



# RASMALAI : Resources for Adaptive Speech Modeling in Indian Languages with Accents and Intonations

Ashwin Sankar<sup>1</sup>, Yoach Lacombe<sup>2</sup>, Sherry Thomas<sup>1</sup>, Praveen Srinivasa Varadhan<sup>1</sup>, Sanchit Gandhi<sup>3</sup>, Mitesh M Khapra<sup>1</sup>

<sup>1</sup>AI4Bharat, WSAI, Indian Institute of Technology Madras, India

<sup>2</sup>Cantina Labs, USA

<sup>3</sup>Mistral AI, France

sankar.ashwin01@gmail.com, yoach@cantina.ai, 21f3001449@ds.study.iitm.ac.in, cs21d201@cse.iitm.ac.in, sanchit.gandhi@mistral.ai, miteshk@dsai.iitm.ac.in

## Abstract

We introduce RASMALAI, a large-scale speech dataset with rich text descriptions, designed to advance controllable and expressive text-to-speech (TTS) synthesis for 23 Indian languages and English. It comprises 13,000 hours of speech and 24 million text-description annotations with fine-grained attributes like speaker identity, accent, emotion, style, and background conditions. Using RASMALAI, we develop INDICPARLERTTS, the first open-source, text-description-guided TTS for Indian languages. Systematic evaluation demonstrates its ability to generate high-quality speech for named speakers, reliably follow text descriptions and accurately synthesize specified attributes. Additionally, it effectively transfers expressive characteristics both within and across languages. INDICPARLERTTS consistently achieves strong performance across these evaluations, setting a new standard for controllable multilingual expressive speech synthesis in Indian languages.

**Index Terms:** TTS, Text-Prompted TTS, Expressive TTS

## 1. Introduction

Text-to-Speech (TTS) synthesis has evolved beyond generating speech from a fixed set of voices using high-quality recordings, with recent advancements exploring speech- and text-prompted TTS systems. Speech-prompted TTS [1, 2, 3, 4] synthesizes speech by replicating the style and speaker characteristics from a provided audio prompt. In contrast, text-prompted TTS [5, 6, 7, 8, 9] generates speech based on textual descriptions that can specify various attributes such as pitch, speed, noise levels, gender, and emotional expression. It is more user-friendly, as it removes the dependency on the availability of speech samples, allowing users to specify their desired speech characteristics via text. Training such models requires datasets containing triplets of {audio, transcript, description}, where the description guides the synthesis model to produce controlled and expressive speech with the specified attributes.

Much of the progress in text-prompted TTS has been limited to English and a few high-resource languages, while Indian languages have lagged behind due to the lack of suitable datasets. While several efforts [10, 11, 12] have introduced studio-recorded TTS datasets, their scale remains limited. Some TTS datasets [11] provide detailed annotations for emotions and speaking styles, while others [13], though not recorded in controlled environments, offer greater scale, speaker diversity, and broader linguistic and acoustic coverage. These datasets vary in size, quality, and style diversity but none provide structured text descriptions that explicitly characterize speaker voice, style, environment and expressive attributes. This gap highlights the need for a dataset that integrates linguistic diversity,

Table 1: Example showcasing (a) Descriptive, (b) Concise, and (c) Attribute Robust captions, for a set of speech attributes.

Attributes	Captions
Jaya, female speaker, Slightly close sounding, High pitch, Expressive tone, Slightly fast pace, Great speech quality, Speaking style is Anger	(a) Jaya, a female speaker, delivers speech in a slightly roomy environment with a high-pitched, expressive tone. She speaks at a slightly fast pace, with excellent overall speech quality. The intended style is anger.
	(b) Jaya, a female speaker, delivers high-pitched, expressive speech in a slightly enclosed environment at a fast pace, with excellent quality and an angry tone.
	(c) Jaya's angry tone, with a sharp voice, echoes with exceptional quality in a moderately reverberant environment.

large-scale speech data, and rich textual descriptions to enable text-prompted TTS for Indian languages.

