a show that our model is able to properly transform the emotion as well as the ground truth, it needs improvement on speaker similarity and word error rates.

MIDI2Audio: We evaluate our model’s ability to convert midi to audio with the 250 simple 2-bar monophonic melodies from Pati et al., 2020. Note that *Fugatto* has never seen monophonic melodies during training, with the average number of stems present in training being 8. We use the error in pitch estimation on the ground truth MIDI notes as the upper-bound quality. We evaluate the L1 distance of the F0 contours extracted from the input rendered MIDI and the generated audio. During evaluation, we transpose the pitch contours to a common key. We provide examples in Appendix A.5.

Table 5: *Fugatto* performs well on novel tasks such as emotion conversion and MIDI2Audio.

(a) Speech modulation (Emotion Conversion): remarkably high Top-2 accuracy and pearson correlation ρ between ground truth and synthetic samples on VAD. Low speaker similarity.

| Model | WER | SIM-o | ρ_{VAD} | Top-2 Acc. |
|-----------------------------|-------|-------|---------------------|------------|
| Ground Truth | 7.42 | | | 0.62 |
| <i>Fugatto</i> $\gamma = 2$ | 24.17 | 0.19 | 0.68, 0.77, 0.77 | 0.59 |
| <i>Fugatto</i> $\gamma = 3$ | 21.96 | 0.21 | 0.68, 0.77, 0.77 | 0.62 |

(b) MIDI2Audio: surprising zero-shot abilities on monophonic melody to audio.

| Model | L1 ↓ |
|-----------------------------|------|
| F0 estimation error | 0.21 |
| <i>Fugatto</i> $\gamma = 1$ | 1.74 |
| <i>Fugatto</i> $\gamma = 2$ | 1.00 |

3.4 SOUND GALLERY AND EMERGENT CAPABILITIES

In light of *Fugatto*’s artistic capabilities, we invite readers to a guided tour through our sound gallery. Our goal is to provide a proof of existence that a text-to-audio model is able to execute tasks and follow instruction not seen during training, e.g. “A person barking the words: ‘Who let the dogs out?’” or converting monophonic midi to singing voice even though the model was only trained on multi-track MIDI to audio conversion. In line with Radford et al., 2019, we show that such abilities emerge with enough model capacity and data diversity; thus, we call them *emergent abilities*.

3.5 COMPOSITIONALITY

Compositionality enables users to combine attributes and tasks to generate novel input combinations not found in the training data, providing artistic and fine-grained control over the desired output. Since such novel combination rarely exists in the training data or the natural world, evaluating these samples is typically challenging. Furthermore, there are no established baselines and metrics for such sounds. Despite these challenges, we aim to provide a qualitative and quantitative evaluation of our proposed approach and invite readers to listen to compositional samples on our website.

3.5.1 ATTRIBUTE/EVENT COMPOSITION

Control intensity of each instruction (Weighted Combination): Compositional synthesis with instructions gives users a knob to control the intensity of each instruction. In order to evaluate this ability, we create $\binom{10}{2}$ pairs of instructions by leveraging 10 event labels provided in Appendix A.6. Given a pair of events we can generate composite instructions through language or *ComposableART*:

Baseline linguistic composition example input:

input: synthesize Event 1 **and** Event 2

Proposed *ComposableART* composition example inputs:

$w_1 * v(\text{input: synthesize Event 1}) + w_2 * v(\text{input: synthesize Event 2})$

For each pair of events, we first generate samples using *ComposableART* with different convex combination of weights for each instruction, and using the linguistic baseline³. Then, for each method and generated audio, we compute the cosine similarity between the audio and text CLAP embeddings for each event in the event pair used to produce the audio, shown as Instruction 1 and 2 in Figure 1.

Figure 1 shows that as we increase the weight on Instruction 1, the occurrence of the event in the generated clip increases as evidenced by the CLAP scores. Reciprocally, as we decrease the weight on Instruction 2, the occurrence of the event in the generated clip decreases. In the linguistic baseline, we cannot control the weights of each attribute and their independent existence in synthesized samples is lower than the *ComposableART* samples at higher weights.

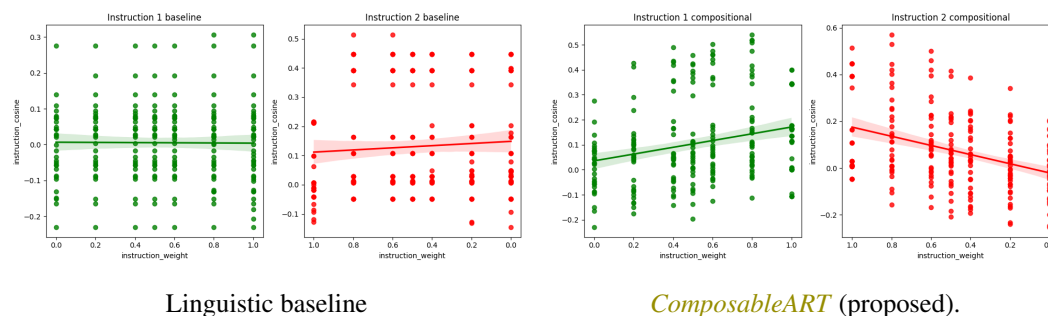


Figure 1: Comparison of CLAP scores between the Linguistic baseline and *ComposableART*'s composition of attributes with *Fugatto*. Instruction is equivalent to event.

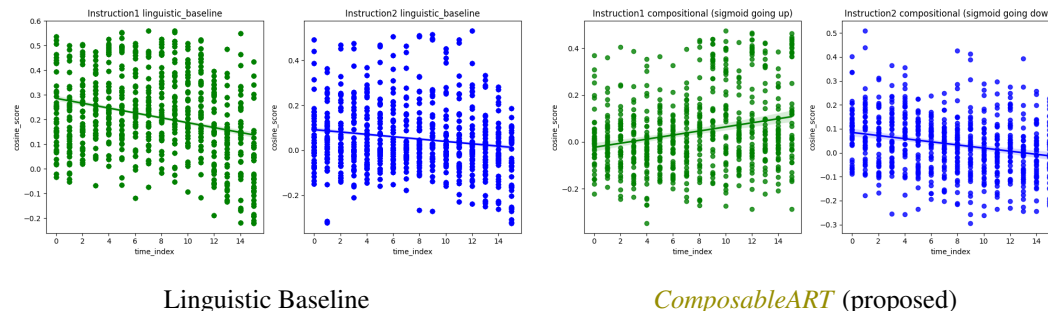


Figure 2: Comparison of Baseline and *ComposableART* Temporal guidance for instruction sequences.

Negation of attributes: With our approach, users can assign negative weights to attributes, producing negative velocity that steers the model away from the attribute. We use the previous $\binom{10}{2}$ event pairs to evaluate our approach alongside linguistic negation with the keyword 'NOT':

Baseline Linguistic Negation example input:

input: synthesize Event 1 and **NOT** Event 2

ComposableART Negation example input:

$v(\text{input: synthesize Event 1}) - v(\text{input: synthesize Event 2})$

Table 6 shows CLAP cosine similarity of the linguistic baseline against the proposed *ComposableART* method. It can be observed that while the positive event remains similar to the baseline, the negative event has a considerably lower cosine similarity than the baseline, indicating the effectiveness in removal of the audio event using *ComposableART* as compared to the linguistic approach. Qualita-

³In Figure 1, the diagonal fit, different weights, and the lack of samples in the linguistic baseline are an implementation and compute timeout byproduct given that the linguistic baseline does not support weights.

tively, we also observe that this can be immensely useful in steering towards negative emotions in speech synthesis or swapping gender, a task which is not trivial.

Interpolation of attributes: *ComposableART* also supports interpolation of attributes by having the ability to independently control one attribute while keeping others. We evaluate our ability to control “pitch” as the interpolation variable by changing the weights on pitch and keeping “text” and “language” attributes fixed. Figure 3 showcases the expected decrease in fundamental frequency as we increase the weight on the “low pitch” attribute.

Table 6: CLAP Cosine Similarity of Attributes

| Instruction Type | CLAP Score |
|--------------------------------------|------------------|
| Linguistic baseline +ve Instruction | 0.02 ± 0.02 |
| Linguistic baseline -ve Instruction | 0.17 ± 0.02 |
| <i>ComposableART</i> +ve Instruction | 0.06 ± 0.01 |
| <i>ComposableART</i> -ve Instruction | -0.04 ± 0.02 |

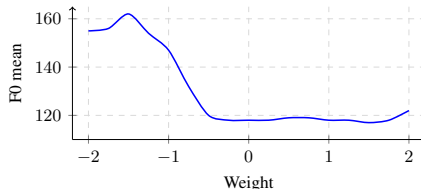


Figure 3: F0 mean given weight on the “low pitch” condition

3.5.2 TASK COMPOSITION

In our website, we provide qualitative samples where we compose a set of tasks involving “electronic music”, “birds chirping”, “dog barking”, and “TTS”. Our results show that the output conforms to the composition of such tasks using the proposed method.

