

Audio is all you need

Adding another modality to our tool box



Task One

A Dataset and Taxonomy for Urban Sound Research

Justin Salamon¹, Christopher Jacoby¹, Juan Pablo Bello¹
¹Music and Audio Research Laboratory, New York University
¹Center for Urban Science and Progress, New York University
 {justin.salamon, cjacoby, jbello}@nyu.edu

ABSTRACT

Automatic urban sound classification is a growing area of research with applications in multimedia retrieval and urban informatics. In this paper we identify two main barriers to research in this area – the lack of a common taxonomy and the existence of large, real-world, annotated data. To address these issues we present a taxonomy of urban sounds and a new dataset, USound, containing 27 hours of audio with 18.5 hours of annotated sound event occurrences across 10 sound classes. The challenge presented by the new dataset are studied through a series of experiments using a baseline classification system.

Categories and Subject Descriptors

H.3.1 [Information Systems]: Content Analysis and Indexing; H.3.3 [Information Systems]: Search and Music Computing

Keywords

Urban sound; dataset; taxonomy; classification

1. INTRODUCTION

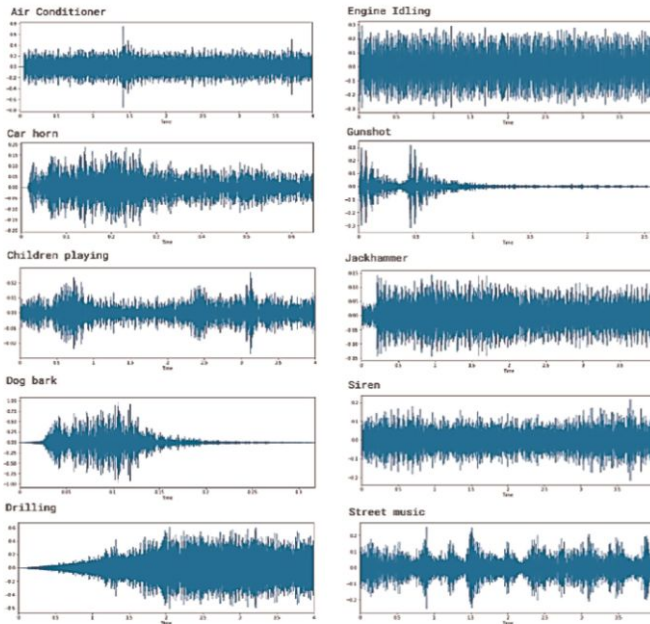
The automatic classification of environmental sound is a growing research field with multiple applications to large-scale, content-based multimedia indexing and retrieval (e.g. [3, 8, 10]). In particular, the audio analysis of urban environments is the subject of increased interest, partly enabled by multimedia sensor networks [15], as well as by large quantities of online multimedia content depicting urban scenes. However, while there is a large body of research in related areas such as speech, music and bioacoustics, work on the analysis of urban acoustic environments is relatively scarce. Furthermore, when existent, it usually focuses on the classification of auditory scene type, e.g. street, park, or airport – to the identification of sound sources in these scenes, e.g. car horn, engine idling, bird tweet. See [5] for an example.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and full citation on the first page. Copyright for components of this work owned by others than the authors is held by their owners. This paper is made available under a CC-BY 4.0 International license. See <http://creativecommons.org/licenses/by/4.0/> for details. This paper is part of the SIGGRAPH Asia 2014 Proceedings, September 28–October 2, 2014, Singapore. Copyright 2014 ACM 978-1-4503-2811-1/14/08...\$10.00. ACM 978-1-4503-2811-1/14/08...\$10.00.

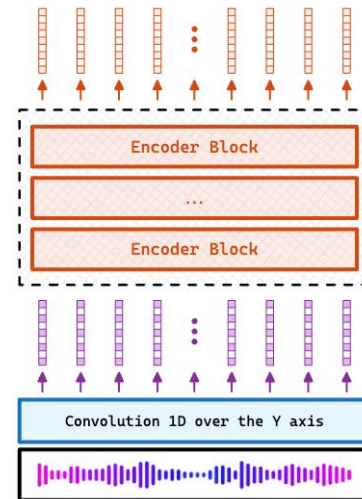
One of the main challenges and hindrances to urban sound research is the lack of labeled audio data. Previous work has focused on audio from carefully produced movies or television tracks [6] from specific environments such as elevator or office space [13, 9], and on commercial or proprietary datasets [11, 4]. The large effort involved in manually annotating real-world data source datasets based on field recordings used to be relatively small (e.g. the event detection dataset of the IEEE AASP Challenge [3] consists of 24 recordings per each of 17 classes). A second challenge faced by the research community is the lack of a common vocabulary when working with urban sounds. This means the classification of sounds into semantic groups may vary from study to study, making it hard to compare results. The goal of this paper is to address the two aforementioned challenges. In Section 2 we propose a taxonomy for urban sound sources to facilitate a common framework for research. Then, in Section 3 we present USound, a dataset of 27 hours of field recordings containing thousands of labeled sound source occurrences. To the best of the authors' knowledge this is the largest free dataset of labeled urban sound events available for research. To understand the complexity and challenges presented by this new dataset, we run a series of baseline sound classification experiments, described in Section 4. The paper concludes with a summary in Section 5.