To address this gap, we propose generating text-based speech descriptions using Large Language Models (LLMs) by leveraging attribute information in the form of structured tags. Specifically, we aggregate multiple existing TTS datasets containing {audio, transcript} pairs and extract three categories of attributes: (i) metadata attributes such as age, gender, and speaker identity, where available; (ii) stylistic attributes, including Ekman emotions when annotated; and (iii) acoustic features such as pitch [14], speaking rate, and clarity index [15, 16], derived using speech processing libraries. These attributes are then provided as input to an LLM to generate detailed text descriptions, creating multiple text prompts per utterance. This methodology is applied to a diverse set of {audio, transcript} datasets spanning multiple Indic languages as well as English with varied accents. The resulting dataset, RASMALAI, comprises over 13,000 hours of transcribed speech, with each utterance annotated with three distinct text-based descriptions of varying specificity and robustness (see Table 1). Additionally, these descriptions are translated into native languages using automated translation techniques [17], further enhancing the dataset's inclusivity for multilingual text-prompted TTS.

Using RASMALAI, we train INDICPARLERTTS, the first multilingual text-prompted TTS system covering 23 Indian languages. We thoroughly evaluate its performance through multiple systematic assessments. First, we assess its ability to function as a general high-quality TTS system by synthesizing natural speech in the voice of named high-quality speakers, as specified in the text prompts. Using MUSHRA tests,

we find that the synthesized speech exhibits high naturalness and strong fidelity to the intended speaker. Second, we evaluate the model’s instruction-following capability by determining whether the generated speech accurately reflects the descriptions provided in the text prompts. Our results show that INDICPARLERTTS effectively adheres to the specified characteristics. Third, we examine the model’s ability to generate expressive speech for speakers whose training data includes expressive variations. Our analysis shows that the model successfully synthesizes the intended emotions with very high accuracy. Fourth, we analyze the model’s capacity to synthesize expressive speech for speakers within a language even when no expressive data for that speaker was seen during training. We find that the model generalizes well in such cases, successfully performing style transfer to match the styles specified in the descriptions. Finally, we test its cross-lingual generalization capability by evaluating if the model can synthesize expressive speech in language  $L_2$  for a speaker in language  $L_1$ , even in cases where no expressive speech data was available for that speaker in either language. INDICPARLERTTS consistently achieves strong performance across all these evaluations, demonstrating its robustness and adaptability. To support further research in multilingual text-prompted TTS, we release all models and code to the community at <https://github.com/AI4Bharat/RASMALAI>.

## 2. RASMALAI : An Annotated Corpus for Controllable Multilingual TTS

Below we describe (i) existing TTS datasets from which {audio, text} pairs were collated (ii) our approach for generating text descriptions for these audios and (iii) statistics of our dataset.

### 2.1. Collating existing TTS datasets

We first collate {audio, text} pairs from multiple existing Indian and English language datasets, as mentioned below:

**RASA [11].** A diverse studio-quality speech dataset that includes neutral readings, expressive speech with six basic Ekman emotions, and command-based interactions from platforms such as Alexa, UMANG, and DigiPay. It also includes spontaneous conversations, news readings, and audiobook narrations. Each utterance is annotated with style, emotion labels, and speaker identity. It comprises 20 speakers across 13 languages, totaling around 400 hours of speech.

**IndicVoices [18].** A large-scale Indian speech dataset with over 7,200 hours of read, extempore, and conversational speech from 16,237 speakers across 22 languages. Originally designed for training ASR systems, it captures diverse acoustic, linguistic, and stylistic variations across various domains with recordings done in natural environments.

**IndicVoices-R [13].** A 4,000-hour subset of IndicVoices, restored using speech enhancement, with high-quality speech from 10,000+ speakers across 22 languages.

**IndicTTS [10] & LIMMITS [12].** High-quality, studio-recorded neutral read speech datasets, with two speakers per language. LIMMITS covers 9 languages with 40 hours per speaker, while IndicTTS covers 13 languages, with 20 hours of native speech and 20 hours of English recordings per language.

**GLOBE. [19]** A high-quality English speech corpus with 535 hours of speech from 23,519 speakers across 164 global accents. It refines Common Voice data through rigorous filtering, to enhance speech quality. The corpus is sampled at 24 kHz and includes detailed metadata with accent labels.