3.5.3 MODEL COMPOSITION

In our website, we provide samples where we consider 2 different *Fugatto* models, one trained on speech datasets and the other trained on general sounds and audio events. We perform model composition to synthesize samples that contain both speech as well as audio events. This can be immensely useful in the future, where each domain specific model is an expert model and a combination of such high-quality experts can be used to synthesize compositional outputs without the need to train a monolithic large generative model.

3.5.4 TEMPORAL COMPOSITION

To evaluate the effectiveness of temporal guidance, we combine the previous set of audio events pairs over time. For *ComposableART*, we use a *sigmoid*-like curve to increase and decrease over time the weights on event 1 and event 2 respectively. For the equivalent linguistic baseline, we create the instruction “Event 2 followed by Event 1”. Figure 2 showcases time-windowed CLAP scores, y-axis, for each approach and event across time, x-axis. We observe that the baseline is unable to establish the temporal trend as well as *ComposableART*, where we observe, as expected, a consistent increase in event 1 and consistent decrease in event 2.

4 DISCUSSION AND LIMITATIONS

We work towards a future where unsupervised multitask learning in audio synthesis and transformation emerges from data and model scale. Our proposed framework *ComposableART Fugatto* establishes our first step towards this direction, and in this step we become aware of our current limitations and challenges. For example, optimizing dataset sampling weights to drive performance on multiple benchmarks is a herculean task, and generative models for audio and video would certainly benefit from research similar to (Albalak et al., 2023; Chung et al., 2023; Xie et al., 2024). Along straighter paths, we plan to replace mels with a latent representation that better supports low frequencies and stereo. We believe this modification should be rather straight forward. Further work is necessary to establish the impact of data and free-form instructions on *Fugatto*’s emergent abilities, especially using language to combine tasks not jointly seen during training. Finally, *ComposableART* requires more analysis on the choice of weights, and how the norm of the vector field can be used for easier and more stable control.

540 REPRODUCIBILITY STATEMENT

541

542 We plan to release our dataset and instruction generation code to facilitate reproducible research.

543

544 REFERENCES

545

546 Alon Albalak, Liangming Pan, Colin Raffel, and William Yang Wang. Efficient online data mixing
547 for language model pre-training. In *RO-FoMo: Robustness of Few-shot and Zero-shot Learning in*
548 *Large Foundation Models*, 2023.549 Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning. In
550 *Proceedings of the 26th annual international conference on machine learning*, pp. 41–48, 2009.

551

552 Mathieu Bernard and Hadrien Titeux. Phonemizer: Text to phones transcription for multiple languages
553 in python. *Journal of Open Source Software*, 6(68):3958, 2021.554 Paul Boersma and Vincent Van Heuven. Speak and unspeak with praat. *Glott International*, 5(9/10):
555 341–347, 2001.

556

557 Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,
558 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
559 few-shot learners. *Advances in Neural Information Processing Systems*, 33:1877–1901, 2020.560 Hyung Won Chung, Noah Constant, Xavier Garcia, Adam Roberts, Yi Tay, Sharan Narang, and
561 Orhan Firat. Unimax: Fairer and more effective language sampling for large-scale multilingual
562 pretraining. *arXiv preprint arXiv:2304.09151*, 2023.

563

564 Shuqi Dai, Ming-Yu Liu, Rafael Valle, and Siddharth Gururani. Expressivesinger: Multilingual and
565 multi-style score-based singing voice synthesis with expressive performance control. In *ACM*
566 *Multimedia 2024*, 2024.567 SeungHeon Doh, Keunwoo Choi, Jongpil Lee, and Juhan Nam. Lp-musiccaps: Llm-based pseudo
568 music captioning. *arXiv preprint arXiv:2307.16372*, 2023.569 Yilun Du, Shuang Li, and Igor Mordatch. Compositional visual generation with energy based
570 models. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Ad-*
571 *vances in Neural Information Processing Systems*, volume 33, pp. 6637–6647. Curran Asso-
572 ciates, Inc., 2020. URL [https://proceedings.neurips.cc/paper_files/paper/](https://proceedings.neurips.cc/paper_files/paper/2020/file/49856ed476ad01fcff881d57e161d73f-Paper.pdf)
573 [2020/file/49856ed476ad01fcff881d57e161d73f-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/49856ed476ad01fcff881d57e161d73f-Paper.pdf).574 Jesse Engel, Cinjon Resnick, Adam Roberts, Sander Dieleman, Mohammad Norouzi, Douglas Eck,
575 and Karen Simonyan. Neural audio synthesis of musical notes with wavenet autoencoders. In
576 *International Conference on Machine Learning*, pp. 1068–1077. PMLR, 2017.

577

578 Sefik Emre Eskimez, Xiaofei Wang, Manthan Thakker, Canrun Li, Chung-Hsien Tsai, Zhen Xiao,
579 Hemin Yang, Zirun Zhu, Min Tang, Xu Tan, et al. E2 tts: Embarrassingly easy fully non-
580 autoregressive zero-shot tts. *arXiv preprint arXiv:2406.18009*, 2024.581 Arushi Goel, Zhifeng Kong, Rafael Valle, and Bryan Catanzaro. Audio dialogues: Dialogues dataset
582 for audio and music understanding. *arXiv preprint arXiv:2404.07616*, 2024.

583

584 Yuan Gong, Hongyin Luo, Alexander H Liu, Leonid Karlinsky, and James Glass. Listen, think, and
585 understand. *arXiv preprint arXiv:2305.10790*, 2023.586 Zhifang Guo, Yichong Leng, Yihan Wu, Sheng Zhao, and Xu Tan. Prompttts: Controllable text-to-
587 speech with text descriptions. In *ICASSP 2023-2023 IEEE International Conference on Acoustics,*
588 *Speech and Signal Processing (ICASSP)*, pp. 1–5. IEEE, 2023.589 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*,
590 2022.

591

592 Tero Karras, Miika Aittala, Tuomas Kynkäänniemi, Jaakko Lehtinen, Timo Aila, and Samuli Laine.
593 Guiding a diffusion model with a bad version of itself, 2024. URL [https://arxiv.org/](https://arxiv.org/abs/2406.02507)
[abs/2406.02507](https://arxiv.org/abs/2406.02507).