2. URBAN SOUND TAXONOMY

The taxonomical categorization of environmental sounds, a common first step in sound classification, has been extensively studied in the context of perceptual soundscapes research [14]. Specific efforts to describe urban sounds have often been limited to subsets of broader taxonomies of acoustic environments [2], and thus only partially address the needs of systematic urban sound analysis. For an extensive review of previous work the reader is referred to [12]. In our view, an urban sound taxonomy should satisfy the following three requirements: (1) it should factor in previous research and proposed taxonomies; (2) it should aim to be as detailed as possible, going down to low-level sound sources such as "car horn" (versus "transportation") and "jackhammer" (versus "construction"); (3) it should, in the first instance, focus on sounds that are of specific relevance to urban sound research, such as sounds that contribute to urban noise pollution. To address (1), we decided to base our taxonomy on the effort of [2] dedicated to the urban acoustic environment. We define 6 top-level groups: human, nature, mechanical and music, which are common to most previous



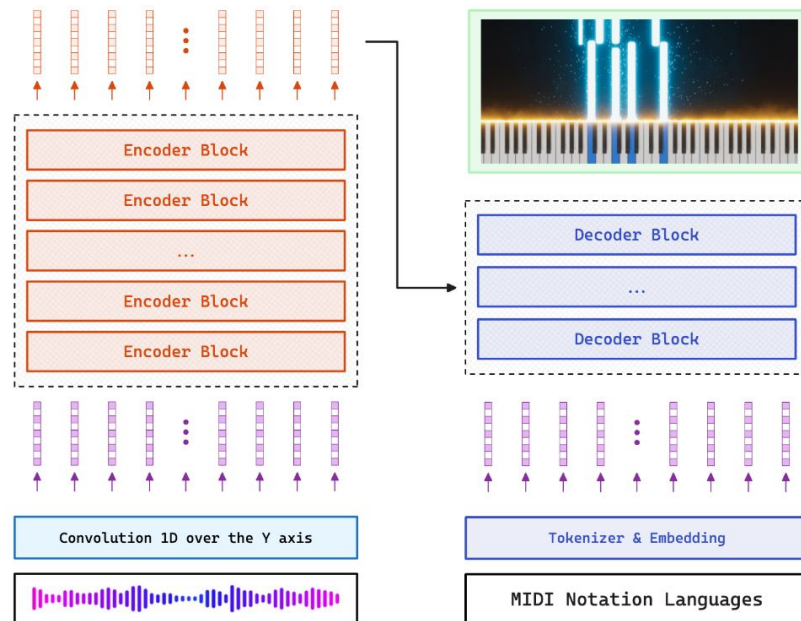
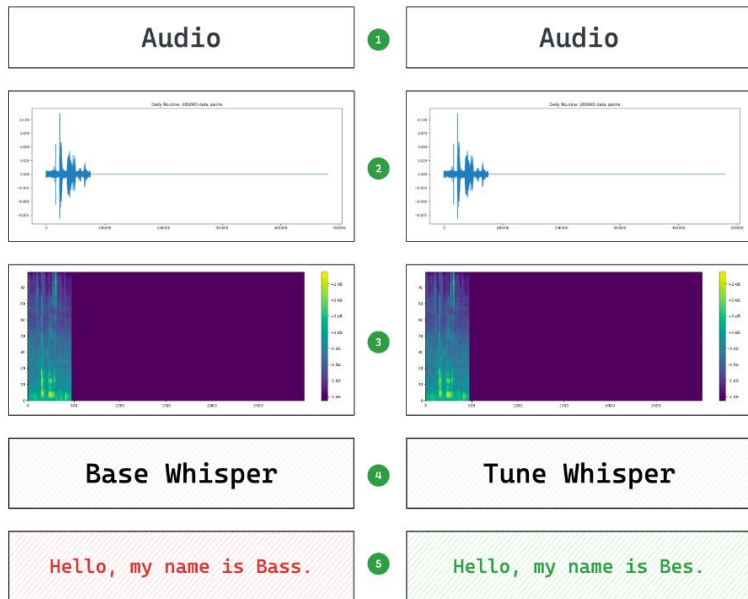
Convolutional Neural Network (CNN)



Salamon et al. 2014

Clearly images! :)

Task Two



Towards Controllable Speech Synthesis in the Era of Large Language Models: A Survey

Tianxin Xie*, Yan Rong*, Pengfei Zhang*, Wenwu Wang, Li Liu

arXiv:2412.06602v2 [cs.CL] 27 Mar 2025

Abstract—Text-to-speech (TTS), also known as speech synthesis, is a prominent research area that aims to generate natural-sounding human speech from text. Recently, with the increasing industrial demand, TTS technologies have evolved beyond synthesizing human-like speech to enabling controllable speech generation. This includes fine-grained control over various attributes of synthesized speech such as emotion, prosody, timbre, and duration. In addition, advancements in deep learning, such as diffusion and large language models, have significantly enhanced controllable TTS over the past several years. In this work, we conduct a comprehensive survey of controllable TTS, covering approaches ranging from basic control techniques to methods utilizing natural language prompts, aiming to provide a clear understanding of the current state of research. We examine the general controllable TTS pipeline, challenges, model architectures, and control strategies, offering a comprehensive and clear taxonomy of existing methods. Additionally, we provide a detailed summary of datasets and evaluation metrics and shed some light on the applications and future directions of controllable TTS. To the best of our knowledge, this survey paper provides the first comprehensive review of emerging controllable TTS methods, which can serve as a beneficial resource for both academic researchers and industrial practitioners.

Index Terms—Text-to-speech, controllable TTS, speech synthesis, TTS survey, large language models, diffusion models.

