

Generating Data with Text-to-Speech and Large-Language Models for Conversational Speech Recognition

Samuele Cornell^{*1}, Jordan Darefsky^{*2}, Zhiyao Duan², Shinji Watanabe¹

¹Carnegie Mellon University, USA

²University of Rochester, USA

samuele.cornell@ieee.org, jdarefsk@u.rochester.edu

Abstract

Currently, a common approach in many speech processing tasks is to leverage large scale pre-trained models by fine-tuning them on in-domain data for a particular application. Yet obtaining even a small amount of such data can be problematic, especially for sensitive domains and conversational speech scenarios, due to both privacy issues and annotation costs. To address this, synthetic data generation using single speaker datasets has been employed. Yet, for multi-speaker cases, such an approach often requires extensive manual effort and is prone to domain mismatches. In this work, we propose a synthetic data generation pipeline for multi-speaker conversational ASR, leveraging a large language model (LLM) for content creation and a conversational multi-speaker text-to-speech (TTS) model for speech synthesis. We conduct evaluation by fine-tuning the Whisper ASR model for telephone and distant conversational speech settings, using both in-domain data and generated synthetic data. Our results show that the proposed method is able to significantly outperform classical multi-speaker generation approaches that use external, non-conversational speech datasets. **Index Terms:** generative synthetic data, multi-talker speech recognition, text-to-speech, conversational speech processing

1. Introduction

Current robust speech processing methods are considerably data hungry. For example, state-of-the-art automatic speech recognition (ASR) systems require tens or even hundreds of thousands of hours of training data in order to achieve enough robustness in different domains [1–3]. Such a vast amount of training data is leveraged either explicitly by training from scratch on a large amount of data or implicitly by fine-tuning/adapting a pre-trained “foundation” model that was originally trained, in a supervised or unsupervised manner [1, 4–6], on a large dataset.

Nevertheless, for some domains, obtaining even a small portion of in-domain supervised data for fine-tuning can be problematic due to potential privacy concerns or prohibitive expense.

This is especially true for sensitive application scenarios, including medical, government, and law enforcement settings. Moreover, due to increasing regulatory attention, even scaling in-domain training data is potentially becoming more difficult.

Aside from privacy issues, applications that require recordings with multiple speakers are also inherently difficult, time-consuming and costly to annotate and thus obtain in scale. Prominent examples are meeting scenarios [7, 8] including doctor-patient recordings, speech captioning, speech analytics and so on.

^{*}These authors contributed equally to this work.

Despite the difficulties associated with obtaining data for multi-speaker scenarios, there are speech processing approaches that require multi-speaker conversational data for training. These approaches have proven to be effective as demonstrated in recent speech processing challenges [7–9]. Prominent examples are end-to-end neural diarization (EEND) and most target speaker voice activity detection (TS-VAD) approaches [10–14], as well as multi-speaker ASR [15–19]. Lack of annotated in-domain conversational data at scale is a significant issue for these techniques, which is only partly mitigated by leveraging foundation models [17–19]. Consequently, many of these approaches have to rely on synthetic data to increase dataset size. This is commonly achieved by artificially overlapping clips from existing datasets and adding noise and reverberation.

While several toolkits have been proposed to ease the workload [20, 21], creating synthetic datasets remains more art than science, as it often requires extensive hand-tuning, domain knowledge, heuristics, and significant trial and error. Crucially, this process is highly prone to the introduction of unwanted biases in the resulting dataset, leading to a performance drop due to domain mismatch [12].

The aforementioned difficulties motivate the development of more automated, machine learning based approaches for synthetic data creation. Several methods have in fact explored this direction, primarily focusing on improving ASR performance by leveraging synthetic data created with text-to-speech (TTS) models [22–27, 27–31] or leveraging ASR and TTS cycle-consistency during training [32, 33] for semi-supervised training. However, these approaches focus on single-speaker scenarios and thus cannot be directly applied to domains where multi-speaker conversational ASR is required. In parallel, recent works [34, 35] on speech summarization and audio captioning have shown how large-language models (LLM)s can be leveraged effectively for synthetic audio data augmentation.

