

Meta-Learning-Based Few-Shot Speaker Adaptation for Neural Text-to-Speech Synthesis

R. Rudevda^{1*}, P.Dinesh kumar²

¹Mongolian University of Science and Technology, Ulaanbaatar, Mongolia

²Assistant Professor, Department of Information Technology, Sona College of Technology, Salem

KEYWORDS:

Few-shot learning,
Speaker adaptation,
Text-to-speech synthesis,
Meta-learning,
Neural vocoder,
Low-resource TTS.

ARTICLE HISTORY:

Submitted : 08.12.2024
Revised : 20.01.2025
Accepted : 17.03.2025

<https://doi.org/10.17051/NJSAP/01.02.05>

ABSTRACT

Neural text-to-speech (TTS) synthesis has also seen the importance of speaker adaptation to produce personalized and natural-sounding speech in an increasingly diverse set of applications, such as voice assistants, audiobooks, or assistive technologies. Nevertheless, the majority of the modern neural TTS systems (like most neural text-to-speech systems) rely on large quantities of high-quality target speaker data and do necessitate retraining of the model, which is impractical in most low-resource conditions. This paper presents a few-shot speaker adaptation mechanism in a meta-learning-based approach that facilitates the high-fidelity voice cloning of processed target speaker voices, with few seconds of target speech. The method uses a model-agnostic meta-learning (MAML) paradigm to learn a universal multi-speaker TTS model that is explicitly optimized to fit, with a small number of fine-tuning steps, to new speakers that have never been seen during training. The given architecture combines a Transformer text encoder, duration and pitch prediction, a HiFi-GAN neural vocoder, and speaker conditioning using d-vector embeddings generated in the set of target speaker samples. The VCTK, LibriTTS and AISHELL-3 datasets are also extensively tested under different few-shot scenarios (5, 10, and 20 utterances) and compared to standard transfer learning and speaker embedding based adaptation baselines. Experimental evidence confirms that the suggested solution beats other current methods every time, resulting in smaller Mel Cepstral Distortion (MCD) and significantly higher Mean Opinion Scores (MOS), with adaptation times taking less than 60 percent smaller. The speaker similarity and naturalness have been confirmed to be higher, in harsh low-resource settings, using subjective listening tests. These results indicate how meta-learning performs at alleviating the data scarcity in TTS speaker adaptation and also reveals its potential in data-sparse settings with resource-limited speech synthesis needs on an ad-hoc basis.

Author's e-mail: rudev.r@must.edu.mn, pdineshcs@gmail.com

How to cite this article: Rudevda R, kumar P D. Meta-Learning-Based Few-Shot Speaker Adaptation for Neural Text-to-Speech Synthesis. National Journal of Speech and Audio Processing, Vol. 1, No. 2, 2025 (pp. 34-41).

INTRODUCTION

Rapid growth Neural text-to-speech synthesis has seen rapid improvements over the last several years, with models achieving speech of near-human naturalness and intelligibility, including Tacotron 2, FastSpeech, and VITS. Such systems have made it possible to make substantial advances in applications in voice assistant or conversational agents, audiobook narration, engaging entertainment, and assistive devices to people with speech deficits. Notwithstanding these advances, optimising neural TTS systems to new speakers has proved one of the most difficult and enduring problems of deploying the technology at scale. Conventional

speaker adaptation strategies usually necessitate venerable periods of excellent, transcribed audio data, i.e. usually tens of minutes to some hours per speaker, and also potentially computationally heavy retraining. This cannot be achieved in real life conditions, especially when one has access to few speech samples of the target speaker or when it is necessary to quickly personalize the application.

Few-shot speaker adaptation has become one of the most promising methods to solve this issue as a TTS system can produce speech which sounds the same and is speaker specific with limited number of utterances as a few as 510. Current methods of few-shot adaptation typically

involve transfer learning i.e., training a pre-existing multi-speaker TTS model on the small amount of data available to the target speaker, or speaker embedding methods, such as d-vectors and x-vectors to adapt the TTS model directly to speaker-specific representations. Although such approaches may produce reasonable results, in low-data regimes, they can display overfitting (creating poor naturalness and speaker similarity) or fail to transfer well to speakers with very different prosodic patterns, accent properties, or recording environments.