I. INTRODUCTION

Speech synthesis, also broadly known as text-to-speech (TTS), is a long-time developed technique that aims to synthesize human-like voices from text [1], [2], and it has extensive applications in our daily lives, such as health care [3], [4], personal assistants [5], entertainment [6], [7], and robotics [8], [9]. Recently, TTS has gained significant attention with the rise of large language model (LLM)-powered chatbots, such as ChatGPT [10] and LLaMA [11], due to its naturalness and convenience for human-computer interaction. Meanwhile, the ability to achieve fine-grained control over synthesized speech attributes, such as emotion, prosody, timbre, and duration, has become a hot research topic in both academia and industry, driven by its vast potential for diverse applications.

Deep learning [12] has made great progress in the past decade due to exponentially growing computational resources like GPUs [13], leading to the explosion of numerous exciting works on TTS [14]–[17]. These methods can synthesize human speech with improved quality [14] and can achieve

fine-grained control of the generated voice [18]–[22]. In addition, some recent works synthesize speech given multi-modal input, such as face images [23], [24], cartoons [7], and videos [25]. Moreover, with the fast development of open-source LLMs [11], [26]–[29], some researchers propose to synthesize fine-grained controllable speech with natural language description [30]–[32], offering a new way to generate custom speech voices. Meanwhile, powering LLMs with speech synthesis has also been a hot topic in the last few years [33]–[35]. In recent years, a wide range of TTS methods has emerged, making it essential for researchers to gain a comprehensive understanding of current research trends, particularly in controllable TTS, and to identify promising future directions in this rapidly evolving field. Consequently, there is a pressing need for an up-to-date survey of TTS techniques. While several existing surveys address parametric approaches [36]–[41] and deep learning-based approaches [42]–[48], they largely overlook the controllability of TTS. Additionally, these surveys do not cover recent advancements, such as natural language description-based TTS methods.

This paper provides a comprehensive and in-depth survey of existing and emerging TTS technologies, with a particular focus on controllable TTS methods. Fig. 1 demonstrates the development of controllable TTS methods in recent years, showing their backbones, feature representations, and control abilities. The remainder of this section begins with a brief comparison between this survey and previous ones, followed by an overview of the history of controllable TTS technologies, ranging from early milestones to state-of-the-art advancements. Finally, we introduce the taxonomy and organization of this paper. We have posted a version of our paper on arXiv.org (<https://arxiv.org/abs/2412.06602>).

A. Comparison with Existing Surveys

Several survey papers have reviewed TTS technologies, spanning early approaches from previous decades [36], [37], [40], [49] to more recent advancements [42], [43], [50]. However, to the best of our knowledge, this paper is the first to focus specifically on controllable TTS. The key differences between this survey and prior work are summarized as follows:

Different Scope. Klatt et al. [36] provided the first comprehensive survey on formant, concatenative, and articulatory TTS methods, with a strong emphasis on text analysis. In the early 2010s, Tabet et al. [49] and King et al. [40] explored rule-based, concatenative, and Hidden Markov Models (HMM)-based techniques. Later, the advent of deep learning catalyzed the emergence of numerous neural model-based TTS methods.

Tianxin Xie, Yan Rong, Pengfei Zhang and Li Liu are with the Hong Kong University of Science and Technology (Guangzhou), Guangzhou 511458, China. Wenwu Wang is with Sun Yat-sen University, UY. Corresponding author: Li Liu, liuliu@hkust-gz.edu.cn.
* Equal contribution.
Readers can check this GitHub repository (<https://github.com/taoxu/awesome-controllable-speech-synthesis>) for updates and discussion.

Xie et al. 2025

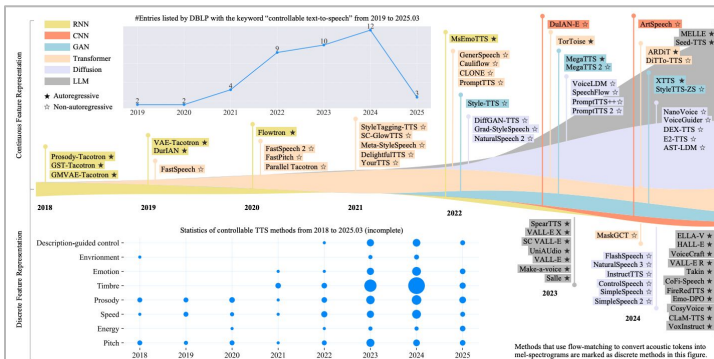


Fig. 1. A summary of representative controllable TTS methods in recent years and their model architectures, feature representations, and control abilities. Additional network structures, such as VAE and flow-based models, are not included in this figure. For more details, refer to Tables IV and III.