- 594 Rafal Kocielnik, Shrimai Prabhumoye, Vivian Zhang, R Michael Alvarez, and Anima Anandkumar.
595 Autobiastest: Controllable sentence generation for automated and open-ended social bias testing in
596 language models. *arXiv preprint arXiv:2302.07371*, 2023.
- 597
598 Zhifeng Kong, Wei Ping, Amrith Dantrey, and Bryan Catanzaro. Cleanunet 2: A hybrid speech
599 denoising model on waveform and spectrogram. *arXiv preprint arXiv:2309.05975*, 2023.
- 600
601 Zhifeng Kong, Arushi Goel, Rohan Badlani, Wei Ping, Rafael Valle, and Bryan Catanzaro. Audio
602 flamingo: A novel audio language model with few-shot learning and dialogue abilities. *arXiv
603 preprint arXiv:2402.01831*, 2024a.
- 604
605 Zhifeng Kong, Sang-gil Lee, Deepanway Ghosal, Navonil Majumder, Ambuj Mehrish, Rafael Valle,
606 Soujanya Poria, and Bryan Catanzaro. Improving text-to-audio models with synthetic captions.
607 *arXiv preprint arXiv:2406.15487*, 2024b.
- 608
609 Matthew Le, Apoorv Vyas, Bowen Shi, Brian Karrer, Leda Sari, Rashel Moritz, Mary Williamson,
610 Vimal Manohar, Yossi Adi, Jay Mahadeokar, et al. Voicebox: Text-guided multilingual universal
611 speech generation at scale. *Advances in neural information processing systems*, 36, 2024.
- 612
613 Sang-gil Lee, Wei Ping, Boris Ginsburg, Bryan Catanzaro, and Sungroh Yoon. Bigvgan: A universal
614 neural vocoder with large-scale training. In *The Eleventh International Conference on Learning
615 Representations*, 2023. URL https://openreview.net/forum?id=iTtGCMDEzS_.
- 616
617 Yeonghyeon Lee, Inmo Yeon, Juhan Nam, and Joon Son Chung. Voiceldm: Text-to-speech with
618 environmental context. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech
619 and Signal Processing (ICASSP)*, pp. 12566–12571. IEEE, 2024.
- 620
621 Yichong Leng, Zhifang Guo, Kai Shen, Xu Tan, Zeqian Ju, Yanqing Liu, Yufei Liu, Dongchao Yang,
622 Leying Zhang, Kaitao Song, et al. Promptts 2: Describing and generating voices with text prompt.
623 *arXiv preprint arXiv:2309.02285*, 2023.
- 624
625 Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching
626 for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.
- 627
628 Haohe Liu, Ke Chen, Qiao Tian, Wenwu Wang, and Mark D Plumbley. Audiosr: Versatile audio
629 super-resolution at scale. In *ICASSP 2024-2024 IEEE International Conference on Acoustics,
630 Speech and Signal Processing (ICASSP)*, pp. 1076–1080. IEEE, 2024.
- 631
632 Nan Liu, Shuang Li, Yilun Du, Antonio Torralba, and Joshua B. Tenenbaum. Compositional visual
633 generation with composable diffusion models, 2023. URL <https://arxiv.org/abs/2206.01714>.
- 634
635 Steven R Livingstone and Frank A Russo. The ryerson audio-visual database of emotional speech
636 and song (ravdess): A dynamic, multimodal set of facial and vocal expressions in north american
637 english. *PLoS one*, 13(5):e0196391, 2018.
- 638
639 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Confer-
640 ence on Learning Representations*, 2019. URL <https://openreview.net/forum?id=Bkg6RiCqY7>.
- 641
642 Justin Lovelace, Soham Ray, Kwangyoum Kim, Kilian Q Weinberger, and Felix Wu. Simple-tts:
643 End-to-end text-to-speech synthesis with latent diffusion. 2023.
- 644
645 Weili Nie, Arash Vahdat, and Anima Anandkumar. Controllable and compositional generation with
646 latent-space energy-based models, 2021. URL <https://arxiv.org/abs/2110.10873>.
- 647
648 OpenAI. Gpt-4o: A powerful multimodal language model. [https://openai.com/research/
649 hello-gpt-4o](https://openai.com/research/hello-gpt-4o), 2024. Accessed: 2024-09-21.
- 650
651 Ashis Pati, Siddharth Kumar Gururani, and Alexander Lerch. dmelodies: A music dataset for
652 disentanglement learning. In *Proceedings of the 21th International Society for Music Information
653 Retrieval Conference, ISMIR 2020, Montreal, Canada, October 11-16, 2020*, pp. 125–133, 2020.
654 URL <http://archives.ismir.net/ismir2020/paper/000300.pdf>.

- 648 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language
649 models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
650
- 651 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi
652 Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text
653 transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
- 654 Zafar Rafii, Antoine Liutkus, Fabian-Robert Stöter, Stylianos Ioannis Mimilakis, and Rachel Bittner.
655 The musdb18 corpus for music separation. 2017.
656
- 657 Jonathan Shen, Ruoming Pang, Ron J Weiss, Mike Schuster, Navdeep Jaitly, Zongheng Yang, Zhifeng
658 Chen, Yu Zhang, Yuxuan Wang, Rj Skerrv-Ryan, et al. Natural tts synthesis by conditioning
659 wavenet on mel spectrogram predictions. In *2018 IEEE international conference on acoustics,
660 speech and signal processing (ICASSP)*, pp. 4779–4783. IEEE, 2018.
- 661 Spotify. Pedalboard documentation. [https://spotify.github.io/pedalboard/index.
662 html](https://spotify.github.io/pedalboard/index.html), 2024. Accessed: 2024-09-23.
663
- 664 Alexander Tong, Nikolay Malkin, Guillaume Hugué, Yanlei Zhang, Jarrid Rector-Brooks, Kilian
665 Fatras, Guy Wolf, and Yoshua Bengio. Conditional flow matching: Simulation-free dynamic
666 optimal transport. *arXiv preprint arXiv:2302.00482*, 2023.
- 667 Rafael Valle, Kevin Shih, Ryan Prenger, and Bryan Catanzaro. Flowtron: an autoregressive flow-based
668 generative network for text-to-speech synthesis. *arXiv preprint arXiv:2005.05957*, 2020.
669
- 670 A Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017.
- 671 Apoorv Vyas, Bowen Shi, Matthew Le, Andros Tjandra, Yi-Chiao Wu, Baishan Guo, Jiemin Zhang,
672 Xinyue Zhang, Robert Adkins, William Ngan, et al. Audiobox: Unified audio generation with
673 natural language prompts. *arXiv preprint arXiv:2312.15821*, 2023.
674
- 675 Chengyi Wang, Sanyuan Chen, Yu Wu, Ziqiang Zhang, Long Zhou, Shujie Liu, Zhuo Chen, Yanqing
676 Liu, Huaming Wang, Jinyu Li, et al. Neural codec language models are zero-shot text to speech
677 synthesizers. *arXiv preprint arXiv:2301.02111*, 2023.
- 678 Xiaofei Wang, Manthan Thakker, Zhuo Chen, Naoyuki Kanda, Sefik Emre Eskimez, Sanyuan Chen,
679 Min Tang, Shujie Liu, Jinyu Li, and Takuya Yoshioka. Speechx: Neural codec language model
680 as a versatile speech transformer. *IEEE/ACM Transactions on Audio, Speech, and Language
681 Processing*, 2024.
- 682 Sang Michael Xie, Hieu Pham, Xuanyi Dong, Nan Du, Hanxiao Liu, Yifeng Lu, Percy S Liang,
683 Quoc V Le, Tengyu Ma, and Adams Wei Yu. Doremi: Optimizing data mixtures speeds up
684 language model pretraining. *Advances in Neural Information Processing Systems*, 36, 2024.
685
- 686 Jingjing Xu, Xu Sun, Zhiyuan Zhang, Guangxiang Zhao, and Junyang Lin. Understanding and
687 improving layer normalization. *Advances in neural information processing systems*, 32, 2019.
- 688 Linting Xue, Aditya Barua, Noah Constant, Rami Al-Rfou, Sharan Narang, Mihir Kale, Adam
689 Roberts, and Colin Raffel. Byt5: Towards a token-free future with pre-trained byte-to-byte models.
690 *Transactions of the Association for Computational Linguistics*, 10:291–306, 2022.
691
- 692 Dongchao Yang, Jinchuan Tian, Xu Tan, Rongjie Huang, Songxiang Liu, Xuankai Chang, Jiatong
693 Shi, Sheng Zhao, Jiang Bian, Xixin Wu, et al. Uniaudio: An audio foundation model toward
694 universal audio generation. *arXiv preprint arXiv:2310.00704*, 2023.
- 695 Jinhyeok Yang, Junhyeok Lee, Hyeong-Seok Choi, Seunghun Ji, Hyeongju Kim, and Juheon Lee.
696 Dualspeech: Enhancing speaker-fidelity and text-intelligibility through dual classifier-free guidance.
697 *arXiv preprint arXiv:2408.14423*, 2024.
698
699
700
701

702 A APPENDIX

703 A.1 DATASETS

704 In this section we provide information about our dataset generation strategy, tasks, instruction
705 generators, and a full list of datasets created and used during training, including their sampling
706 weights and task probabilities.

707 A.1.1 APPROACH TO DATASET GENERATION

708 In this subsection we provide further descriptions of the procedures in each of the five pillars described
709 in Section 2.1. Overall, the captions and instructions generation process is centered around leveraging
710 LLMs to transform tags into richer descriptions, always being careful that the description relates to
711 the audio content, and to create what we call *instruction generators*, i.e. python methods that create
712 free-form instructions based on the synthetic descriptions and task at hand.