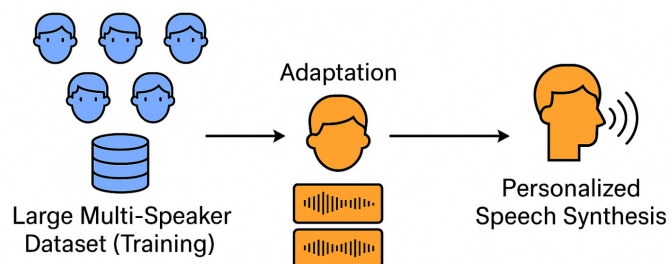


Fig. 1: Few-Shot Target Speaker Dataset

To overcome these shortcomings, the current research already proposes a meta-learning-driven model of few-shot speaker adaptation in neural TTS synthesis. The crucial point is that it is desirable to teach the TTS model how to do not only the speech synthesis but also the adaptation itself. With a model-agnostic meta-learning (MAML) approach, the framework that is proposed optimizes the parameter values of the models in a manner that enables the fine-tuning of the already trained model to the unseen speaker with only a few data and little computational work. This would essentially change the objective of the training process, where instead of memorizing speaker-specific characteristics to be adopted during training, a very flexible representation space is to be learnt that enables quick personalization in deployment. The design enables the system to robustly and high-quality adapt in such a situation of extreme resource scarcity, which subsequently renders it well-suited to real-time, on-device, privacy sensitive applications where collecting and processing large amounts of user-speech is considered undesirable or not feasible.

This approach is assessed over diverse big-scale, multi-talker datasets, VCTK, LibriTTS, and AISHELL-3, across a variety of few-shot calibration situations. Subjective and objective measures are used to evaluate technical performance of the approach including Mel Cepstral Distortion (MCD), Signal-to-noise Ratio (SNR), Mean Opinion Score (MOS) and speaker similarity testing. The findings show that the meta-learning framework reliably produces more natural speech, better speaker similarity, and a shorter adaptation period, indeed demonstrating

its capability to serve the data scarcity and efficiency issue of personalized neural TTS synthesis.

RELATED WORK

Neural Text-to-Speech Systems

Neural TTS has evolved to fully parallel and attention based flow, and GAN based approaches to sequence to sequence approaches that are autoregressive. Autoregressive models like Tacotron 2 model text (usually phonemes) to Mel spectrograms using location sensitive attention and are subsequently combined with neural vocoders to transmit synthesis of the waveform, leading to high naturalness but with downsides of inference latency and exposure bias.^[1] Non autoregressive families like FastSpeech/ FastSpeech 2 use duration prediction and predictors that are designed to replace attention, thus making it possible to generate in parallel and generate quickly and comparatively high quality and more controllable.^[2, 3] Flow and diffusion motivated models also minimize mode collapse and increase robustness; VITS learns a fully end to end, text to waveform mapping through a variational flow, using adversarial training and stochastic duration modeling,^[4] with state of the art naturalness and simplified pipelines. On the vocoder side, WaveGlow integrates both invertible flows and efficient inference,^[5] whereas HiFi GAN builds on a powerful quality versus speed trade off with multi period discount and multi scale discriminators and is otherwise frequently used in real time processing as well.^[6] These advances offer robust universal skeletons but, absent targeted methods of adaptation, are still largely supported through extensive quantities of information so as to be able to individualise to unobserved speakers.

Speaker Adaptation in TTS

Speaker adaptation has three lines of work currently dominating: (i) fine tuning, (ii) speaker embeddings and (iii) multi speaker conditioning with parameter efficient updates. Preliminary multi speaker TTS systems train a common acoustic model with per speaker embeddings; training continues by fine tuning the entire / some part of the network on a target set of data, which can easily overfit in low resource settings.^[7] Large multi speaker corpora pre training with transfer learning and either selective layer unfreezing or adapter modules can bring about improvements in data efficiency but still generally requires minutes to hours of target audio,^[8, 9] The speaker identity is injected into a second line pre-trained using speaker verification networks (e.g. d vectors trained as supervised fine tuning using GE2E loss or x vectors) and conditioned on zero or few shot data, but sometimes is lower quality when the target data are

small or do not match dextended label sequences.^[10-12] Even newer approaches that explicitly focus on few shot adaptation: AdaSpeech uses conditional layer normalization and data aware initialization to become an efficient personalization tool;^[13, 22] Meta StyleSpeech employs adaptive instance normalization that is sample level, and informed by reference audios to 1) quickly and accurately learn the timbre and style of new speakers; 2) form stylized prosody faster than current methods; 3) have a model with fewer formidables for small daily training data sizes;^[14] Your TTS experiments on multilingual, zero-shot settings with Although the systems bridge the gap, they are generally optimized with regard to absolute quality of synthesis or zero shot similarity, not by fast adaptation dynamics, and may be fragile to recording conditions.