Table 2: Attributes coverage of different prompt-TTS Datasets

Attributes	AUDIOBOX	PROMPTTTS2	LIBRITTS-P	RASMALAI
SNR	✓	✗	✓	✓
c50	✗	✗	✓	✓
Speaking Rate	✓	✓	✓	✓
PESQ	✓	✗	✓	✓
FO <sub>mean</sub>	✓	✓	✓	✓
FO <sub>std</sub>	✗	✗	✓	✓
Age	✓	✗	✓	✓
Speaker ID	✗	✗	✓	✓
Gender	✓	✓	✓	✓
Style	✓	✓	✓	✓
Accent	✓	✗	✗	✓
Env-tags	✓	✗	✗	✓
Multi-lingual	✗	✗	✗	✓

### 2.2. Generating textual descriptions

For all the {audio, text} pairs in the above datasets, we generate descriptions using a two-stage approach as described below:

**Extracting attributes from the audio.** Building on the pipeline from [8], we extract key acoustic features, such as pitch, C50, SNR, and speaking rate, and incorporate PESQ to assess overall speech quality. When available, we utilize the general metadata attributes like age and gender, as well as stylistic attributes such as expression and speech style. To ensure structured and consistent text descriptions, the continuous features, like SNR and C50, are discretized into bins, and all attributes are categorized.

**Generating description from prompts using an LLM.** For each utterance, we provide LLAMA-3.1-8B-INSTRUCT with the categorized attributes, and prompt it to generate three types of textual descriptions (see Table 1) - (i) *Descriptive Prompt*: a detailed summary covering all labeled attributes, (ii) *Concise Prompt*: a brief description of the sample, and (iii) *Attribute-Robust Prompt*: which omits selected attributes to improve model robustness. We also translate the Descriptive Prompts into the respective language using IndicTrans2 [17], generating a *Native Prompt* to improve accessibility for multilingual users.

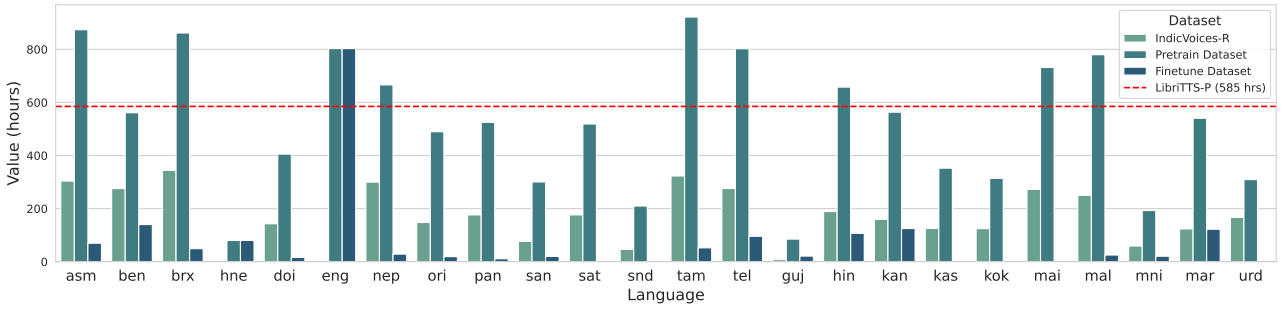
### 2.3. Dataset Statistics

We provide two versions of the dataset: RASMALAI-PRETRAIN, which includes the entire dataset, and RASMALAI-FINETUNE, a studio-quality subset with named speakers. The finetuning set is designed for downstream tasks, focusing on improving speaker consistency and expression rendering. Figure 1 presents a breakdown of dataset duration per language, comparing RASMALAI against IndicVoices-R. Notably, while prior open-source text-prompted TTS datasets [7] for English contain only 585 hours of speech, our dataset substantially expands coverage to 23 Indian languages. Moreover, in the pre-training set, it surpasses the scale of English data for 8 languages, further enhancing multilingual speech synthesis resources. Additionally, as shown in Table 2, our dataset offers richer attribute coverage in comparison to existing text-prompted datasets, providing fine-grained control for accents, environments, and emotions. Overall, RASMALAI contains **13,000 hours** of speech data and **24 million text-description annotations**.

## 3. Experimental Setup

**Model:** Building on Parler-TTS mini v1, we develop INDICPARLERTTS to support Indian languages by replacing the default tokenizer in Parler-TTS with an expanded Llama2 tok-

Figure 1: Comparison of durations for RASMALAI-PRETRAIN, RASMALAI-FINETUNE and INDICVOICES-R across 24 languages.



enizer, following the approach in [20], enabling better subword segmentation in Indian languages.