Building upon this previous research, in this work we explore using TTS models along with LLMs to generate multi-speaker conversational data. We focus on two-speaker ASR on real-world telephone (Fisher [36]) and distant speech recognition settings (Mixer 6 Speech [37]) by fine-tuning Whisper [1]. The contributions of this work are the following: 1) We propose a synthetic data generation pipeline for conversational ASR using LLMs for content generation and a conversational multi-speaker TTS model for speech generation; 2) We perform a systematic investigation on the use of synthetic data for training multi-speaker ASR models with three different approaches: using “classical” LibriSpeech based multi-speaker simulation, using a conventional state-of-the-art (SotA) TTS model, and using a recently proposed conversational TTS model [38].

2. Method under study

Our approach is summarized in Figure 1. We explore the use of a pre-trained chat-optimized LLM for creating short conversation transcripts between two participants from scratch for when in-domain conversational transcriptions are not available or would be costly to obtain. Specifically, we use the recently released Llama 3 8B Instruct model and few-shot prompt it with 8 text-prompt examples randomly selected from a 1000-example subset of Spotify Podcasts dataset [39] used to train Parakeet (the text data was transcribed using Whisper-D, described in [38]). That is, for each new example we want to generate, we randomly select a subset of eight text samples from our Parakeet subset to use as the few-shot prompt. This procedure could also be used to augment existing in-domain text-only data. It could also be worth exploring fine-tuning on in-domain data instead of prompting.

These LLM obtained transcripts are then used to generate synthesized speech through a multi-speaker TTS model. The resulting data, consisting of ground truth multi-speaker transcripts and the synthesized multi-speaker mixture can then be used for training or fine-tuning purposes, i.e. in Sec. 4 for adapting Whisper to perform multi-speaker ASR.

2.1. Conversational TTS generation

The effectiveness of this approach will heavily depend on the capability of the TTS model used. While we expect LLMs will be proficient in generating conversational transcripts as shown in previous work on summarization [34], most TTS models are not capable of synthesizing multi-speaker conversational data. Although one could naively generate each speaker’s utterances independently and then stitch them together, such an approach would fail to capture real conversational speech turn-taking dynamics and para-linguistic subtleties such as changes in intonation, etc., and would therefore potentially introduce a domain mismatch in the generated audio.

Recently, in [38] a conversational TTS model, Parakeet, has been proposed. Parakeet’s training dataset includes 60,000 hours of Spotify Podcasts data, much of which is multi-speaker. It is therefore able to directly generate two-speaker short conversations of up to 30 seconds when given a text prompt in the style of the one in Figure 1, i.e. with speaker-id related tags [S1] and [S2]. We use a diffusion version of Parakeet that, similar to [40] autoregressively generates blocks of continuous latents using latent diffusion on each block. The autoencoder is trained to map 44,100 Hz audio to 16-channel dimensional latents, with a time downsampling factor of 1024. Each diffusion block consists of 128 (time-wise) latent vectors, which correspond to approximately three seconds of audio.

LLM-generated transcripts and speech examples are available online¹.

3. Experimental setup

3.1. Evaluation data

In this work, we focus on two-speaker multi-speaker conversational ASR. This focus is due to the limitations of Parakeet, whose generations tend to lose correctness as the number of unique speakers in the text prompt increases. Furthermore, we also consider scenarios with relatively high signal-to-noise ratio (SNR); tackling more complex settings such as CHiME-6 [7] requires modeling of background noise and dynamic acoustic

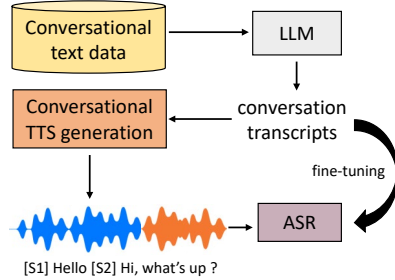


Figure 1: Block diagram of the proposed approach.

conditions (as the participants move, reverberation can change significantly). We thus perform our experiments using two conversational speech datasets with these characteristics: Fisher Corpus (both Part 1 and Part 2) and Mixer 6 Speech.

3.1.1. Fisher

Fisher consists of 11699 telephone conversations between two English speakers sampled at 8 kHz. Each conversation is around 10 minutes long. We use the train, validation, and test split from [41] (11577, 61 and 61 conversations of respectively 1960 h, 7 h and 7 h). The Fisher recordings originally separate each of the speakers into different channels; however, since our focus is on general single-channel conversational speech processing, we mixdown the two channels to mono. We also resample the signal to 16 kHz as we use Whisper which was trained on 16 kHz data (see Sec. 3.3).