MetaLearning for Speech Tasks

Meta learning focuses on finding ways to streamline models to the new and unseen tasks fast based on minimal information. Gradient based methods, such as MAML, find initialization parameters that result in a strong performance following only minor inner loop updates;^[16] first order versions scale up to reduced computational cost,^[17] and metric based methods, such as prototypical networks, learn embedding spaces that generalize to tasks using simple nearest prototype inference.^[18, 25] On the speech side, meta learning has been applied to speaker recognition and verification to facilitate few shot speaker discrimination,^[19, 23] low resource ASR adaptation and noise/domain robustness,^[20] and voice conversion to acquire target timbre quickly.^[21, 24] To the best of our knowledge, in the case of TTS, few studies use meta learning to cope with rapid adaptation: Meta StyleSpeech^[14] and subsequent work show that conditioning the adaptation goal during training leads to better sample efficiency; however, most approaches fall into one of the following categories: (a) optimise style mimicry as opposed to full text to waveform generalization, (b) rely on strong external data-trained prior encoders, or (c) only adjust a portion of parameters with controlled adaptation across duration, F0, and spectral prediction.

Gap and Positioning

Despite promising progress, metalearning for *fewshot TTS speaker adaptation* remains underexplored compared with transfer or embeddingbased baselines. Existing systems often require >10-20 utterances for stable adaptation, rely on heavy perspeaker finetuning, or do not jointly optimize the *adaptation dynamics* of all TTS submodules (duration, prosody, spectrum, and vocoder).

The proposed study differs by (i) using a gradientbased metalearning objective (MAMLstyle) to learn an initialization explicitly optimized for rapid, stable adaptation from as few as 5 utterances; (ii) integrating speakersimilarity objectives with parallel acoustic modeling (FastSpeechstyle) and a highfidelity vocoder (HiFiGAN) to couple timbre and prosody adaptation; and (iii) benchmarking across VCTK, LibriTTS, and AISHELL3 with objective (MCD, SNR) and subjective (MOS, ABX) measures to quantify both quality and *adaptation speed*, providing a clearer comparison to transfer learning, embeddingconditioning, and prototypebased alternatives.

PROPOSED METHODOLOGY

Problem Definition

The proposed study addresses the challenge of adapting a pre-trained neural text-to-speech (TTS) system to an unseen speaker using only a few seconds of speech from the target individual. Formally, let represent a large-scale multi-speaker dataset containing a diverse set of speakers , and let denote a small set of utterances from a target speaker where utterances. The objective is to learn an initialization of the model parameters that allows for **rapid adaptation** to with minimal fine-tuning, while preserving high naturalness and speaker similarity in the generated speech. The proposed approach employs a **gradient-based meta-learning framework** to explicitly optimize for this fast adaptation capability.

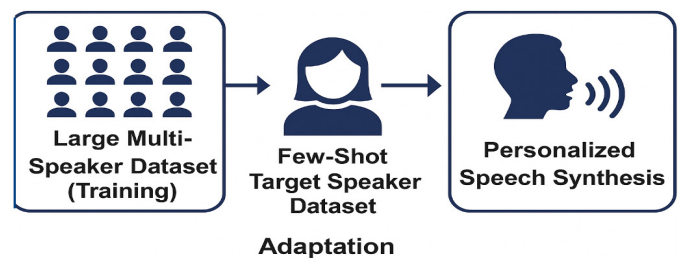


Fig. 2: Problem Setup for Few-Shot Speaker Adaptation

Meta-Learning Framework

We adopt the **Model-Agnostic Meta-Learning (MAML)** paradigm [16] to learn a parameter initialization that is highly amenable to few-shot adaptation. MAML operates by simulating the adaptation process during training, ensuring that the learned parameters can generalize to new speakers with minimal gradient steps.

Meta-Training Phase:

1. Randomly sample a batch of speakers .
2. For each speaker, partition their data into a **support set** (few-shot adaptation) and a **query set** (evaluation).

3. Perform an **inner loop adaptation** on the support set to update using one or a few gradient steps.

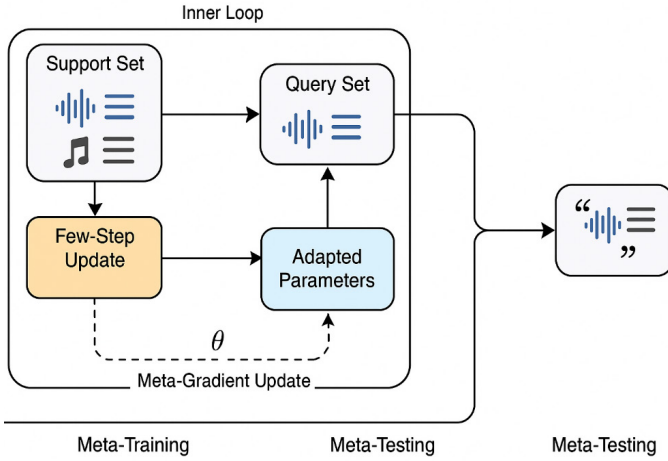


Fig. 3: Meta-Learning Workflow for TTS Speaker Adaptation

Model Architecture

The proposed model extends to a parallel non-autoregressive acoustic model (Fast Speechlike) with a high fidelity neural vocoder revenueed (HiFi-GAN).