**Training Setup:** We train INDICPARLERTTS using 32 H100 GPUs for 1.5 days, leveraging the open-sourced configuration from the original Parler-TTS implementation. We followed the hyperparameters outlined in the configuration, and replace the original scheduler with CosineLR with 3,000 warm-up steps, with an effective batch size of 256.

**Data:** We pretrain the model on RASMALAI’s pretraining set, which spans 24 languages and 13,000 hours of data, and further finetune it on the finetuning set, comprising high-fidelity data from 18 languages with 1,804 hours of speech. This enhances the speech quality while improving speaker consistency.

**Evaluation:** We evaluate our model against the previous state-of-the-art zero-shot TTS system, INDICVOICECRAFT [13] (INDICVC), across 13 Indian languages in the Rasa dataset, using MUSHRA and present the results in Tables 3 and 4. Additionally, we report the consolidated CER, WER [18, 21], MOS [22], speaker similarity<sup>1</sup>, and *instruction adherence* in Table 5. We evaluate instruction adherence by re-annotating acoustic features of the synthesised samples using our pipeline up to the binning stage. We then compute IF-BLEU [7] by constructing prompts from binned values as comma-separated sequences and compare them to the original instruction from which this sample was synthesised. We also report attribute-level accuracy in Table 6, which is calculated directly by comparing the binned values against the original attributes used to create the description. To evaluate emotion rendering, we conduct a perceptual emotion classification test with human raters and plot a confusion matrix. Finally, we extend our evaluation to a cross-lingual and cross-speaker setup to assess INDICPARLERTTS for naturalness and expression generalization in a *MUSHRA-like* evaluation. For all our evaluations, we employ 307 listeners who rated 5078 utterances, across multiple tests.

## 4. Results

We evaluate INDICPARLERTTS on naturalness, style adherence, expressivity, and multilinguality.

### 4.1. Approaching Human-level Synthesis on Seen Speakers

Table 3 highlights languages where INDICPARLERTTS approaches human-level naturalness, while Table 4 reports cases where our model significantly outperforms the previous best-performing system. As shown in Table 3, INDICPAR-

<sup>1</sup><https://huggingface.co/microsoft/wavlm-base-plus-sv>

Table 3: Subjective evaluation of INDICPARLERTTS on seen speakers in the Rasa-13 test set, with MUSHRA scores highlighting four languages reaching near-human naturalness.

lang	brx	mai	mar	tel	Rasa-13
<b>Human</b>	93.4 ± 1.8	87.7 ± 1.8	91.8 ± 1.1	92.4 ± 1.0	89.7 ± 1.8
<b>Ours</b>	85.4 ± 2.6	84.8 ± 1.9	88.0 ± 1.8	85.9 ± 2.3	81.7 ± 2.7

Table 4: MUSHRA scores highlighting languages where our model significantly surpasses previous best-performing system.

lang	asm	ben	brx	nep	Rasa-13
<b>IndicVC</b>	61.6 ± 4.9	73.2 ± 2.9	71.0 ± 3.2	55.0 ± 4.8	73.0 ± 3.4
<b>Ours</b>	82.1 ± 2.4	83.4 ± 2.5	85.4 ± 2.6	75.4 ± 3.8	81.7 ± 2.7

LERTTS demonstrates high naturalness for seen expressive speakers in the Rasa test set covering 13 languages, achieving an average MUSHRA score of 81.7 compared to 89.7 for human speech. Furthermore, it approaches human-level speech synthesis in four languages: *Bodo (brx)*, *Maithili (mai)*, *Marathi (mar)*, and *Telugu (tel)*, demonstrating its ability to generalize across language families while maintaining naturalness and high fidelity. Table 4 further illustrates our model’s consistent improvement over the previous baseline on the Rasa test set, INDICVC with an average MUSHRA score of 81.7 versus 73.0. Notably, our model demonstrates significant improvements in *Assamese (asm)*, *Bengali (ben)*, and *Bodo (brx)*, elevating their MUSHRA scores from the Good (60–80) to the Excellent (80–100) range. Similarly, for *Nepali (nep)*, the score improves from Fair (55.0) to Good (75.4), marking a substantial leap in perceived quality compared to INDICVC (Table 4).