3.1.2. Mixer 6 Speech

As an additional scenario, we consider Mixer 6 Speech. Specifically we use the version re-annotated for the CHiME-7 challenge [8]. It consists of two-speaker interviews of approximately 15 minutes (sampled at 16 kHz) recorded by 14 different far-field recording devices. In this work we only use recordings from the tabletop microphone device (CH04). We use the splitting from [8], where full long-form annotation is only available for the development (59 interviews, 15 h) and evaluation sets (23 interviews, 13 h). Here we further split the development set into an adaptation portion and a validation portion of respectively 2:30 h and 4 h after discarding utterance groups longer than 30 s as done in [19]. This further split allows us to compare the use of synthetic data versus in-domain data for fine-tuning.

3.2. Baseline Methods

3.2.1. NeMo multi-speaker simulation tool

We consider two baseline methods. The first method we consider is a “classical” synthetic speech generation method, where single speaker speech from one high quality speech dataset (e.g. LibriSpeech [42]) is used to construct conversation-style synthetic recordings by artificially overlapping single speaker utterances and contaminating them by adding noise, artificial room impulse response (RIR) or other transforms (e.g. clipping, microphone transfer function etc.). We make use of the SoTA NeMo multi-speaker simulation tool [21] (NeMo MSS in the following). We use LibriSpeech train-clean 360 and 100 portions and generate 100 h of short conversations between two speakers of up to 30 seconds in length. For Mixer 6 Speech experiments, we additionally use the built-in RIR simulation in order to generate simulated far-field speech.

¹popcornell.github.io/SynthConvASRDemo

3.2.2. xTTS-v2

The second baseline method we consider is the approach outlined in Section 2, where a standard TTS model is used to generate the training data. We explore this using the Coqui xTTS-v2 model [43] (denoted simply as xTTS in Sec. 4) In detail, for each utterance group in the training dataset (either LLM-generated or taken from a text-only corpus) we sample two speaker ids from LibriSpeech train-clean 360 and 100 and then two corresponding LibriSpeech enrollment utterances to condition xTTS-v2 for the generated TTS id. We then generate each utterance in the utterance group independently via xTTS-v2 and truncate excessive leading and trailing silence regions using Silero VAD [44]. The generated audio is then resampled to 16 kHz and mixed together by randomly adding start time offsets based on the order of the sentences in the utterance group transcript, ensuring that utterances from the same speaker do not overlap.

3.3. ASR System

In our experiments, which focus on two-speaker conversational speech, we use the method proposed in [19] where Whisper [1] is adapted to perform multi-speaker ASR through fine-tuning with a serialized output training (SOT) [15] objective on utterance groups. This approach aligns with common practices in the field where a model pre-trained on a large amount of data (i.e. a foundation model) is fine-tuned/adapted for a particular domain or application of interest.

Compared to [19], in our experiments we focus only on standard SOT without considering timestamps and use only Whisper medium. We use low-rank adapters (LoRA) [45] while the rest of the model is kept frozen. During each fine-tuning experiment a linear warm-up schedule is employed for the first N epoch, then the learning rate is linearly decayed over a maximum of 20 epochs. The L^2 norm of the gradients is clipped to 5. One LoRA adapter for each linear layer in the model (i.e. for each query, key, value and feed-forward network layer) is used. For each adapter we set the LORA rank to 64, alpha to 128, and dropout to 0.1. In our preliminary experiments on the full Fisher training set, we found that this configuration yields the best results, even when compared to fine-tuning the entire model. If validation loss does not improve for 2 consecutive epochs the training is stopped. We tune the batch size, number of warm-up epochs (N) and the value of the maximum learning rate for each set of experiments. Parakeet synthesized audio is resampled to 16 kHz in our experiments. In Fisher experiments, for all synthetic data, we use on-the-fly resampling to simulate telephone 3400 Hz band-limiting.

In Mixer 6 experiments, only for xTTS and Parakeet, we contaminate the data with reverberation using random RIRs obtained from [46]. This of course is less realistic than the RIR simulation used in NeMo MSS as the RIR is the same for both speakers. We make our fine-tuning code publicly available².

3.4. Evaluation Setup

For each dataset, we run our experiments using the same setup as in [19], where oracle voice activity detection (VAD) is used and the dataset is divided into several utterance groups [3, 19]. Continuing to follow [19], we then perform evaluation for each utterance group independently and accumulate word error rate (WER) statistics over the whole dataset (insertions, deletions

etc.). We choose this evaluation method because we only focus on multi-speaker ASR, and an evaluation which considers the whole conversation (e.g. as in CHiME-6/7) would require a diarization component, which would add significant complexity.