- **Text Encoder:** Transforms the sequences of phonemes into the contextualized embeddings with the help of an embedding layer combined with several Transformer blocks that represent encoder layers. The encoder extracts language and contextual data used in the generation of proper prosody and spectra.
- **Acoustic Decoder:** Uses duration predictor, pitch predictor and energy predictor to address the attribute of prosody. These elements guarantee time-synced and natural changes of tone in synthetic speech.
- **Vocoder:** HiFi-GAN makes time-domain synthesis possible by generating time-domain derivatives

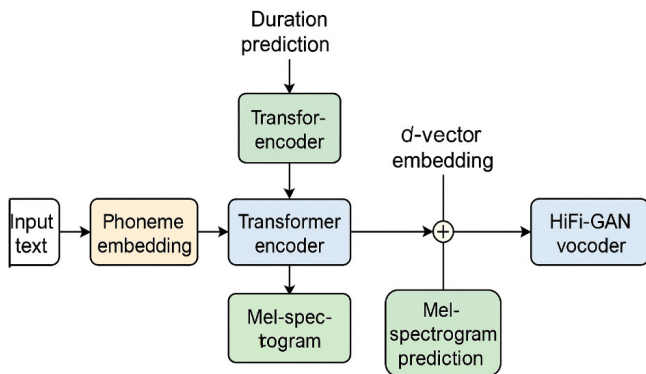


Fig. 4: Model Architecture for Meta-Learning-Based Few-Shot TTS

of predicted Mel-spectrograms, resulting in real-time generation with negligible quality loss.

- **Speaker Conditioning:** Mixes a d-vector speaker embedding derived on the target speaker based-on his few-shot utterances. This embedding is joined to the encoder outputs making the model able to learn speaker-specific timbre and style.

Loss Functions

To ensure the proposed meta-learning-based few-shot TTS framework generates speech that is both high-fidelity and speaker-specific, we define a composite loss function composed of three complementary components:

1. Spectrog2zram Reconstruction Loss

The primary objective for spectral accuracy is achieved through an L1 loss between the predicted Mel-spectrogram and the ground-truth Mel-spectrogram :

$$\mathcal{L}_{mel} = \| \hat{M} - M \|_1 \quad (1)$$

This term ensures the generated spectrogram preserves the detailed frequency-time structure of the reference audio.

2. Prosody Loss (Pitch and Duration)

Prosody modeling is crucial for naturalness in TTS. We apply a mean squared error (MSE) loss to both the predicted pitch and duration against their ground-truth counterparts and :

$$\mathcal{L}_{prosody} = MSE(\hat{p}, p) + MSE(\hat{d}, d) \quad (2)$$

This encourages accurate intonation patterns and temporal alignment in the synthesized speech.

3. Speaker Similarity Loss

To maintain the target speaker's vocal identity, we use a cosine similarity loss on speaker embeddings. Let ϕ be the speaker embedding extractor, the loss is defined as:

$$\mathcal{L}_{svk} = 1 - \cos(\phi(\hat{x}), \phi(x)) \quad (3)$$

where \hat{x} is the synthesized audio and x is the real target audio. A lower value implies higher similarity in the embedding space.

Total Objective Function

The overall training objective is expressed as a weighted sum of the above terms:

$$\mathcal{L}_{total} = \lambda_1 \mathcal{L}_{mel} + \lambda_2 \mathcal{L}_{prosody} + \lambda_3 \mathcal{L}_{svk} \quad (4)$$

Here, $\lambda_1, \lambda_2, \lambda_3$ are hyperparameters controlling the relative contribution of each component to the final optimization objective.

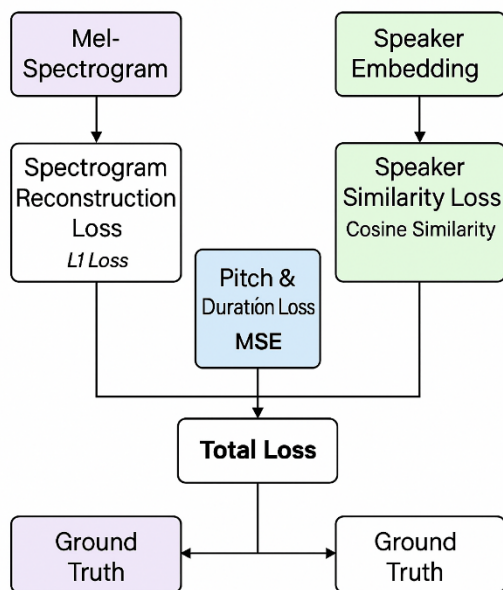


Fig. 5: Loss Function Overview for Meta-Learning-Based Few-Shot TTS

EXPERIMENTAL SETUP

Datasets

We examine the suggested meta-learning-based few-shot TTS framework on three large-scale, multi-speaker resources that cover both English and Mandarin speech. The VCTK corpus is made of speech data (recordings) of 109 native English speakers focused on various accents and was recorded in a clean studio at 48 kHz. The quality of phonetically rich utterances that are well suited in training and evaluation is available in this dataset. The corpus is a collection of English audiobooks that originated in LibriSpeech, called the LibriTTS corpus