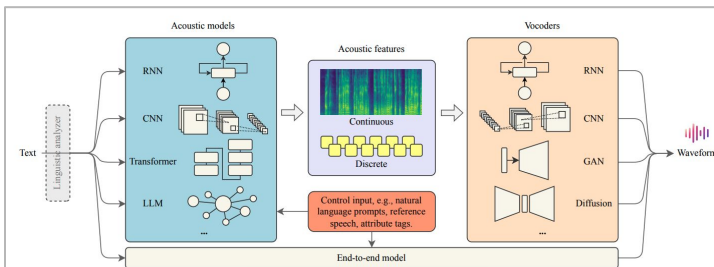
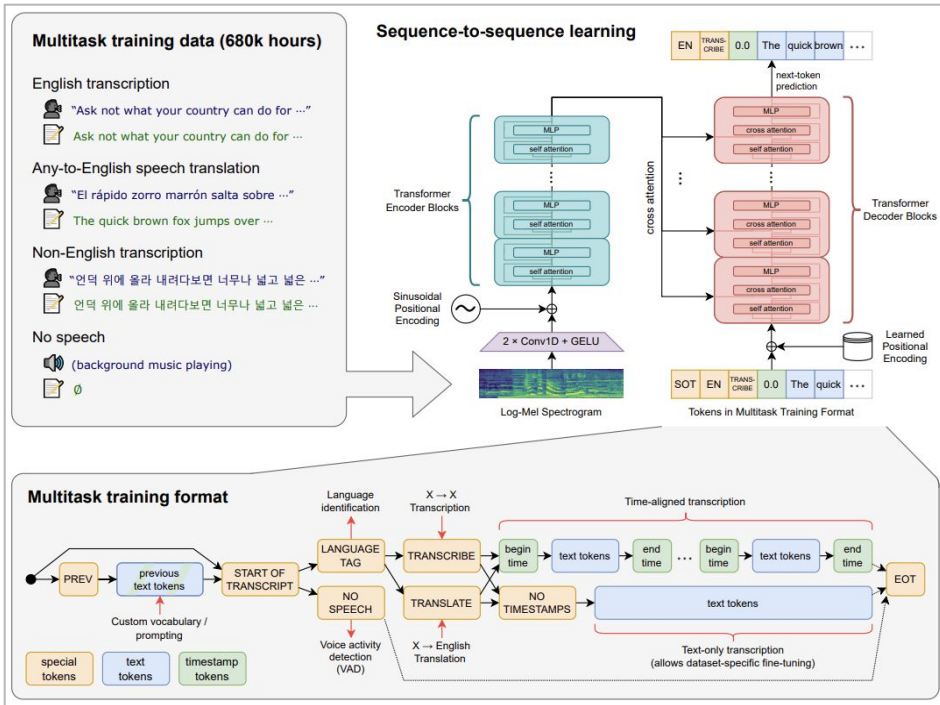


Fig. 2. General pipeline of controllable TTS from the perspective of network structure. Linguistic analysis is necessary for parametric and a few neural methods but is no longer needed for most modern neural methods. In this paper, we only review neural model-based controllable TTS methods and do not investigate acoustic features (e.g., MFCC [107], LSP [108], F0 [109]) used in early TTS methods.



Robust Speech Recognition via Large-Scale Weak Supervision

Alec Radford¹ Jong Wook Kim¹ Tao Xu¹ Greg Brockman¹ Christine McLeavey¹ Ilya Sutskever¹

Abstract

We study the capabilities of speech processing systems trained simply to predict large amounts of transcripts of audio on the internet. When scaled to 680,000 hours of multilingual and multitask supervision, the resulting models generalize well to standard benchmarks and are often competitive with prior fully supervised results but in a zero-shot transfer setting without the need for any fine-tuning. When compared to humans, the models approach their accuracy and robustness. We are releasing models and inference code to serve as a foundation for further work on robust speech processing.

methods are exceedingly adept at finding patterns within a training dataset which boost performance on held-out data from the same dataset. However, some of these patterns are brittle and spurious and don't generalize to other datasets and distributions. In a particularly disturbing example, Radford et al. (2021) documented a 9.2% increase in object classification accuracy when fine-tuning a computer vision model on the ImageNet dataset (Russakovsky et al., 2015) without observing any improvement in average accuracy when classifying the same objects on seven other natural image datasets. A model that achieves "superhuman" performance when trained on a dataset can still make many basic errors when evaluated on another, possibly precisely because it is exploiting those dataset-specific quirks that humans are oblivious to (Ceirios et al., 2020).

1. Introduction

Progress in speech recognition has been energized by the development of unsupervised pre-training techniques exemplified by Wav2Vec 2.0 (Bauvict et al., 2020). Since these methods learn directly from raw audio without the need for human labels, they can predictively use large datasets of unlabeled speech and have been quickly scaled up to 1,000,000 hours of training data (Zhang et al., 2021), far more than the 1,000 or so hours typical of an academic supervised dataset. When fine-tuned on standard benchmarks, this approach has improved the state of the art, especially in a low-data setting.

This suggests that while unsupervised pre-training has improved the quality of audio encoders dramatically, the lack of an equivalently high-quality pre-trained decoder, combined with a recommended protocol of dataset-specific fine-tuning, is a crucial weakness which limits their usefulness and robustness. The goal of a speech recognition system should be to work reliably "out of the box" in a broad range of environments without requiring supervised fine-tuning of a decoder for every deployment distribution.

As demonstrated by Narayanan et al. (2018), Likhomanenko et al. (2020), and Chai et al. (2021) speech recognition systems that are pre-trained in a supervised fashion across many datasets/domains exhibit higher robustness and generalize much more effectively to held-out datasets than models trained on a single source. These works achieve this by combining a many existing high-quality speech recognition datasets as possible. However, there is still only a moderate amount of this data easily available: SpeechShew (Chai et al., 2021) mixes together 7 pre-existing datasets totalling 5,140 hours of supervision. While not insignificant, this is still tiny compared to the previously mentioned 1,000,000 hours of unlabeled speech data utilized in Zhang et al. (2021).

¹Equal contribution ²OpenAI, San Francisco, CA 94110, USA. Correspondence to: Alec Radford <alec@openai.com>, Jong Wook Kim <jongwook@openai.com>.