We thus consider concatenated minimum permutation WER (cpWER) [7]. This is the same as WER in [19], with the best permutation evaluated for each utterance group independently. We also consider multi-input multi-output WER (MIMO-WER), which is more tolerant than cpWER to speaker assignment errors. We use the Meeteval toolkit [47] to compute both scores. Whisper text normalization is used both during training and scoring.

4. Experiments

4.1. Fisher

In Table 1 we report results obtained on the Fisher test set as defined in Sec. 3.1.1 with different data used for fine-tuning. As a baseline, in the first row, we report the results with no adaptation. In the second panel, we report results on in-domain Fisher training data adaptation. We observe only a modest difference between using the full training set or a 80 h data subset, which is likely because we are leveraging a strong pre-trained model. In the third and fourth panels, we report results obtained with synthetic data approaches. In particular, for the two TTS approaches (xTTS and Parakeet), we consider two opposite situations: a best-case/oracle scenario where we use in-domain conversation transcriptions and another one where we suppose we have none and thus we use as input Llama-3 random generated utterance groups transcripts (LLM_{rnd}) as described in Sec. 2.

We observe that xTTS-based generation outperforms NeMo MSS when Fisher only transcriptions (Fisher) are used. When LLM generated transcriptions are used (LLM_{rnd}), xTTS performance is on par/slightly worse than NeMo MSS. In contrast, when using Parakeet, the difference between using LLM generated transcripts versus the Fisher training set transcriptions is modest, and interestingly, the generated transcripts yield the best performance. In general, while the performance gain compared to the baseline synthetic data approaches (xTTS and NeMo MSS) is significant, there remains a substantial gap compared to using in-domain data (Fisher). It appears that this gap cannot be bridged solely by scaling the amount of synthetic data.

In Figure 2 we report cpWER on Fisher for different amounts of adaptation data, both from Fisher training set and from synthetic approaches. For modest amounts of data (less than 5 h) the proposed approach is competitive to using in-domain data; however, as the amount of adaptation data is scaled, performance saturates quickly: The improvement between 50 h and 5 h is marginal when compared to the one afforded by using in-domain data. This trend is also observed for the other synthetic data approaches and suggests that there is some inherent mismatch in all of the synthetic data approaches tested that prevents effective scaling. At least for Parakeet, results suggest that this mismatch seems to be more related to the signal/acoustic content rather than the transcription semantic content as the gap between using Fisher transcriptions and LLM-generated transcription is modest.

4.2. Mixer 6 Speech

In Table 2, we show results obtained on Mixer 6. The trends observed are consistent with the Fisher experiments, despite the rather naive artificial reverberation strategy used for xTTS

²github.com/popcornell/ASRLightingFT

Table 1: Multi-speaker ASR results on Fisher test set with different adaptation data.

Adaptation Data	amount (hours)	cpWER (%)	MIMO-WER (%)
-	0	44.94	26.15
Fisher	1960	13.76	13.58
Fisher	80	15.43	14.94
NeMo MSS	80	34.37	26.51
xTTS (Fisher)	80	24.88	24.07
xTTS (LLM _{rnd})	80	34.65	28.31
Parakeet (Fisher)	80	21.44	21.00
Parakeet (LLM _{rnd})	80	20.41	19.48
Parakeet (LLM _{rnd})	160	19.93	19.45

■ Fisher ■ xTTS ■ Parakeet+LLM_{rnd}
■ NeMo Sim ■ Parakeet

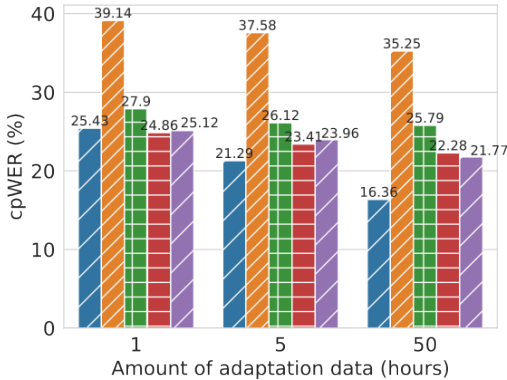


Figure 2: Multi-speaker ASR results on Fisher test set for different adaptation data sources and quantity.