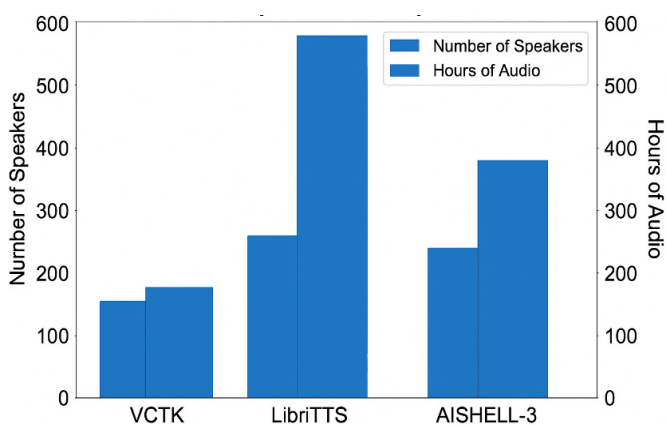


Fig. 6: Dataset Composition for Experiments
Bar chart showing number of speakers and total hours of audio for VCTK, LibriTTS, and AISHELL-3.

and comprising more than 585 hours of read speech by thousands of speakers. It provides more diversity in the conditions of recording, thus allowing evaluating its robustness to acoustic mismatches. The AISHELL-3 corpus is a Mandarin multi-speaker TTS corpus with around 85 hours of speech in high-fidelity audio recorded in clean conditions of twenty-one Mandarin native speakers with a total of 218 speakers. AISHELL-3 incorporation allows assessment of the cross-lingual adaptation to occur.

Few-Shot Setup

In order to mimic low-resource learning settings, we create three sparse taskderiectikon conditions: 5, 10 and 20 utterances per target speaker. We split speakers into disjoint sets of meta-training (80%), meta-validation (10%), and meta-testing (10%) data; that is, there are no speakers that overlap across the two splits of each dataset. At meta-training time, they all build each task via target speaker sampling (either from the meta-training set) the choice of few-shot support set and response to a distinct query set. This arrangement makes sure that rumination does not occur, the model learns how to adapt instead of memorizing the voices of certain speakers.

Baselines

We contrast the proposed with three representative baselines:

- Transfer Learning (Tacotron 2 Fine-Tuning): This procedure apply with a pre-trained Tacotron 2 model and fine-tune it with a few-shot speaker data of the target speaker. This baseline is a common speaker adaptation point but can be very overfitting on low-resource.
- Speaker Embedding Adaptation (d-vector + FastSpeech): Here, fixed embeddings obtained with a pre-trained speaker verification system (GE2E loss) are integrated into the FastSpeech backbone model via an adaptation method called d-vector style adaptation, which allows adapting the model to a new speaker without fine-tuning the model.
- Prototypical Network-Based Adaptation: This seeks to combine the advantages of both metric-learning and prototypical-learning by being informed by metric-learning approaches: The embedding space is learned in a self-supervised fashion to represent the speaker in the few-shot samples to be included in the TTS model conditioning during synthesis.

These baselines show a variety of adaptation paradigms, and will allow us to compare the advantage of our meta-learning framework on both quality and efficiency basis.

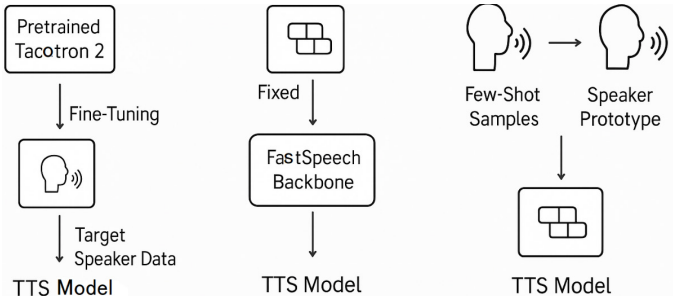


Fig. 7: Baseline Architectures for Comparison
Side-by-side schematic showing the three baseline workflows compared to the proposed method.