713 **I** – The python code snippet below is an LLM generated script that outputs free-form instructions for
714 the task *reverse sound*, for example. Once such a script exists, it can be used to prompt an LLM to
715 generate scripts for other tasks with increasing levels of complexity.

```

720 class ReverseAudioInstructor:
721     audio_references = {
722         'standard': {
723             'verbs': ['reverse', 'play backward', 'invert'],
724             'gerunds': ['reversing', 'playing backward', 'inverting'],
725             'contexts': ['audio', 'sound', 'recording', 'clip', 'track'],
726             'asks': ['Can you', 'Please', 'Could you', 'I need you to'],
727             'styles': ['completely', 'precisely', 'accurately'],
728             'mentions': ['I provided', 'I sent', 'I attached']
729         },
730         ...
731     }
732
733     @staticmethod
734     def generate_instruction(persona='standard'):
735         ref = ReverseAudioInstructor.audio_references[persona]
736
737         templates = [
738             "{ask} {verb} the {context}.",
739             "{verb} the {context} {style}.",
740             "{ask} {verb} the {context} {mention}.",
741             "Your task is to {verb} the {context}.",
742             "We need the {context} {mention} to be {gerund}.",
743             "{gerund} the {context} is the goal.",
744             "Please focus on {gerund} the {context}.",
745             "The objective is {gerund} the {context} {mention}."
746         ]
747
748         template = random.choice(templates)
749
750         instruction = template.format(
751             ask=random.choice(ref['asks']),
752             verb=random.choice(ref['verbs']),
753             gerund=random.choice(ref['gerunds']),
754             context=random.choice(ref['contexts']),
755             style=random.choice(ref['styles']),
756             mention=random.choice(ref['mentions'])
757         )
758
759         return instruction.capitalize()
760
761 if __name__ == '__main__':
762     for persona in ReverseAudioInstructor.audio_references.keys():
763         print(f"\n{persona.capitalize()} instructions:")

```

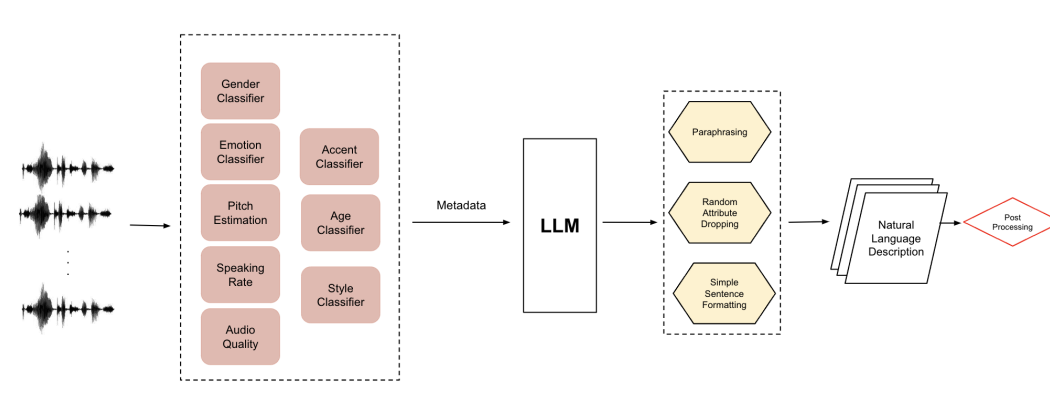

756 **II** – The python code snippet below is an LLM generated script that outputs free-form instructions
 757 for the task *speech modulation*, for example.. Note that the instructions refer to relative changes such
 758 as increase and decrease with different magnitudes such small or large changes. Once such a script
 759 exists, it can be used to prompt an LLM to generate scripts for other tasks with increasing levels of
 760 complexity.

```

761 class SpeechModulationInstructor:
762     instructions = {
763         'scale_formant': {
764             'increase': {
765                 'small': [
766                     "add a touch more resonance",
767                     "slightly enhance the formant frequencies",
768                     ...
769                 ],
770                 ...
771                 'large': [
772                     "dramatically enhance the resonance",
773                     "massively boost the formant frequencies",
774                     ...
775                 ]
776             },
777             'decrease': {
778                 'small': [
779                     "tone down the resonance just a little",
780                     "slightly reduce the formant frequencies",
781                     ...
782                 ],
783                 ...
784                 'large': [
785                     "dramatically reduce the resonance",
786                     "massively lower the formant frequencies",
787                     ...
788                 ]
789             }
790         },
791         ...
792     }
793
794     @staticmethod
795     def get_instruction(modulation, direction, intensity):
796         return random.choices(SpeechModulationInstructor.instructions[modulation][direction][intensity])[0]
797
798     @staticmethod
799     def combine_instructions(modulations):
800         parts = []
801         for modulation_i in modulations:
802             modulation = modulation_i['modulation']
803             direction = modulation_i['direction']
804             intensity = modulation_i['intensity']
805             instruction_part = SpeechModulationInstructor.get_instruction(modulation, direction, intensity)
806             parts.append(instruction_part)
807         instruction = ""
808         if parts:
809             # Combine parts into a single instruction
810             if len(parts) > 1:
811                 combined_instruction = ", and ".join(parts[:-1]) + ", and " + parts[-1]
812             else:
813                 combined_instruction = parts[0]
814             instruction = "Let's " + combined_instruction + "."
815         return instruction
816
817 if __name__ == '__main__':
818     print("main")

```

810 **III** - Figure 4 provides a visual depiction of our speech captioning, dubbed Prompt-2-Voice
 811 (P2V) pipeline. We applied this approach to several open-source soeech datasets. Additionally, we
 812 incorporated existing speaker descriptions from datasets like PromptSpeech (Guo et al., 2023), which
 813 provide details on gender, pitch, volume, and speaking rate, and enriched them by adding emotion and
 814 quality information. This automation not only streamlines the process but also allows us to efficiently
 815 label in-the-wild datasets, significantly scaling the available training data wigh high quality captions.
 816



830

831 Figure 4: Synthetic caption generation pipeline for Prompt-to-Voice (P2V).

832

833 **V** - We use audio effects libraries to create synthetic paired data where some factors are held constant
 834 while others change. We use Pedalboard (Spotify, 2024), a Python library for audio effects, to apply
 835 different audio effects to sounds believed not to have effects. For each effect, there are a number
 836 of parameters that can be altered. We generate on the fly audio segments with the same effect but
 837 different levels of a single parameter while keeping the others fixed, so that we can not only create
 838 instructions with synthetic pairs of unaffected vs. effected data, but also synthetic pairs where we ask
 839 the model to increase or decrease the intensity of a given parameter (e.g. "increase the compression
 840 rate a little bit" or "reduce the room size for the reverb moderately"). For each parameter, we
 841 determined a reasonable range for the parameter, and increments that would correspond to "a little",
 842 "moderate", and "a lot" of that parameter.
 843

| | |
|-----|------------------------|
| 864 | A.1.2 LIST OF DATASETS |
| 865 | |
| 866 | |
| 867 | |
| 868 | |
| 869 | |
| 870 | |
| 871 | |
| 872 | |
| 873 | |
| 874 | |
| 875 | |
| 876 | |
| 877 | |
| 878 | |
| 879 | |
| 880 | |
| 881 | |
| 882 | |
| 883 | |
| 884 | |
| 885 | |
| 886 | |
| 887 | |
| 888 | |
| 889 | |
| 890 | |
| 891 | |
| 892 | |
| 893 | |
| 894 | |
| 895 | |
| 896 | |
| 897 | |
| 898 | |
| 899 | |
| 900 | |
| 901 | |
| 902 | |
| 903 | |
| 904 | |
| 905 | |
| 906 | |
| 907 | |
| 908 | |
| 909 | |
| 910 | |
| 911 | |
| 912 | |
| 913 | |
| 914 | |
| 915 | |
| 916 | |
| 917 | |

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971**Open Source Vocal Datasets**

Table 7: Open source vocal datasets. tts is Text-To-Speech and tta is Text-To-Audio (Audio from Caption)

| Dataset | Sampling Weights | Task and Probabilities |
|---|------------------|--|
| AISHELL-3 | 2.46 | tts |
| CML-Dutch | 0.769 | tts |
| CML-French | 0.334 | tts |
| CML-German | 1.762 | tts |
| CML-Italian | 0.158 | tts |
| CML-Polish | 0.051 | tts |
| CML-Portuguese | 0.086 | tts |
| CML-Spanish | 0.512 | tts |
| common-accent-AccentClassification | 0.131 | tta |
| CREMA-D-EmotionClassification | 0.074 | tta |
| DAPS-Enhancement | 0.21 | enhancement paired 0.99 deenhancement paired 0.01 |
| DNS-Challenge-2020 | 14.247 | source separation |
| emov-db-EmotionClassification | 0.068 | tta |
| IEMOCAP-EmotionClassification | 0.059 | tta |
| jl-corpus-EmotionClassification | 0.024 | tta |
| LibriTTS-Clean-100 | 1.23 | tts |
| LibriTTS-Clean-360 | 4.31 | tts |
| LibriTTS-Other-500 | 0.643 | tts |
| LibriVox-English | 52.54 | tts |
| LibriVox-French | 0.909 | tts |
| LibriVox-German | 1.147 | tts |
| LibriVox-Italian | 0.114 | tts |
| LibriVox-Portuguese | 0.12 | tts |
| LibriVox-Spanish | 0.276 | tts |
| LIMMITS2024-* | 0.768 | tts 0.80 inpainting 0.10 upsampling 0.10 |
| MSP-PODCAST-Publish-1.9-EmotionClassification | 0.451 | tta |
| NonSpeech7k-EventClassification | 0.063 | tta |
| OMGEmotion-EmotionClassification | 0.017 | tta |
| ravdess-EmotionClassification | 0.014 | tta |
| SongDescriber-AudioCaptioning | 0.077 | tta |
| SONYC-UST-EventClassification | 0.279 | tta |
| tess-EmotionClassification | 0.028 | tta |
| VCTK-VoiceConversion | 2.89 | voice conversion paired vctk |
| VCTK-TTS | 0.137 | tts |
| VocalSound-VocalClassification | 0.155 | tta |