### 4.2. Faithful Execution of Instructions

In Table 5, we evaluate our model across key speech synthesis metrics, including intelligibility, naturalness, speaker consistency, and instruction adherence. The low CER (12%) and WER (24%) confirm that the model generates highly intelligible speech with minimal errors. A Noresqa-MOS score of 4.14, further indicates that the synthesized speech is natural and human-like. The speaker similarity score (0.95) demonstrates that our model is exceptionally consistent in preserving speaker identity across generated utterances. To assess instruction adherence, we report IF-BLEU, which reaches a high score of 93.18,

showing the model’s strong ability to follow given instructions accurately. Additionally, the attribute accuracy scores presented in Table 6 show strong performance in adhering to the specified attributes across most categories. Together, these metrics highlight our model as a highly faithful executor of instructions. However, we note that PESQ accuracy is relatively lower, likely due to the finetuning dataset’s focus on studio-quality recordings, limiting its ability to synthesize lower-quality speech.

Table 5: Automatic evaluation results for INDICPARLERTTS.

CER (%)	WER (%)	MOS	S-SIM	IF-BLEU
12	24	4.14	0.95	93.18

Table 6: Table shows the per attribute accuracy (%) of synthesized samples from the TTS subset of the test set.

c50	F0 <sub>std</sub>	F0 <sub>mean</sub>	SNR	Speaking Rate	PESQ
99.26	96.61	86.2	96.83	97.12	80.44

### 4.3. Evaluation of Expressive TTS

Our model exhibits strong emotion synthesis capabilities, as evaluated through a perceptual classification test. In this test, human listeners hear synthesized speech samples and classify them into one of seven target emotions. We report the accuracy of these human classifications for each emotion, providing insight into how well the synthesized speech conveys distinct emotional cues ( see Figure 2). The confusion matrix shows that key emotions such as *anger* (72.54%), *fear* (85.09%), *happy* (84.74%), and *surprise* (78.83%) are recognized with high accuracy, indicating that the model effectively conveys distinct emotional cues in synthesized speech. However, some emotions exhibit higher confusion rates. Notably, *sadness* (65.35%) is misclassified as *fear* 15.18% of the time. This confusion likely arises from their shared acoustic features, such as lower pitch, slower speech rate, and reduced energy. Similarly, neutral speech is often confused with other emotions, as seen in Figure 2. This warrants further investigation to better understand the underlying factors contributing to this ambiguity.

### 4.4. Zero-Shot Expressive Synthesis

INDICPARLERTTS demonstrates strong zero-shot expressive transfer, synthesizing expressive speech even for speakers with no expressive data in the training set. This is evident from the high MUSHRA scores (86.73) for Native speakers (expressing in their own language) in Table 7. Beyond native performance, our model also shows effective cross-lingual and cross-speaker expression transfer. Proximal speakers (from related languages) achieve a MUSHRA score of 80.86, while Distal speakers (from unrelated languages) score 76.01. Notably, both groups generate expressive speech in a new language, despite no training data for these speakers in the target language or emotion. The slight degradation in expressiveness between Proximal and Distal speakers suggests that the model may struggle to generalize effectively to speakers from more distant language families, where differences in accent, dialect, and phonetic structure could contribute to lower expressivity scores.

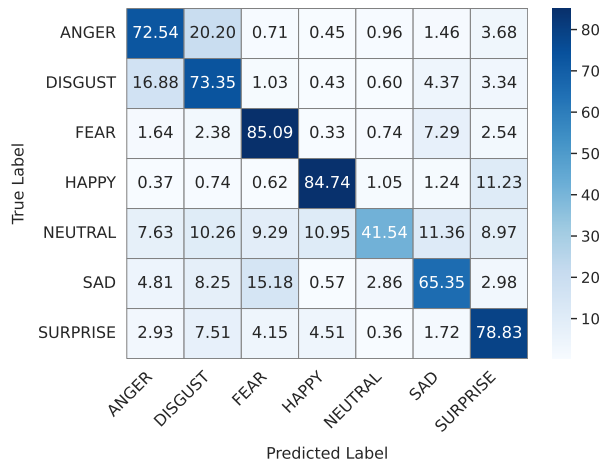


Figure 2: Confusion plot for Perceptual Emotion Classification

Table 7: MUSHRA scores for zero-shot expressive synthesis across three different speaker groups.