and Parakeet experiments. This confirms that the proposed approach can also be effective for far-field multi-speaker synthetic data, at least when compared to the classical approach (NeMo MSS results) and when available in-domain data is very scarce (here 2:30h). Parakeet (LLM_{rnd}, 80h) also compares favorably with the third and fourth rows, where we report the results of using the Fisher full 1960h training set and a 80h subset respectively for adaptation. For these Fisher experiments, to reduce the mismatch due to the telephone lower sampling frequency, we apply telephone band-limiting to Mixer 6 in the inference phase. We also contaminate the Fisher 6 training data with reverberation as done for Parakeet and xTTS as described in Sec. 3.3.

4.3. Further discussion & remarks

Considering both Fisher and Mixer 6 experiments, the fact that Parakeet+LLM_{rnd} improves considerably over NeMo MSS while xTTS fails suggests that turn-taking and para-linguistics may play a considerable role for multi-talker ASR.

Finally, for both Mixer 6 Speech and Fisher scenarios, we tried using 50h of synthetic LLM_{rnd} data to augment a portion of in-domain data (5h and 50h) by mixing the two or by training on synthetic data and then fine-tuning on in-domain data. However, in most instances, this approach does not result in any

Table 2: Multi-speaker ASR results on Mixer 6 Speech eval set with different adaptation data.

Adaptation Data	amount (hours)	cpWER (%)	MIMO-WER (%)
-	0	43.67	32.16
Mixer6	2.30	20.36	19.77
Fisher	1960	20.83	20.33
Fisher	80	22.12	21.36
NeMo MSS	80	36.71	28.21
xTTS (Mixer6)	2.30	25.99	24.47
xTTS (LLM _{rnd})	80	35.65	30.18
Parakeet (Mixer6)	2.30	23.52	22.82
Parakeet (LLM _{rnd})	2.30	23.70	22.12
Parakeet (LLM _{rnd})	80	21.25	20.17

improvement over using solely the in-domain data; in the xTTS and NeMo MSS cases we even observe performance degradation. For example, by combining 50h of Parakeet (LLM_{rnd}) and 50h of original Fisher training data the model achieved a cpWER of 15.74% which is only marginally better than the 16.36% obtained with only 50h of Fisher (Figure 2). Interestingly, negligible or no improvement was also observed when the in-domain data was more modest (5h). This may be due to the fact that we are leveraging a strong pre-trained model, and thus the quality of adaptation data rather than quantity matters most. Future work should explore adaptation of the TTS model to generate synthetic audio that better matches distribution of in-domain data.

5. Conclusions

In this work, we study the use of synthetically generated data for multi-speaker ASR, focusing on the two-speaker case. We explore different strategies of generating synthetic data, comparing artificially overlapped data and SotA conventional TTS models with a novel conversational TTS model, Parakeet, capable of natively generating multi-speaker utterances. Our results show that our approach using Parakeet significantly outperforms previous SotA multi-speaker simulation techniques. Furthermore, when in-domain data is limited to only a few hours, our approach achieves performance reasonably close to that of using in-domain data; however, when more in-domain data is available, our approach lags behind using real data. For Mixer 6, our approach also obtains results comparable to using external real-world multi-speaker data (Fisher). Overall, our experiments suggest that the LLM generated transcripts are reliable but that there is currently a performance gap compared to using in-domain audio data (when enough in-domain data exists).

Limitations of our work include that we only consider two-speaker conversational speech, short 30-second conversations, and relatively high SNR scenarios. These constraints were primarily imposed by the current limitations of the Parakeet TTS model, and thus improvement of TTS capabilities is crucial to increasing synthetic data viability. For example, to tackle more complex noisy/reverberant scenarios, the TTS model needs to incorporate acoustic scenario modeling, e.g. via acoustic style transfer techniques or even few-shot adaptation on some in-domain data (e.g. via [48]). Another possible limitation is that Parakeet itself is trained on text-audio pairs where the text is “synthetic”, i.e. Whisper-D [38] is used to generate multi-

speaker transcriptions for Spotify podcast audio which is then used for Parakeet training. Since Whisper-D is fine-tuned from Whisper using a small number of annotated multi-speaker examples (and Whisper itself is likely trained on a sizeable quantity of multi-speaker data), there is an indirect but somewhat circular dependency on the existence of ground-truth annotations. Also, Parakeet’s weakness in generating consistent 3/4-speaker conversational data could in part be due to limitations of Whisper-D. Future work could potentially explore the joint bootstrapping of audio-to-text and text-to-audio models.