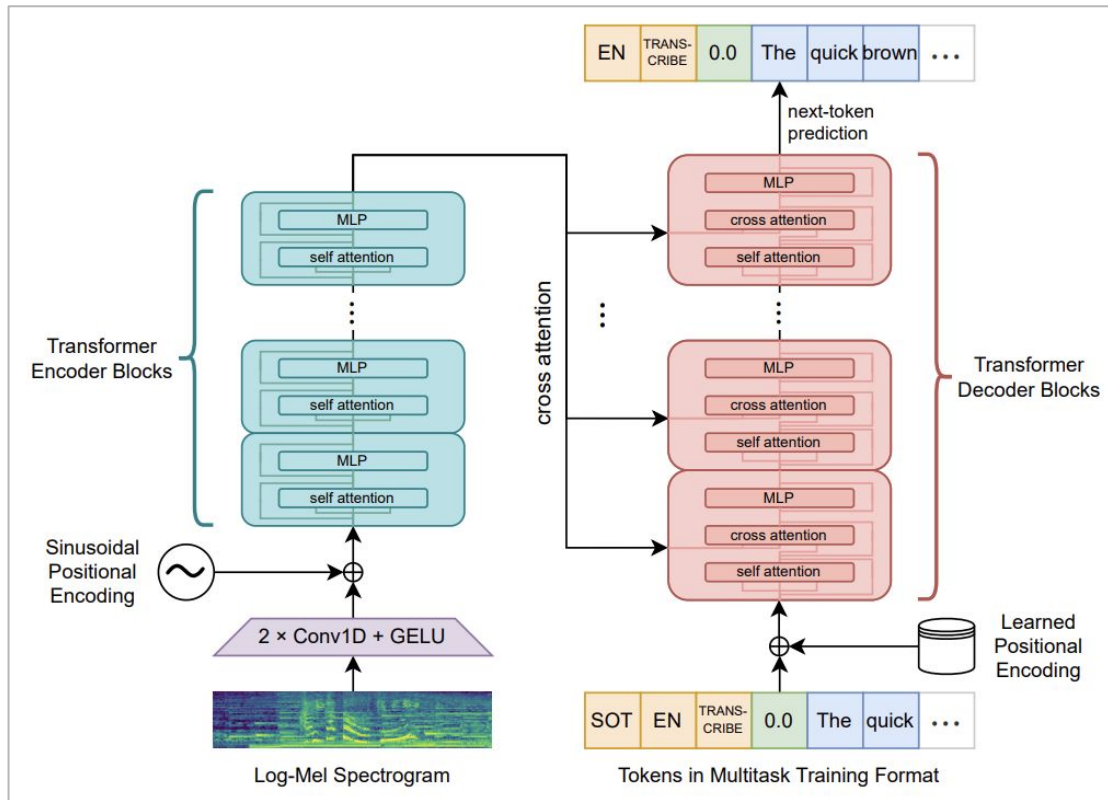
³Bauvict et al. (2021) is an exciting exception - having developed a fully unsupervised speech recognition system

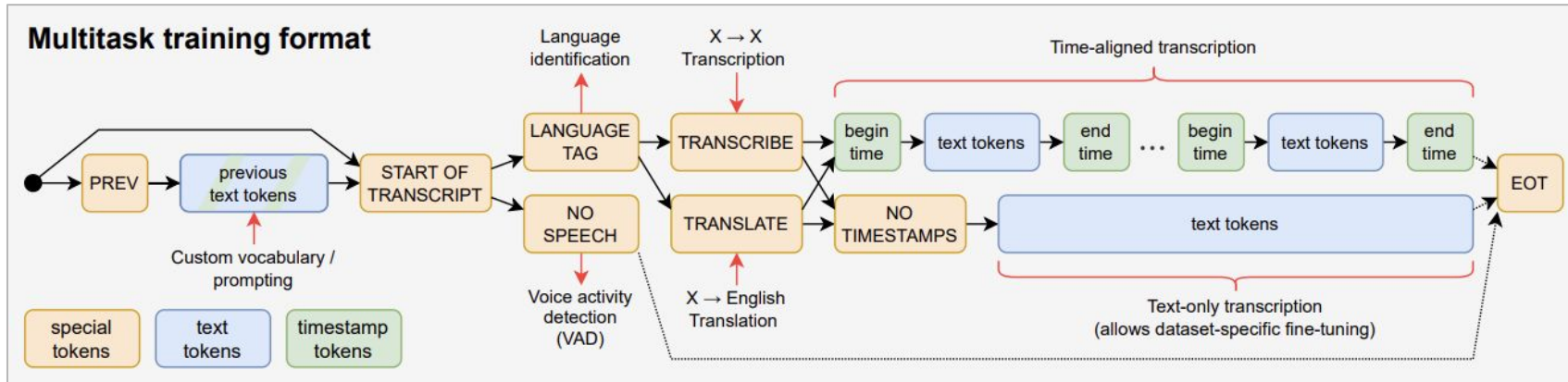
Recognizing the limiting size of existing high-quality supervised datasets, recent efforts have created larger datasets for speech recognition. By relaxing the requirement of gold-standard human-validated transcripts, Chen et al. (2021) and Galvez et al. (2021) make use of sophisticated automated