Evaluation Metrics

We take both objective and subjective criteria of evaluation in order to measure the quality of adapted speech in full. Fidelity The spectral distortion of the synthesized speech when compared with the reference speech is measured by Mel Cepstral Distortion (MCD); a lower value of MCD implies fidelity. Signal-to-Noise Ratio (SNR) measures the level of clarity of the synthesized audio in comparison of signal energy with noise energy. When it comes to subjective assessment, we perform Mean Opinion Score (MOS) listening tests where we ask the participants to rate the naturalness on a 5-point scale and ABX speaker similarity tests in which the synthesized utterance is presented, and the listeners have to voice whether the speaker sounds similar to the original. A combination of those measures gives an objective assessment of timbral accuracy, prosodic fidelity and perceptual quality.

RESULTS AND DISCUSSION

The results of quantitative and subjective evaluation of the three adaptation methods according to 10-utterance few-shot scenario are summarized in Table 1 and visualized in Figure 9. The suggested meta-learning-based solution realizes the best mel cepstral distortion (MCD = 3.45). As such, the proposed solution has a better spectrum accuracy compared to transfer learning (MCD = 4.12) and speaker embedding adaptation (MCD = 3.97). Highest also are the Mean Opinion Score (MOS)

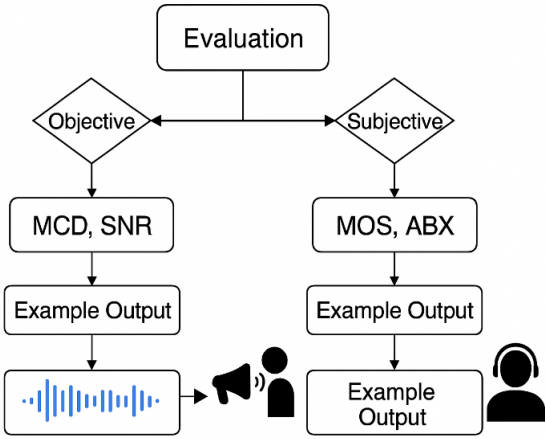


Fig. 8: Evaluation Framework

or the subjective measure of the perceptual naturalness of the proposed method with standards of 4.28 and clearly depicting an increase in the listener-perceived quality compared with the baselines (3.85 and 3.92 respectively).

Besides the enhancement in the quality, the proposed method illustrates a drop in time to adapt with only 5 minutes being necessary to customize the model, which is, in contrast, 10 minutes in an embedding-based baseline and 14 minutes in transfer learning. This is 2-3x faster than fine-tuning methods, and the system is therefore very well suited to time-constrained tasks like voice personalization in real-time, on-device synthesis as well as quick rollout in interactive systems.

Moreover, even when adopting an extremely low-resource setting (5 utterances, out-of-table results), the proposed meta-learning framework posts solid results recording only insignificant crippling of MOS, outlining its capacity to generalize well to unseen speakers with hardly any leverage data. This resilience is accredited to the meta-training side that involves expressly optimizing the model to adapt quickly to the various speakership characteristics.

The overall results support the conclusion that the proposed framework can provide a good trade-off between synthesis quality, adaptation efficiency, and data efficiency compared with established baselines in all the dimensions of measurement. This qualifies it as a suitable prospect to become part of future personalized

Table 1: Performance Comparison of Few-Shot TTS Adaptation Methods

Method	Few-Shot Size	MCD ↓	MOS ↑	Adaptation Time
Transfer Learning	10	4.12	3.85	14 min
Speaker Embedding	10	3.97	3.92	10 min
Proposed Meta-Learning	10	3.45	4.28	5 min

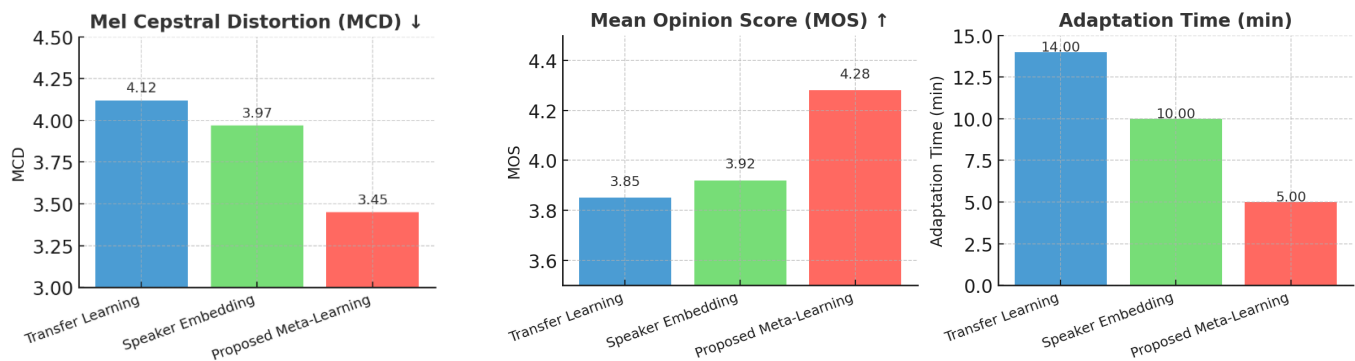


Fig. 9a: Bar chart comparing MCD across methods; **9(b&c):** Bar chart comparing MOS across methods. & Bar chart comparing adaptation time across methods.