6. Acknowledgments

S. Cornell was supported by IC Postdoctoral Research Fellowship Program at CMU via ORISE through an agreement between U.S. DoE and ODNI. We’d like to thank Google’s TPU Research Cloud (TRC), which provided compute for generating synthetic Parakeet samples and Llama synthetic text utterances. Our work would not have been possible without their support.

7. References

- [1] A. Radford *et al.*, “Robust speech recognition via large-scale weak supervision,” in *ICML*. PMLR, 2023.
- [2] Y. Peng *et al.*, “Reproducing whisper-style training using an open-source toolkit and publicly available data,” in *Proc. of ASRU*. IEEE, 2023.
- [3] N. Kanda *et al.*, “Large-scale pre-training of end-to-end multi-talker asr for meeting transcription with single distant microphone,” *arXiv preprint arXiv:2103.16776*, 2021.
- [4] A. Baeviski *et al.*, “wav2vec 2.0: A framework for self-supervised learning of speech representations,” *Advances in neural information processing systems*, vol. 33, 2020.
- [5] W.-N. Hsu *et al.*, “Hubert: Self-supervised speech representation learning by masked prediction of hidden units,” *IEEE/ACM TASLP*, vol. 29, 2021.
- [6] S. Chen *et al.*, “WavLM: Large-scale self-supervised pre-training for full stack speech processing,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 16, no. 6, 2022.
- [7] S. Watanabe *et al.*, “CHiME-6 challenge: Tackling multispeaker speech recognition for unsegmented recordings,” in *CHiME Workshop*, 2020.
- [8] S. Cornell *et al.*, “The CHiME-7 DASR challenge: Distant meeting transcription with multiple devices in diverse scenarios,” *CHiME Workshop*, 2023.
- [9] N. Ryant *et al.*, “The third dihard diarization challenge,” *Proc. of Interspeech*, 2021.
- [10] Y. Fujita *et al.*, “End-to-end neural speaker diarization with self-attention,” in *Proc. of ASRU*. IEEE, 2019.
- [11] K. Kinoshita, M. Delcroix, and N. Tawara, “Integrating end-to-end neural and clustering-based diarization: Getting the best of both worlds,” in *Proc. of ICASSP*. IEEE, 2021.
- [12] F. Landini *et al.*, “From simulated mixtures to simulated conversations as training data for end-to-end neural diarization,” 2022.
- [13] I. Medennikov *et al.*, “Target-speaker voice activity detection: a novel approach for multi-speaker diarization in a dinner party scenario,” *Proc. of Interspeech*, 2020.
- [14] N. Tawara *et al.*, “Ntt speaker diarization system for chime-7: multi-domain, multi-microphone end-to-end and vector clustering diarization,” *CHiME Workshop*, 2023.
- [15] N. Kanda *et al.*, “Serialized output training for end-to-end overlapped speech recognition,” *Proc. of Interspeech*, 2020.
- [16] —, “Investigation of end-to-end speaker-attributed asr for continuous multi-talker recordings,” in *Proc. of SLT*. IEEE, 2021.
- [17] Z. Huang *et al.*, “Adapting self-supervised models to multi-talker speech recognition using speaker embeddings,” in *Proc. of ICASSP*. IEEE, 2023.
- [18] S. Cornell *et al.*, “One model to rule them all? towards end-to-end joint speaker diarization and speech recognition,” in *Proc. of ICASSP*. IEEE, 2024.
- [19] C. Li *et al.*, “Adapting multi-lingual asr models for handling multiple talkers,” *Proc. of Interspeech*, 2023.
- [20] T. Cord-Landwehr *et al.*, “Mms-msg: A multi-purpose multi-speaker mixture signal generator,” in *Proc. of IWAENC*. IEEE, 2022.
- [21] T. J. Park *et al.*, “Property-aware multi-speaker data simulation: A probabilistic modelling technique for synthetic data generation,” *Proc. of Interspeech*, 2023.
- [22] A. Rosenberg *et al.*, “Speech recognition with augmented synthesized speech,” in *Proc. of ASRU*. IEEE, 2019.