Open Source Non-Vocal Datasets

Table 8: Open Source Non-Vocal Datasets. tts is Text-To-Speech and tta is Text-To-Audio (tta)

| Dataset | Sampling Weights | Task and Probabilities |
|--|------------------|--|
| audiocaps-AudioCaptioning | 5.21 | tta 0.90 inpainting 0.10 |
| BBCSoundEffects-AudioDescription | 0.15 | tta |
| chime-home-EventClassification | 0.05 | tta |
| Clotho-AQA-EventClassification | 0.01 | tta |
| Clotho-AQA singlelabel-EventClassification | 0.07 | tta |
| Clotho-v2-AudioCaptioning | 0.19 | tta |
| CochlScene-SceneClassification | 0.61 | tta |
| Epidemic sound-AudioCaptioning | 0.41 | tta |
| ESC-50 | 1.12 | tta 0.90 inpainting 0.10 |
| FMA-GenreClassification | 1.04 | tta |
| FSD50k-EventClassification | 0.41 | tta |
| GTZAN-GenreClassification | 0.01 | tta |
| LP-MusicCaps-MC-AudioCaptioning | 0.07 | tta |
| LP-MusicCaps-MSD-0 | 8.56 | tta 0.90 inpainting 0.10 |
| LP-MusicCaps-MSD-1 | 8.62 | tta 0.90 inpainting 0.10 |
| LP-MusicCaps-MSD-AudioCaptioning | 11.82 | tta |
| LP-MusicCaps-MTT-AudioCaptioning | 0.47 | tta |
| MACS-AudioCaptioning | 0.17 | tta |
| Maestro | 0.08 | tta 0.90 upsampling 0.10 |
| Medley-solos-DB | 0.19 | tta 0.90 upsampling 0.10 |
| MSD | 12.38 | tta 0.80 upsampling 0.10 inpainting 0.10 |
| MTG-Jamendo | 2.16 | tta 0.90 inpainting 0.10 |
| musdbhq-InstrClassification | 0.10 | tta |
| MusicCaps-AudioCaptioning | 0.54 | tta 0.90 inpainting 0.10 |
| NonSpeech7k-EventClassification | 0.06 | tta |
| NSynth-MIR | 2.31 | tta |
| SongDescriber-AudioCaptioning | 0.08 | tta |
| SONYC-UST-EventClassification | 0.28 | tta |
| SoundDescs-AudioDescription | 0.23 | tta |
| SoundVE-Caps | 37.00 | tta |
| UrbanSound8K-EventClassification | 0.09 | tta |
| WavText5K-AudioCaptioning | 0.04 | tta |
| WavText5K-Tagging | 0.02 | tta |

1026 **New Datasets generated from Open Source Data**
 1027

1028 Table 9: New datasets generated from open source data. tts is Text-To-Speech, tta is Text-To-Audio
 1029 (Audio from Caption), and P2V refers to prompt to voice, in which LLMs were used to create full
 1030 form captions given speech attributes. The -AF suffix indicates synthetic captions generated with
 1031 Audio Flamingo.
 1032

| 1033 | Dataset | Sampling Weights | Task and Probabilities |
|------|---|------------------|--|
| 1034 | audiocaps-AudioCaptioning-AF | 5.79 | tta 0.90 inpainting 0.10 |
| 1035 | AudioSet-AF | 37.24 | tta 0.80 inpainting 0.05 inpainting random mask 0.05 upsampling 0.09 downsampling 0.01 |
| 1036 | AISHELL-3-AddRemove-Sound-Effects | 3.23 | add sound effects 0.50 remove sound effects 0.50 |
| 1037 | AISHELL-3-SoundEffectsModulation | 3.84 | sound effects modulation |
| 1038 | AISHELL-3-SpeechModulationPraat | 3.84 | speech modulation praat |
| 1039 | CLAP freesound-AF | 2.82 | tta |
| 1040 | CREMA-D-P2V | 0.62 | tts 0.80 inpainting 0.10 upsampling 0.10 |
| 1041 | EGFxSet-AddSoundEffects | 0.00 | add sound effects paired |
| 1042 | EGFxSet-RemoveSoundEffects | 0.00 | remove sound effects paired |
| 1043 | emov-db-EmotionClassification | 0.07 | tta |
| 1044 | ESD-ChangeEmotion | 1.90 | speech modulation paired |
| 1045 | ESD-ENGLISH-P2V | 0.93 | tts 0.80 inpainting 0.10 upsampling 0.10 |
| 1046 | ESD-MANDARIN-P2V | 0.63 | tts 0.80 inpainting 0.10 upsampling 0.10 |
| 1047 | IEMOCAP-EmotionClassification | 0.06 | tta |
| 1048 | jl-corpus-EmotionClassification | 0.02 | tta |
| 1049 | LibriTTS-Clean-100-Add-Remove-Sound-Effects | 7.85 | add sound effects 0.50 remove sound effects 0.50 |
| 1050 | LibriTTS-Clean-100-SoundEffectsModulation | 0.06 | sound effects modulation |
| 1051 | LibriTTS-Clean-100-SpeechModulationPraat | 9.34 | speech modulation praat |
| 1052 | LibriTTS-Clean-100-Enhancement | 0.85 | enhancement paired 0.99 deenhancement paired 0.01 |
| 1053 | LibriTTS-Clean-360-Enhancement | 2.71 | enhancement paired 0.99 deenhancement paired 0.01 |
| 1054 | LibriTTS-Other-500-Enhancement | 6.11 | enhancement paired |
| 1055 | JL-CORPUS-P2V | 0.62 | tts 0.80 inpainting 0.10 upsampling 0.10 |
| 1056 | NVYT-English | 27.23 | tts 0.80 inpainting 0.05 inpainting random mask 0.05 upsampling 0.04 reverse sound 0.04 downsampling 0.01 |
| 1057 | RAVDESS-P2V | 0.62 | tts 0.80 inpainting 0.10 upsampling 0.10 midi2audio |
| 1058 | LMD-Aligned | 1.17 | tta |
| 1059 | musdbhq-InstrClassification | 0.10 | tta |
| 1060 | musdbhq-add-sound | 3.81 | add sound 0.50 add sound to mixture 0.50 |
| 1061 | musdbhq-singing | 0.86 | singing voice synthesis |
| 1062 | musdbhq-singing-aggregated | 0.86 | singing voice synthesis |
| 1063 | musdbhq-source-separation | 3.81 | source separation 0.50 remove sound from mixture 0.50 |
| 1064 | NonSpeech7k-EventClassification | 0.06 | tta |
| 1065 | NSynth-MIR | 2.31 | tta |
| 1066 | NVYT-English | 27.23 | tts 0.80 inpainting 0.05 inpainting random mask 0.05 upsampling 0.04 reverse sound 0.04 downsampling 0.01 |
| 1067 | OMGEmotion-EmotionClassification | 0.02 | tta |
| 1068 | ravdess-EmotionClassification | 0.01 | tta |
| 1069 | RAVDESS-CreateVariation | 0.02 | speech modulation paired |
| 1070 | RAVDESS-ChangeIntensity | 0.15 | speech modulation paired |
| 1071 | RAVDESS-ChangeEmotion | 0.22 | speech modulation paired |
| 1072 | ExpressiveSinger-P2V | 2.84 | singing voice synthesis nolanguage |
| 1073 | tess-EmotionClassification | 0.03 | tta |
| 1074 | VGG-AF | 0.92 | tta |
| 1075 | WavCaps-AF | 7.83 | tta |

1080 Below we provide descriptions for the suffixes associated with our generated datasets presented in
 1081 Table 7, Table 8, and Table 9.

1082 **-AF:** Refers to generating synthetic captions with the strategy described in (Kong et al., 2024b). In
 1083 summary, the strategy consists of using an audio understanding model, here Audio Flamingo Chat,
 1084 to caption sounds in the wild, and then filtering out synthetic captions based on the CLAP cosine
 1085 similarity between the audio and the synthetic caption. In this paper, we used this strategy to create
 1086 synthetic captions for AudioCaps, AudioSet, CLAP Freesound, and WavCaps.

1087 **-{Add, Remove} Sound Effects:** These refer to applying, on the fly, audio effects modifications to
 1088 existing audio data to create paired data that describe tasks related to adding and removing sound
 1089 effects. We focus on speech, which we assume has the least amount of audio effects applied to it,
 1090 especially when compared to music data. In this iteration, the audio modulations are performed
 1091 with (Spotify, 2024) and include a large list of effects such as chorus, reverb, distortion, amongst
 1092 others.