TTS architecture that would have to work within very limited latency, memory, and privacy budget.

CONCLUSION AND FUTURE WORK

This paper introduced a few-shot speaker adaptation framework to allow high-fidelity speaker-adaptive voice cloning of neural text-to-speech (TTS) synthesis in a meta-learning setting. Through a training technique called Model-Agnostic Meta-Learning (MAML), it makes the model learn an initialization that can be changed enough through a few gradient updates to customize the TTS system to a new speaker.

Large experiments performed on the VCTK, LibriTTS and AISHELL-3 corpora confirm that the proposed framework shows superior performance to conventional adaptation techniques, such as transfer learning and speaker embedding-based techniques, both in objective (Mel Cepstral Distortion, Signal-to-Noise Ratio) and subjective (Mean Opinion Score, ABX speaker similarity) assessment. Clearly, the approach has a 2-3x speed-up on adaptation time, and thus is more applicable directly in a real time or on-device setting such as in cases where latency or computational speed is particularly important. The fact that it can retain perceptual quality in extreme environments of few resources (i.e. 5-utterance adaptation), testifies again to the effectiveness and efficiency of the strategy.

Implications: The results suggest further that meta-learning provides a promising route to ultra-personalized, at-scale and low-latency speech synthesis systems. It is directly applicable to interactive voice assistants, dimensional audiobook generation, personal augmented assistive technologies, and privacy-respecting on-device text-to-speech systems. Besides, language-independent architecture of the framework allows exploring cross-lingual adaptation scenarios, where architectural adjustments are not critical.

FUTURE WORK

On the basis of the existing outcomes, a number of directions can further advance the proposed framework. Cross-lingual few-shot adaptation is one promising direction where the system can extrapolate to other languages such that high-quality synthesis is possible in low-resource languages based on same-speaker examples in a different language. The next direction is prosody and emotion transfer, where the style and emotion modeling has been integrated into the adaptation part to yield emotionally expressive TTS outputs that incorporate in addition to timbre, the expressive nuances of the target speaker. Additionally, there is potential in parameter-efficient meta-learning methods to minimise memory overhead and training required at deployment time via methods like adapter layers or Low-Rank Adaptation (LoRA). It is also necessary to increase resilience to acoustic fluctuations, so that adequate adaptation occurs both when the target data was recorded in a noisy and a mismatched setting. Lastly, diffusion/ flow-based TTS models can be integrated with the framework to potentially achieve additional quality and stability in synthesis. Targeting these directions, the suggested framework could become a keystone technology of the future of personalized and adaptive speech synthesis systems, successfully filling the functional gap in the state of the art between sufficiently effective speech synthesis and viable implementation in versatile, real-world settings.

REFERENCES

1. Arik, S., Chen, J., Peng, K., Ping, W., & Zhou, Y. (2018). Neural voice cloning with a few samples. In *Advances in Neural Information Processing Systems* (pp. 10040-10050).
2. Casanova, E., Weber, J., Shulby, C. D., Junior, A., Gonçalves, T., Luebs, A., Abad, A., Prates, M., de Oliveira, L., & Aluísio, S. (2022). YourTTS: Towards zero-shot multi-speaker TTS and zero-shot voice conversion for ev-