arXiv:2212.04356v1 [eess.AS] 6 Dec 2022

Radford et al. 2022

Transformer

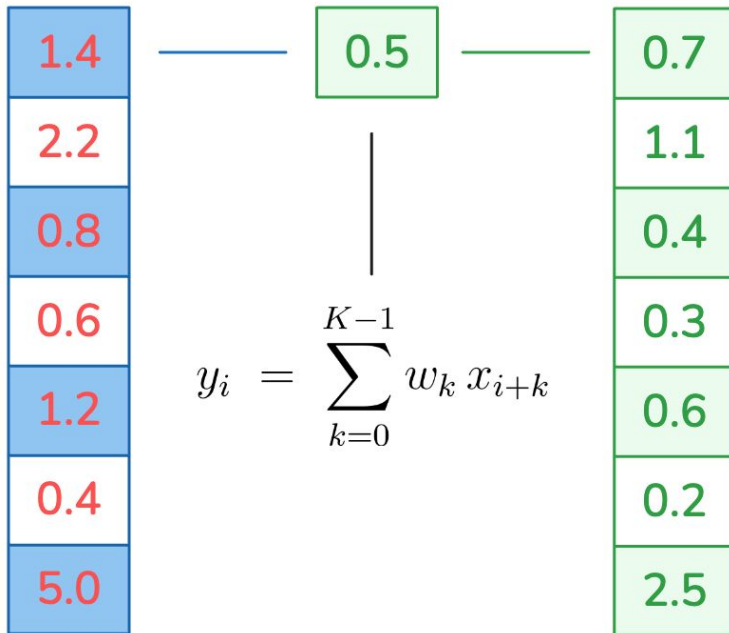




`<|startoftranscript|><|en|><|transcribe|><|notimestamps|>Hello, my name is Bes.<|endoftext|>`

`[50258, 50259, 50359, 50363, 15947, 11, 452, 1315, 307, 8190, 13, 50257]`

Convolutions

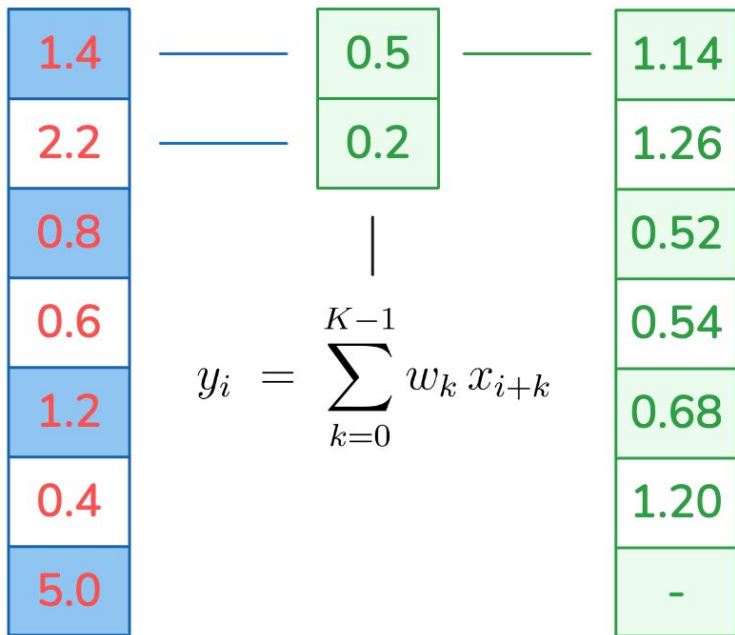


```

1 #
2 #
3 import torch
4
5 #
6 #
7 cov = torch.nn.Conv1d(
8     in_channels=1,
9     out_channels=1,
10    kernel_size=1,
11    padding=0,
12    bias=False
13 )
14
15 #
16 #
17 with torch.no_grad():
18     cov.weight[:] = 0.5
19
20 #
21 #
22 bar = torch.tensor([[1.4, 2.2, 0.8, 0.6, 1.2, 0.4, 5.0]])
23 out = cov(bar)
24 print("Out:", out)
25

```


Convolutions

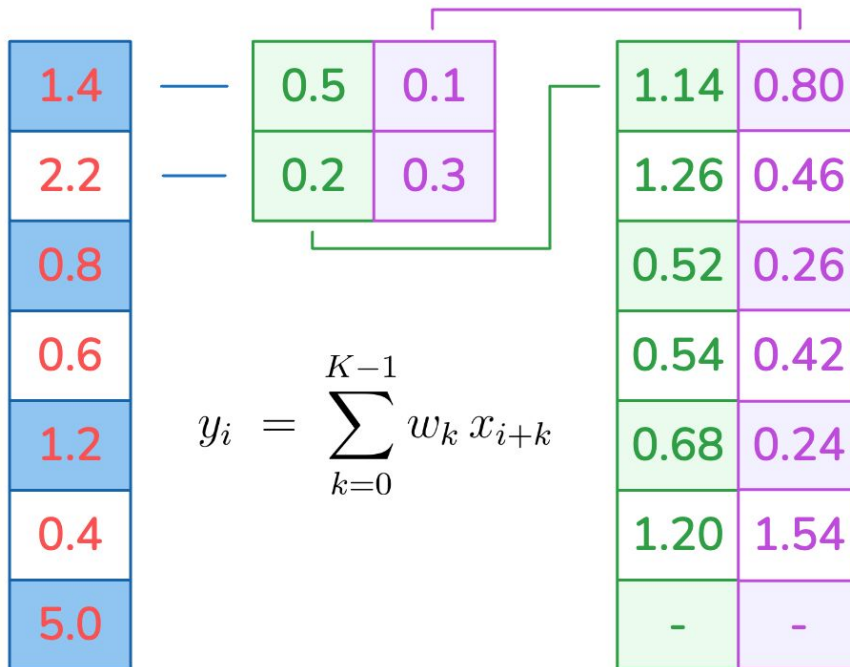


```

1 #
2 #
3 import torch
4
5 #
6 #
7 cov = torch.nn.Conv1d(
8     in_channels=1,
9     out_channels=1,
10    kernel_size=2,
11    padding=0,
12    bias=False
13 )
14
15 #
16 #
17 with torch.no_grad():
18     cov.weight[:] = torch.tensor([[0.5, 0.2]])
19
20 #
21 #
22 bar = torch.tensor([1.4, 2.2, 0.8, 0.6, 1.2, 0.4, 5.0])
23 out = cov(bar)
24 print("Out:", out)
25