1093 **-{Add, Remove}:** This refers to splitting an audio mixture into separate audio stems and creating
 1094 manifests that add one track given another track. In this *opus*, we have not explored creating artificial
 1095 mixtures by adding random waveforms but imagine this strategy can yield good results.

1096 **-P2V:** P2V, or prompt-to-voice, is a task that enables control of speech synthesis through textual
 1097 prompts, allowing for the description of speaker characteristics when an appropriate audio prompt that
 1098 matches the desired persona is unavailable. The strategy has several components. First, we extract
 1099 and curate all possible tags from existing metadata. For example, the expressive singer (Dai et al.,
 1100 2024) dataset includes, implicitly and explicitly, information about vocal range, accent, language,
 1101 style and others. Then, following the strategy in Section 2.1 we first prompt LLMs to produce long
 1102 form descriptions of each attribute, then we prompt an LLM to create a instruction generator that
 1103 combines the modified attributes in different ways.

1104 **-Singing:** refers to leveraging music stems dataset with vocal tracks to create singing datasets. As
 1105 usual, we leverage all the metadata available, combined with transcriptions obtained from speech
 1106 transcription models and song captions obtained by prompting LLMs. In cases where an accompa-
 1107 niment or backing track is available, we associate the timestamps on each lyrics snippet with the
 1108 respective accompaniment. Once the metadata types are available, we prompt an LLM to create a
 1109 python method that produces instructions given the metadata for a singing voice synthesis task with
 1110 or without accompaniment.

1111 **-Singing-Aggregated:** -Singing-Aggregated adds to -Singing the combination of adjacent sentences,
 1112 given criteria such as max gap between sentences and max sentence length.

1113 **-Source Separation and Add To Mixture:** refers to leveraging existing In this iteration, assuming
 1114 doing such would produce samples undesirable out-of-distribution samples, we do not create mixtures
 1115 by combining random samples together. Instead, we leverage mixtures that are already exist.

1116 **-Speech Modulation:** refers to applying speech modulations to existing speech data to create paired
 1117 data that describe a task in which a relative change is being applied to an attribute of speech, e.g.
 1118 "Increase the speaking rate in this sample". In this iteration, the speech modulations are performed
 1119 with Praat (Boersma & Van Heuven, 2001) and we focus on transformations such as formant scaling,
 1120 F0 mean scaling (equivalent to transposition in music), F0 variance scaling (equivalent to flattening
 1121 or expanding the F0 contour), and speaking rate scaling (equivalent to making a person speak faster
 1122 or slower). Due to artifacts that can be introduced in the audio as a result of the modulation, we limit
 1123 the scaling values to ranges that we find acceptable in terms of audio quality.

1124 **-Speech Modulation Paired:** This refers to leveraging available paired data to create new tasks.
 1125 For example, an emotional speech dataset with paired data can be used to enable an emotion
 1126 transformation tasks. The procedure is straight forward and consists of first grouping samples by
 1127 speaker and transcript, then creating pairs that establish relationships between two samples. For
 1128 example, given a speaker and transcript, two samples with different emotion can be used to define
 1129 a "convert emotion task", two samples with the same emotion and intensity can be used to create a
 1130 "create a variation of this speech", and two samples with the same emotion but different intensity
 1131 can be used to define a "increase the intensity in the emotion task". Once the pairs and new tasks
 1132 are defined, we prompt an LLM to create a script that takes in the attributes and tasks to generate
 1133 task-specific instructions.

1134 **–Sound Effects Modulation:** refers to applying sound effects modulations to existing data to create
 1135 paired data for relative change in audio effects in a file. In this iteration, the speech modulations are
 1136 performed with Praat (Boersma & Van Heuven, 2001) and focus on formant scaling, F0 mean scaling,
 1137 F0 variance scaling and speaking rate scaling. We limit the values to the range below
 1138

1141 A.1.3 TASKS

1144 Below we provide the list of tasks that are known to be supported by our model, and the respective
 1145 minimal instruction for that task. Please note that we provide minimal instructions due to layout
 1146 reasons, not due to model limitations.

| 1147 Task | Minimal Instruction |
|---|---|
| 1148 Add Sounds to Instrument Track | |
| 1149 using audio example | add drums like this <audio> to this guitar track <audio> |
| 1150 using text description | add rock drums to this guitar track <audio> |
| 1151 Add Sounds to Sound Mixture | |
| 1152 using audio example | add guitar like this <audio> to this backing track <audio> |
| 1153 using text description | add guitar to this backing track <audio> |
| 1154 Apply Sound Effects | add long reverb to this audio <audio> |
| 1155 Copy Audio | copy this <audio> |
| 1156 Downsample Audio | downsample this to 16kHz <audio> |
| 1157 Enhance Audio | enhance this sound <audio> |
| 1158 Generate Audio From Captions | generate a saxophone barking |
| 1159 Generate Music From Captions | generate a calm sound track |
| 1160 Generate Speech From Captions | generate an angry voice |
| 1161 Inpaint Audio (Continuation) | inpaint this sound <audio> |
| 1162 MIDI2Audio | |
| 1163 using audio example | turn this MIDI track <audio> into natural audio. |
| 1164 using text description | turn this MIDI track <audio> into a heavy metal track. |
| 1165 Remove Sound Effects | remove sound effects <audio> |
| 1166 Remove Sound from Sound Mixture | |
| 1167 from audio examples | remove this <audio> from these sounds |
| 1168 from text captions | remove the piano from this music track <audio> |
| 1169 Reverse Sound | reverse this sound |
| 1170 Singing-Voice Synthesis | |
| 1171 given Speech Captions | sing this 'At last my love' with a male voice |
| 1172 given Language | sing this 'At last my love' in English |
| 1173 given Accent | sing this 'At last my love' in English with an Italian accent |
| 1174 Separate Audio into Sources | |
| 1175 given audio examples | give me this <audio> from this music track <audio> |
| 1176 given text description | give me the piano from this music track <audio> |
| 1177 Sound Effects Modulation | (Increase and Decrease Attributes) |
| 1178 Chorus | increase the chorus a bit <audio> |
| 1179 Compressor | decrease the compressor threshold <audio> |
| 1180 Delay | decrease the delay time <audio> |
| 1181 Distortion | increase the distortion <audio> |
| 1182 Limiter | make the limiter less strong <audio> |
| 1183 Phaser | increase the phaser speed <audio> |
| 1184 Reverb | increase the reverb room size in this <audio> |
| 1185 Speech Modulation (Paired) | |
| 1186 Voice Conversion | convert from this <audio> to this <audio> |
| 1187 Accent Conversion | convert from British English to American English |
| Emotion Conversion | make this calm sample <audio> sound angry |
| Speech Variation | give me a different take on this voice <audio> |
| 1188 Speech Modulation (Praat) | |
| 1189 Scale Speaking Rate | increase the speaking rate <audio> |
| 1190 Scale F0 Mean | increase the average pitch <audio> |
| 1191 Scale F0 Variance | increase the variance in pitch <audio> |
| 1192 Scale Formants | scale the formants here <audio> |
| 1193 Text-To-Speech | |
| 1194 from Speech Prompt | say this 'May the force!' given this voice <audio> |
| 1195 from Speech Captions | say this 'May the force!' with a male voice |
| 1196 from Language | say this 'May the force!' in English |
| 1197 from Accent | say this 'May the force!' in English with an Italian accent |
| 1198 Upsample Audio | upsample this sound <audio> |

1188 A.1.4 INSTRUCTIONS

1189

1190 After an informal evaluation, we concluded that, for our purpose, Claude Sonnet is better than GPT4-o
1191 at producing instructions and following prompts. Below we provide examples of template-based
1192 and free-form instructions for a handful of tasks and datasets. For clarity, we remove the redundant
1193 'input:' and 'output:' parts from all instructions.

1194 AISHELL-3-AddRemove-Sound-Effects

1195 Eradicating the audio from the provided audio material with the
1196 Distortion, and Phaser effect is your task. Focus on eradicating
1197 it systematically.

1198 AISHELL-3-SoundEffectsModulation

1199 We need to minimize the delay time.</caption>

1200

1201 AISHELL-3-SpeechModulationPraat

1202 Let's subtly enhance the pitch for a brighter sound, and
1203 dramatically slow down the speech for a very relaxed
1204 delivery.</caption>

1205 audiocaps-AudioCaptioning

1206 synthesize A consistent, loud mechanical motor</caption>

1207 AudioSetFullwoAudioMusicCaps-EventClassification

1208 synthesize This is a sound of Speech</caption>

1209

1210 AudioSet-AF

1211 Yo, mind fill in the missing bits in this tune? Please do this
1212 smoothly.]

1213 CREMA-D-P2V

1214 Manifest a voice reproduction in American English verbalizing "The
1215 airplane is almost full.", with the timbre is often monotone and
1216 lacks energy, and it's like a middle-aged who is not Hispanic.