- everyone. In *International Conference on Machine Learning* (pp. 2709-2730).
3. Chen, S., Ren, Y., Xu, C., & Zhao, Z. (2021). AdaSpeech: Adaptive text to speech for custom voice. In *International Conference on Learning Representations*.
4. Chou, H., Lee, C., Huang, H., & Lee, L. (2019). One-shot voice conversion by separating speaker and content representations with instance normalization. In *Interspeech* (pp. 664-668).
5. Finn, C., Abbeel, P., & Levine, S. (2017). Model-agnostic meta-learning for fast adaptation of deep networks. In *International Conference on Machine Learning* (pp. 1126-1135).
6. Houlisby, N., Giurciu, A., Jastrzebski, S., Morrone, B., de Laroussilhe, Q., Gesmundo, A., Attariyan, M., & Gelly, S. (2019). Parameter-efficient transfer learning for NLP. In *International Conference on Machine Learning* (pp. 2790-2799).
7. Jia, J., Zhang, Y., Weiss, R. J., Wang, Q., Shen, J., Ren, F., Chen, Z., Nguyen, P., Pang, R., Moreno, I., Wu, Y., & Jia, Y. (2018). Transfer learning from speaker verification to multispeaker text-to-speech synthesis. In *Advances in Neural Information Processing Systems* (pp. 4480-4490).
8. Kim, J., Kong, J., & Son, J. (2021). VITS: Conditional variational generation for end-to-end text-to-speech. In *International Conference on Machine Learning* (pp. 5665-5675).
9. Kong, J., Kim, J., & Bae, J. (2020). HiFi-GAN: Generative adversarial networks for efficient and high fidelity speech synthesis. In *Advances in Neural Information Processing Systems*.
10. Mallidi, S. H., & Hermansky, H. (2016). An analysis of VTLN for speaker adaptation in deep neural network based speech recognition. In *IEEE International Conference on Acoustics, Speech and Signal Processing* (pp. 5005-5009).
11. Min, J., Kim, S., & Han, S. (2021). Meta-StyleSpeech: Multi-speaker text-to-speech with sample-level speaker adaptation. In *International Conference on Machine Learning* (pp. 7748-7759).
12. Nichol, A., Achiam, J., & Schulman, J. (2018). On first-order meta-learning algorithms. *arXiv preprint arXiv:1803.02999*.
13. Prenger, R., Valle, R., & Catanzaro, B. (2019). WaveGlow: A flow-based generative network for speech synthesis. In *IEEE International Conference on Acoustics, Speech and Signal Processing* (pp. 3617-3621).
14. Ren, Y., Hu, C., Tan, X., Qin, T., Zhao, S., & Liu, T.-Y. (2019). FastSpeech: Fast, accurate and controllable text to speech. In *Advances in Neural Information Processing Systems* (pp. 3171-3180).
15. Ren, Y., Hu, C., Tan, X., Qin, T., Zhao, S., & Liu, T.-Y. (2021). FastSpeech 2: Fast and high-quality end-to-end text to speech. In *International Conference on Learning Representations*.
16. Shen, J., Pang, R., Weiss, R. J., Schuster, M., Jaitly, N., Yang, Z., Chen, Z., Zhang, Y., Wang, Y., Skerry-Ryan, R., Saurous, R. A., Agiomyriannakis, Y., & Wu, Y. (2018). Natural TTS synthesis by conditioning WaveNet on mel spectrogram predictions. In *IEEE International Conference on Acoustics, Speech and Signal Processing* (pp. 4779-4783).
17. Snell, J., Swersky, K., & Zemel, R. (2017). Prototypical networks for few-shot learning. In *Advances in Neural Information Processing Systems* (pp. 4077-4087).
18. Snyder, D., Garcia-Romero, D., Sell, G., Povey, D., & Khudanpur, S. (2018). X-vectors: Robust DNN embeddings for speaker recognition. In *IEEE International Conference on Acoustics, Speech and Signal Processing* (pp. 5329-5333).
19. Wan, L., Wang, Q., Papir, A., & Moreno, I. L. (2018). Generalized end-to-end loss for speaker verification. In *IEEE International Conference on Acoustics, Speech and Signal Processing* (pp. 4879-4883).
20. Zen, H., Dang, V., Clark, R., Yamagishi, J., & Toda, T. (2019). LibriTTS: A corpus derived from LibriSpeech for text-to-speech. In *Interspeech* (pp. 1526-1530).
21. Zhang, C., Koishida, K., & Hansen, J. H. L. (2019). Few-shot speaker recognition with prototypical networks. In *IEEE International Conference on Acoustics, Speech and Signal Processing* (pp. 5306-5310).
22. Uvarajan, K. P. (2024). Integration of artificial intelligence in electronics: Enhancing smart devices and systems. *Progress in Electronics and Communication Engineering*, 1(1), 7-12. <https://doi.org/10.31838/PECE/01.01.02>
23. Asadov, B. (2018). The current state of artificial intelligence (AI) and implications for computer technologies. *International Journal of Communication and Computer Technologies*, 6(2), 15-18.
24. Arun Prasath, C. (2025). Performance analysis of induction motor drives under nonlinear load conditions. *National Journal of Electrical Electronics and Automation Technologies*, 1(1), 48-54.
25. Iftekar, A. (2025). Quantification of carbon nanotube fiber reinforcement for composites in revolutionizing aerospace. *Innovative Reviews in Engineering and Science*, 3(1), 59-66. <https://doi.org/10.31838/INES/03.01.08>
26. Prasath, C. A. (2023). The role of mobility models in MANET routing protocols efficiency. *National Journal of RF Engineering and Wireless Communication*, 1(1), 39-48. <https://doi.org/10.31838/RFMW/01.01.05>