```

Convolutions

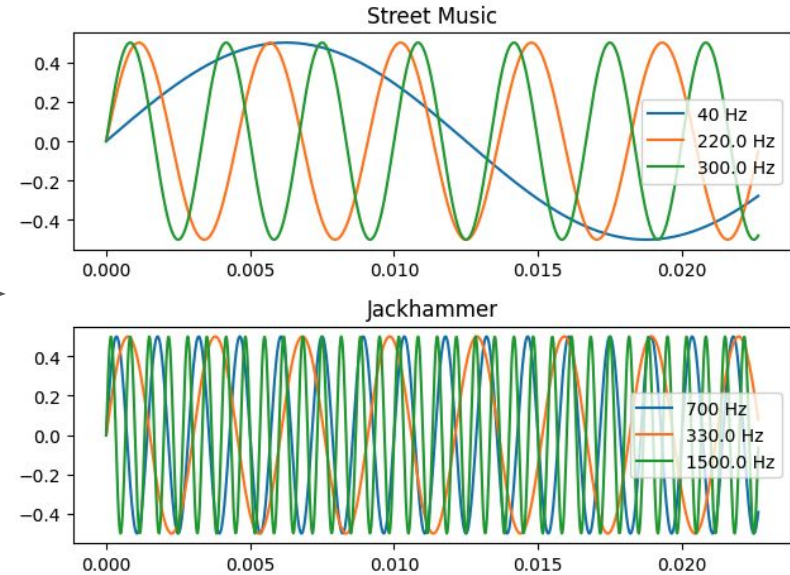
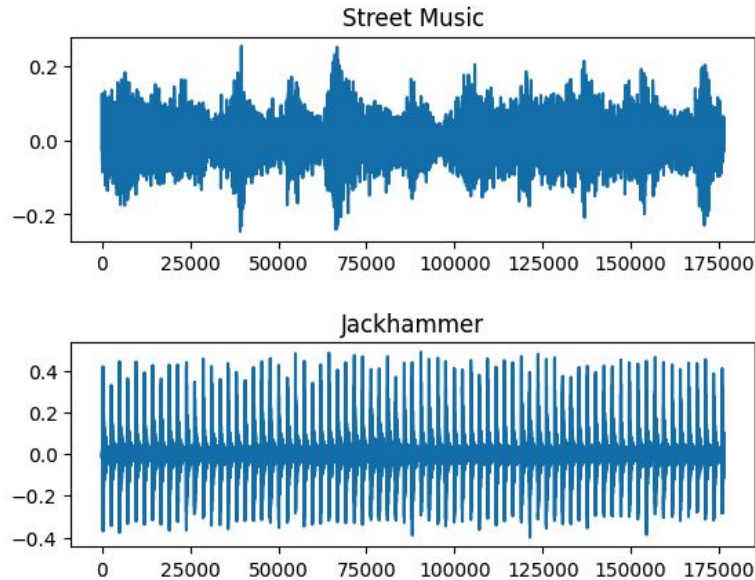


```

1 #
2 #
3 import torch
4
5 #
6 #
7 cov = torch.nn.Conv1d(
8     in_channels=1,
9     out_channels=2,
10    kernel_size=2,
11    padding=0,
12    bias=False
13 )
14
15 #
16 #
17 with torch.no_grad():
18     cov.weight[0, 0] = torch.tensor([0.5, 0.2])
19     cov.weight[1, 0] = torch.tensor([0.1, 0.3])
20
21 #
22 #
23 bar = torch.tensor([1.4, 2.2, 0.8, 0.6, 1.2, 0.4, 5.0])
24 out = cov(bar)
25

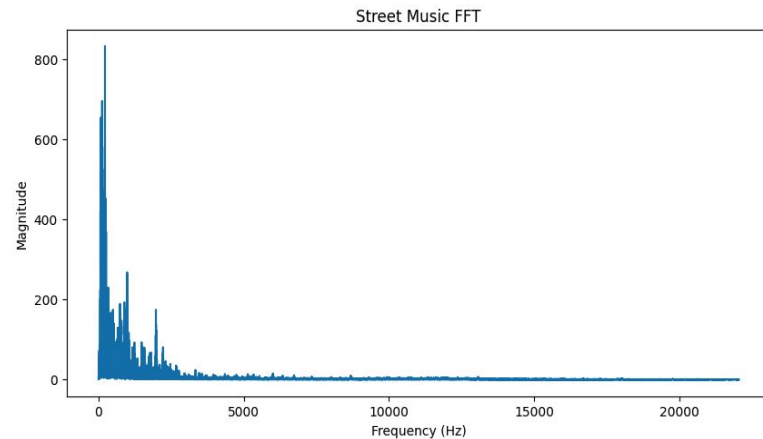
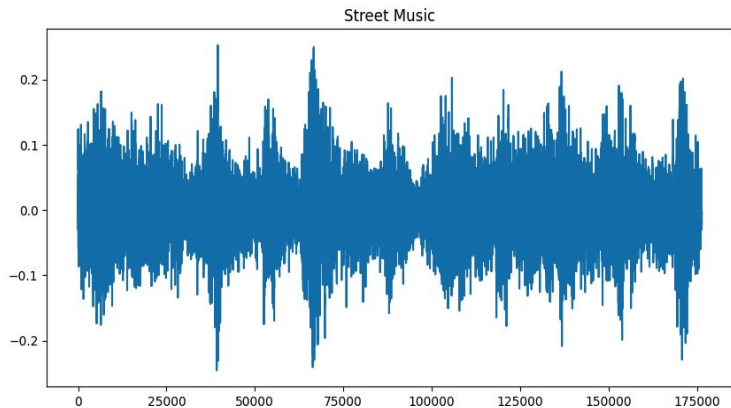
```

What is audio?