1217 LibriTTS-Clean-100

1218 Stitch up a spiel in en declaring ""But there was a passenger
1219 dropped off for you-a little girl.", with a female speaker with
1220 a moderate pitch and intonation delivers her words quite rapidly
1221 in a confined, clear acoustic environment.

1222 RAVDESS-ChangeEmotion

1223 modulate I want to switch from angry to fearful.</caption> given
1224 example:

1225

1226

1227

1228

1229

1230

1231

1232

1233

1234

1235

1236

1237

1238

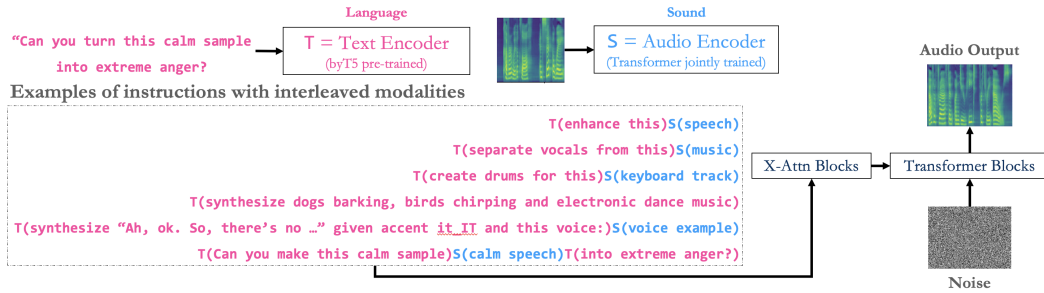
1239

1240

1241

A.2 MODEL AND TRAINING

Model visualization: Figure 5 provides a visual description of *Fugatto*'s architecture and input handling and Algorithm 1 provides a pseudo-algorithm for the optimal-transport conditional flow matching loss. Essentially, it minimizes the mean squared error between the estimator's prediction and a linear interpolation between the data and the Gaussian noise sample used to condition the model on x_t , scaled by $(1 - \sigma)$, where σ is a small enough value.

Figure 5: A description of *Fugatto*'s architecture and input handling.**Algorithm 1** Optimal Transport Conditional Flow Matching Loss Pseudo-Algorithm

- 1: Sample $\mathbf{x}_1 \sim \mathcal{X}$, where \mathcal{X} is the data distribution
- 2: $\sigma \leftarrow 0.01$
- 3: $\mathbf{x}_0 \leftarrow \text{randn_like}(\mathbf{x}_1)$
- 4: $\mathbf{x}_t \leftarrow (1 - (1 - \sigma) \cdot t) \cdot \mathbf{x}_0 + t \cdot \mathbf{x}_1$
- 5: $\mathbf{v}_t \leftarrow \text{estimator}_\theta(\mathbf{x}_t, t;)$
- 6: $\mathbf{u}_t \leftarrow \mathbf{x}_1 - \mathbf{x}_0 \cdot (1 - \sigma)$
- 7: $\text{loss} \leftarrow \text{mse}(\mathbf{v}_t, \mathbf{u}_t)$

Optimal Transport Conditional Flow Matching Pseudo-Algorithm:

1296 **Model hyperparameters:** We provide a list of *Fugatto*'s hyperparameters in Table 10.
 1297

1298 Table 10: Model Hyperparameters for the main *Fugatto* evaluated in this paper.
 1299

| 1300 Hyperparameter | Value |
|---|-------------------|
| 1301 t_schedule | uniform |
| 1302 n_mel_channels | 80 |
| 1303 n_hidden | 1536 |
| 1304 sigma | 0.01 |
| 1305 text_encoder_config.name | google/byt5-large |
| 1306 text_encoder_config.scale | 1.0 |
| 1307 text_encoder_config.n_hidden | 1536 |
| 1308 mel_encoding_strategy | separate |
| 1309 mel_encoder_config.is_causal | false |
| 1310 mel_encoder_config.pos_emb.name | rope |
| 1311 mel_encoder_config.pos_emb.base | 16384 |
| 1312 mel_encoder_config.use_flash_attention | true |
| 1313 mel_encoder_config.deterministic | false |
| 1314 mel_encoder_config.n_layers | 3 |
| 1315 mel_encoder_config.p_dropout | 0.1 |
| 1316 mel_encoder_config.p_dropout_out | 0.0 |
| 1317 mel_encoder_config.n_heads | 16 |
| 1318 mel_encoder_config.has_xattn | false |
| 1319 mel_encoder_config.apply_norm_to_cond | false |
| 1320 mel_encoder_config.layer_norm_method | pre |
| 1321 mel_encoder_config.kernel_size | 3 |
| 1322 mel_encoder_config.use_layer_scale | true |
| 1323 mel_encoder_config.layer_scale_init | 0.1 |
| 1324 mel_encoder_config.layer_scale_decay | 0.95 |
| 1325 mel_encoder_config.d_model | 1536 |
| 1326 decoder_config.d_time | 128 |
| 1327 decoder_config.transformer_hparams.is_causal | false |
| 1328 decoder_config.transformer_hparams.pos_emb.name | rope |
| 1329 decoder_config.transformer_hparams.pos_emb.base | 16384 |
| 1330 decoder_config.transformer_hparams.use_flash_attention | true |
| 1331 decoder_config.transformer_hparams.deterministic | false |
| 1332 decoder_config.transformer_hparams.n_layers | 24 |
| 1333 decoder_config.transformer_hparams.p_dropout | 0.1 |
| 1334 decoder_config.transformer_hparams.n_heads | 16 |
| 1335 decoder_config.transformer_hparams.has_xattn | true |
| 1336 decoder_config.transformer_hparams.kernel_size | 3 |
| 1337 decoder_config.transformer_hparams.context_xattn.n_heads | 16 |
| 1338 decoder_config.transformer_hparams.context_xattn.d_heads | 1536 |
| 1339 decoder_config.d_data | 80 |
| 1340 decoder_config.d_model | 1536 |

1340 **Training and Inference** During the first phase, *Fugatto* is trained on at least 32 NVIDIA A100
 1341 GPU for approximately 1M iterations with template-based instructions and a subset of tasks. During
 1342 the second phase, we restart the optimizer and train for approximately 250k iterations, sampling
 1343 from template-based and free-form instructions uniformly and adding all tasks. We use the AdamW
 1344 optimizer (Loshchilov & Hutter, 2019) with a learning rate of 1e-4, annealing the learning rate to
 1345 1e-6 during the second phase. A G2P model (Bernard & Titeux, 2021) pre-processes the text into
 1346 the International Phonetic Alphabet (IPA) format. During inference, we generate mel-spectrograms
 1347 using 50 function evaluations, 100 in practice, with Heun's Solver and task-specific guidance scale γ .
 1348 Mel-spectrogram to waveform conversion is performed using the pre-trained universal BigVGAN V2
 1349 vocoder, available in the BigVGAN (Lee et al., 2023) repository⁴.

⁴BigVGAN: <https://github.com/nvidia/bigvgan>

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

A.3 ABLATIONS

t-sampling We draw inspiration from (Shen et al., 2018; Valle et al., 2020), where, due to optimization issues, the model would be stuck in a local minima and not make use of the text conditioning variable, rendering the model useless during inference. We believe the same is true with $t \sim \text{sigmoid}(\mathcal{N}(0, 1))$, and that with such a distribution the model is stuck in a local minima and does not leverage the text to minimize the loss.

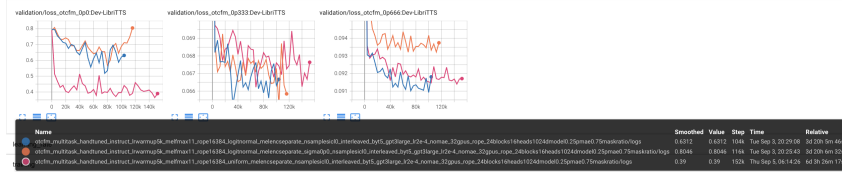


Figure 6: A description of the image.

Model capacity: Exponentially smoothed validation loss per task at different t -values from smaller models with 0.8 B, 1.4 B params, to larger models with 2.5 B parameters.