- [23] Z. Chen *et al.*, “Improving speech recognition using gan-based speech synthesis and contrastive unspoken text selection,” in *Proc. of Interspeech*, 2020.
- [24] N. Rossenbach *et al.*, “Generating synthetic audio data for attention-based speech recognition systems,” in *Proc. of ICASSP*. IEEE, 2020.
- [25] A. Tjandra, S. Sakti, and S. Nakamura, “Machine speech chain,” *IEEE/ACM TASLP*, vol. 28, 2020.
- [26] A. Fazel *et al.*, “SynthASR: Unlocking synthetic data for speech recognition,” *Proc. of Interspeech*, 2021.
- [27] X. Zheng *et al.*, “Using synthetic audio to improve the recognition of out-of-vocabulary words in end-to-end asr systems,” in *Proc. of ICASSP*. IEEE, 2021.
- [28] S. Ueno *et al.*, “Data augmentation for asr using tts via a discrete representation,” in *Proc. of ASRU*. IEEE, 2021.
- [29] T.-Y. Hu *et al.*, “Synt++: Utilizing imperfect synthetic data to improve speech recognition,” in *Proc. of ICASSP*. IEEE, 2022.
- [30] M. Soleymanpour *et al.*, “Synthesizing dysarthric speech using multi-speaker tts for dysarthric speech recognition,” in *Proc. of ICASSP*. IEEE, 2022.
- [31] E. Casanova *et al.*, “Asr data augmentation in low-resource settings using cross-lingual multi-speaker tts and cross-lingual voice conversion,” in *Proc. of Interspeech*, 2023.
- [32] T. Hori *et al.*, “Cycle-consistency training for end-to-end speech recognition,” in *Proc. of ICASSP*. IEEE, 2019.
- [33] M. K. Baskar *et al.*, “Eat: Enhanced asr-tts for self-supervised speech recognition,” in *Proc. of ICASSP*. IEEE, 2021.
- [34] J.-w. Jung *et al.*, “Augsumm: towards generalizable speech summarization using synthetic labels from large language model,” *Proc. of ICASSP*, 2024.
- [35] S.-L. Wu *et al.*, “Improving audio captioning models with fine-grained audio features, text embedding supervision, and llm mix-up augmentation,” in *Proc. of ICASSP*. IEEE, 2024.
- [36] C. Cieri, D. Miller, and K. Walker, “The Fisher corpus: A resource for the next generations of speech-to-text,” in *LREC*, 2004.
- [37] L. Brandschain *et al.*, “The Mixer 6 corpus: Resources for cross-channel and text independent speaker recognition,” in *LREC*, 2010.
- [38] J. Darefsky, G. Zhu, and Z. Duan, “Parakeet,” 2024. [Online]. Available: <https://jordandarefsky.com/blog/2024/parakeet/>
- [39] A. Clifton *et al.*, “100,000 podcasts: A spoken english document corpus,” in *Proceedings of the 28th International Conference on Computational Linguistics*, 2020.
- [40] Z. Liu *et al.*, “Autoregressive diffusion transformer for text-to-speech synthesis,” *arXiv preprint arXiv:2406.05551*, 2024.
- [41] G. Morrone *et al.*, “End-to-end integration of speech separation and voice activity detection for low-latency diarization of telephone conversations,” *Speech Communication*, 2024.

- [42] V. Panayotov *et al.*, “Librispeech: an asr corpus based on public domain audio books,” in *Proc. of ICASSP*, 2015.
- [43] E. Casanova *et al.*, “Xtts: a massively multilingual zero-shot text-to-speech model,” *arXiv e-prints*, 2024.
- [44] S. Team, “Siltero vad: pre-trained enterprise-grade voice activity detector (vad), number detector and language classifier,” <https://github.com/snakers4/siltero-vad>, 2021.
- [45] E. J. Hu *et al.*, “Lora: Low-rank adaptation of large language models,” *arXiv preprint arXiv:2106.09685*, 2021.
- [46] T. Ko *et al.*, “A study on data augmentation of reverberant speech for robust speech recognition,” in *Proc. of ICASSP*. IEEE, 2017.
- [47] T. von Neumann *et al.*, “MeetEval: A toolkit for computation of word error rates for meeting transcription systems,” *CHiME Workshop*, 2023.
- [48] L. Zhang, A. Rao, and M. Agrawala, “Adding conditional control to text-to-image diffusion models,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023.