- Sound is a bundle of vibrations. It is complex
- As usual we need to extract features

Fourier Transform



- Fourier Transform unbundles audio by going through each frequency to test how much each frequency is contributing to the overall signal
- In practice we use something called the Discrete Fourier Transform

Tune Code

```

1 #
2 #
3 #
4 import torch
5 import whisper
6
7 #
8 #
9 #
10 #
11 model = whisper.load_model('tiny')
12 audio = whisper.load_audio('name.wav')
13 audio = whisper.pad_or_trim(audio)
14 lg_ml = whisper.log_mel_spectrogram(audio)
15 tknsr = whisper.tokenizer.get_tokenizer(multilingual=True)
16
17 #
18 #
19 #
20 #
21 opt = whisper.DecodingOptions()
22 res = whisper.decode(model, lg_ml.to(model.device), opt)
23 print('Baseline:', res.text) # Hello my name is Bass.
24 print('-----')
25
26 #
27 #
28 #
29 #
30 ids = []
31 ids += [tknsr.sot]
32 ids += [tknsr.language_token]
33 ids += [tknsr.transcribe]
34 ids += [tknsr.no_timestamps]
35 ids += tknsr.encode(' Hello, my name is Bes.')
36 ids += [tknsr.eot]
37
38 #
39 #
40 #
41 #
42 optimizer = torch.optim.Adam(model.parameters(), lr=0.00001)
43 criterion = torch.nn.CrossEntropyLoss()
44
45 #
46 #
47 #
48 #
49 model.train()
50 tks = torch.tensor(ids).unsqueeze(0).to(model.device)
51 mel = whisper.log_mel_spectrogram(audio).unsqueeze(0).to(model.device)

```

```

52 #
53 #
54 #
55 #
56 #
57 pred = model(tokens=tks, mel=mel)
58 trgt = tks[:, 1:].contiguous()
59 pred = pred[:, :-1, :].contiguous()
60
61 #
62 #
63 #
64 #
65 print('Ids Target:', trgt.squeeze().tolist())
66 print('Ids Output:', torch.argmax(pred, dim=-1).squeeze().tolist())
67 print('Txt Target:', tknsr.decode(trgt.squeeze().tolist()))
68 print('Txt Output:', tknsr.decode(torch.argmax(pred, dim=-1).squeeze().tolist()))
69
70 #
71 #
72 #
73 #
74 loss = criterion(pred.transpose(1, 2), trgt)
75 print('Loss:', loss.item())
76 print('-----')
77 optimizer.zero_grad()
78 loss.backward()
79 optimizer.step()
80
81 #
82 #
83 #
84 #
85 model.eval()
86 prd = model(tokens=tks, mel=mel)
87 prd = prd[:, :-1, :].contiguous()
88
89 #
90 #
91 #
92 #
93 print('Ids Target:', trgt.squeeze().tolist())
94 print('Ids Output:', torch.argmax(prd, dim=-1).squeeze().tolist())
95 print('Txt Target:', tknsr.decode(trgt.squeeze().tolist()))
96 print('Txt Output:', tknsr.decode(torch.argmax(prd, dim=-1).squeeze().tolist()))
97 loss = criterion(prd.transpose(1, 2), trgt)
98 print('Loss:', loss.item())
99
100 #
101 #
102 #
103 #
104

```


Tune Result

```
00: torch.Size([1, 80, 3000])
```

```
Baseline: Hello, my name is Bass.
```

```
=====
```

```
00: torch.Size([1, 80, 3000])
```

```
Ids Target: [50259, 50359, 50363, 2425, 11, 452, 1315, 307, 8190, 13, 50257]
```

```
Ids Output: [50259, 50359, 50363, 2425, 11, 452, 1315, 307, 29626, 13, 50257]
```

```
Txt Target: <len|><|transcribel><|notimestamps|> Hello, my name is Bes.<|endoftext|>
```

```
Txt Output: <len|><|transcribel><|notimestamps|> Hello, my name is Bass.<|endoftext|>
```

```
Loss: 0.5395039916038513
```

```
=====
```

```
00: torch.Size([1, 80, 3000])
```

```
Ids Target: [50259, 50359, 50363, 2425, 11, 452, 1315, 307, 8190, 13, 50257]
```

```
Ids Output: [50259, 50359, 50363, 2425, 11, 452, 1315, 307, 8190, 13, 50257]
```

```
Txt Target: <len|><|transcribel><|notimestamps|> Hello, my name is Bes.<|endoftext|>
```

```
Txt Output: <len|><|transcribel><|notimestamps|> Hello, my name is Bes.<|endoftext|>
```

```
Loss: 0.1766674667596817
```

Good luck!

Recurrent Rebels

Ardrit
Yurii
Daniel
Kori

Perceptron Party

Kenton
Maxime
Guillaume
Josh

Activation Aces

Artemis
Amy
David
Gaurav

Gradient Giggers

Milo
Ollie
Dimitar
Nnamdi

Backprop Bunch

Evelyn
Dimitris
Stanley
Filippo

Overfitting Overlords

James
Pry
Liam
Andrea

Dropout Disco

Loredana
Aygün
Neville
Coline