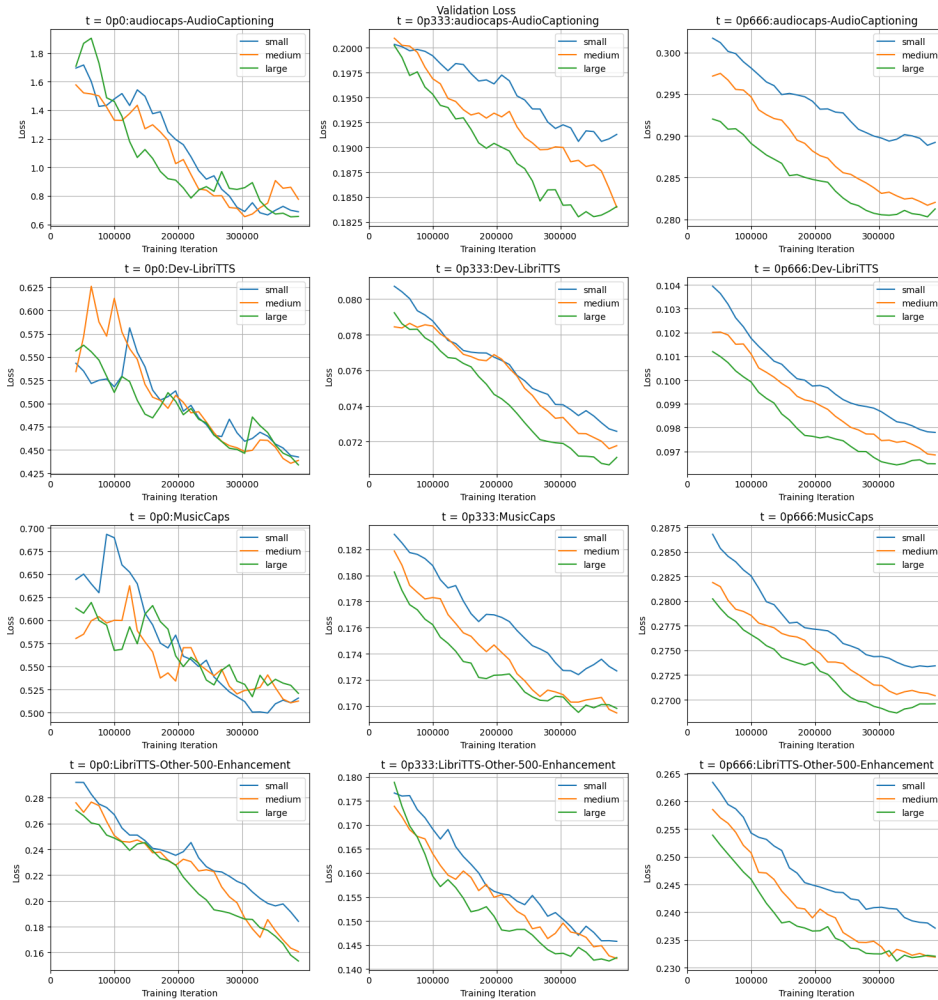


Figure 7: Validation scores for different *Fugatto* sizes on 3 benchmarks at 3 different t -values.

1404 A.4 SINGING VOICE SYNTHESIS MUSIC STYLES AND LYRICS
1405

1406 The Singing-Voice-Synthesis experiments in Section 3.2 evaluates all combinations between the
1407 13 lyrics snippets and 10 music styles below. For each combination, we use the singing-voice-
1408 synthesis instruction generator to create instructions such as: “Showcases a female singer with
1409 an interesting sound, conveys the message through american english lyrics, and infuses country
1410 influences throughout.”

1411 This is the set of 13 lyrics snippets used during evaluation:

1412
1413 "Is this the real life?\nIs this just fantasy?"
1414 "As I walk through the valley of the shadow of death"
1415 "Somebody once told me\nThe world is gonna roll me."
1416 "Look\nIf you had\nOne shot\nOr one opportunity."
1417 "Joy to the world\nThe lord is come."
1418 "Carry on, my wayward son\nThere'll be peace when you are done."
1419 "Please allow me to introduce myself."
1420 "At first I was afraid, I was petrified."
1421 "The world was on fire and no one could save me but you."
1422 "Josie's on a vacation far away."
1423 "She's got a smile that it seems to me."
1424 "She was a fast machine, she kept her motor clean."
1425 "Do you have the time to listen to me whine."

1426 This is the set of 10 music styles used during evaluation:

1427 "Country",
1428 "Electronic",
1429 "Hard Rock",
1430 "Hip-Hop",
1431 "Latin Rock",
1432 "Metal",
1433 "Opera",
1434 "Pop",
1435 "R&B",
1436 "Singer-Songwriter"

1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457

1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511

A.5 MIDI2AUDIO F0 CONTOURS

Figure 8 shows comparisons of 25 monophonic melodies from the test set. As evident from the plot, the model manages to follow the provided MIDI notes and timing well, while inserting nuances, and varying timbres to the performance.

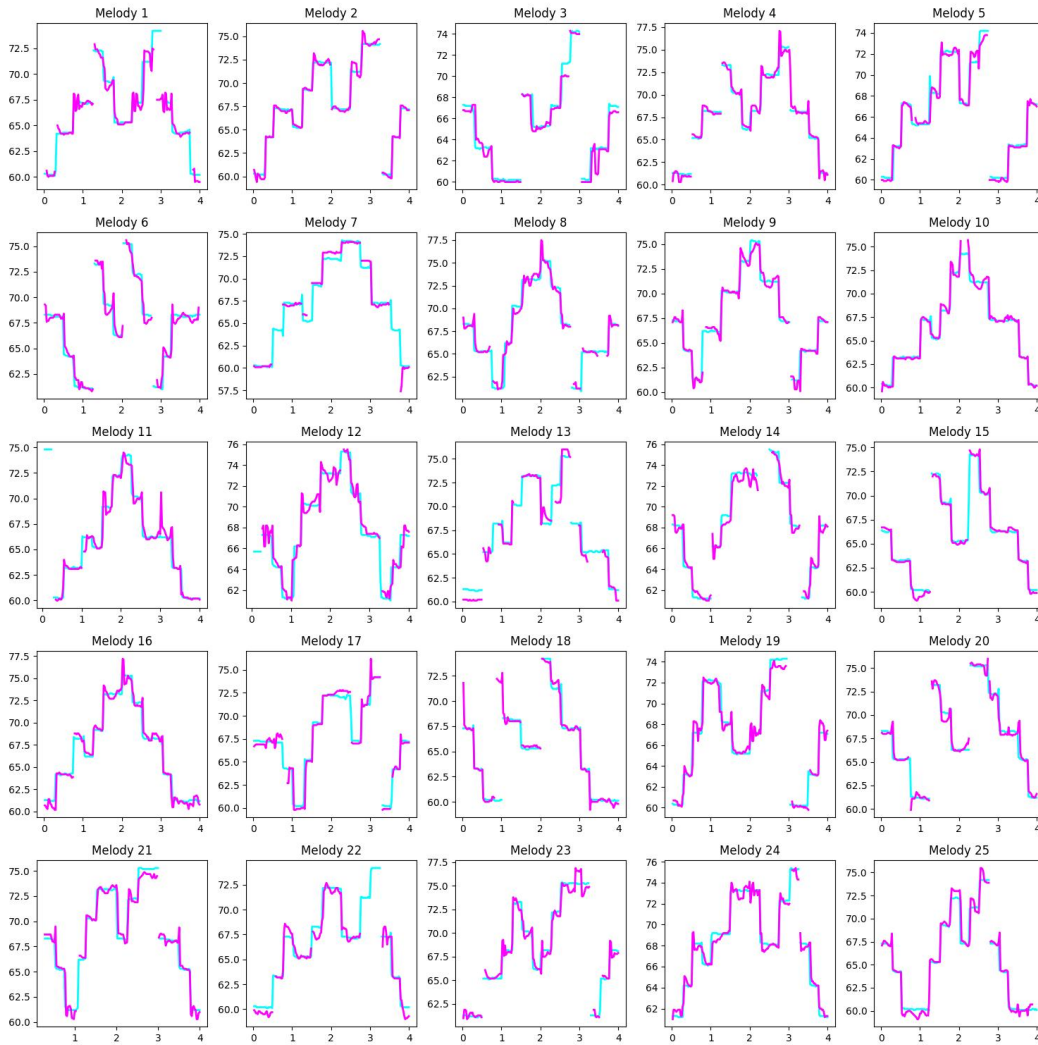


Figure 8: F0 contours of input MIDI (cyan), and generated melodies (magenta). X-axis denotes time in seconds, and Y-axis denotes MIDI pitch.

A.6 COMPOSITIONALITY

Attribute Composition (Audio Events): This is the list of audio events used during Attribute/Event composition in 3.5:

Violin, fiddle; Accelerating, revving, vroom; Water; Acoustic guitar; Afrobeat; Whistle; Air conditioning; Air horn, truck horn; Aircraft; Wind