Native	Proximal	Distal
86.73 ± 2.10	80.86 ± 3.28	76.01 ± 5.25

## 5. Related Work

**Resources and Models for Text-Prompted TTS:** Recent works [5, 6, 8, 7] have explored large-scale dataset annotation to improve TTS control through diverse speech attributes. However, most efforts are limited to English [5, 7], and cover only a narrow set of attributes with restricted domain diversity. Among these, only [7] publicly releases its dataset. [8] introduces a language-agnostic approach for generating text annotations without relying on human annotation or speech captioning models. In this work, we extend that pipeline to incorporate a broader range of attributes, and support more languages.

**Resources for Indian TTS:** Previous efforts [10, 11, 13] have provided valuable datasets for Indian language TTS, supporting the development of single-speaker, expressive, and multilingual multi-speaker models. These resources have driven progress in zero-shot speaker and style-adaptive TTS. However, none of these datasets include the text annotations required for training text-prompted TTS models.

## 6. Conclusion

We present RASMALAI, a dataset with rich textual descriptions for 13,000 hours of speech across 24 languages, covering diverse speakers, emotions, and styles, built using structured attribute extraction and LLM-based text generation. Using this dataset, we train INDICPARLERTTS, the first multilingual text-prompted TTS system for Indic languages. Our evaluations show that INDICPARLERTTS produces highly natural, speaker-faithful speech, effectively follows text prompts, and synthesizes expressive speech with strong emotion accuracy. It also enables robust style transfer, generating expressive speech even for speakers without expressive training data, and demonstrates promising cross-lingual generalization. We release RASMALAI, along with our models and code, to drive progress in multilingual, expressive, and controllable TTS.

## 7. Acknowledgements

We also thank Yotta, EkStep Foundation, and Nilekani Philanthropies for their generous support, which made this work possible. Yoach Lacombe and Sanchit Gandhi contributed to this work during their time at Hugging Face.

## 8. References

- [1] P. Peng, P. Huang, S. Li, A. Mohamed, and D. Harwath, "Voicecraft: Zero-shot speech editing and text-to-speech in the wild," in *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, L. Ku, A. Martins, and V. Srikumar, Eds. Association for Computational Linguistics, 2024, pp. 12 442–12 462. [Online]. Available: <https://doi.org/10.18653/v1/2024.acl-long.673>
- [2] Y. A. Li, C. Han, V. S. Raghavan, G. Mischler, and N. Mesgarani, "StyleTTS 2: Towards human-level text-to-speech through style diffusion and adversarial training with large speech language models," in *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. [Online]. Available: <https://openreview.net/forum?id=m0RbqrUM26>
- [3] Z. Ju, Y. Wang, K. Shen, X. Tan, D. Xin, D. Yang, E. Liu, Y. Leng, K. Song, S. Tang, Z. Wu, T. Qin, X. Li, W. Ye, S. Zhang, J. Bian, L. He, J. Li, and sheng zhao, "Naturalspeech 3: Zero-shot speech synthesis with factorized codec and diffusion models," in *Forty-first International Conference on Machine Learning*, 2024. [Online]. Available: <https://openreview.net/forum?id=dVhrnjZJad>
- [4] S. E. Eskimez, X. Wang, M. Thakker, C. Li, C. Tsai, Z. Xiao, H. Yang, Z. Zhu, M. Tang, X. Tan, Y. Liu, S. Zhao, and N. Kanda, "E2 TTS: embarrassingly easy fully non-autoregressive zero-shot TTS," *CoRR*, vol. abs/2406.18009, 2024. [Online]. Available: <https://doi.org/10.48550/arXiv.2406.18009>
- [5] Y. Leng, Z. Guo, K. Shen, Z. Ju, X. Tan, E. Liu, Y. Liu, D. Yang, leying zhang, K. Song, L. He, X. Li, sheng zhao, T. Qin, and J. Bian, "PromptTTS 2: Describing and generating voices with text prompt," in *The Twelfth International Conference on Learning Representations*, 2024. [Online]. Available: <https://openreview.net/forum?id=NscXDyv2Bn>
- [6] A. Vyas, B. Shi, M. Le *et al.*, "Audiobox: Unified audio generation with natural language prompts," *CoRR*, vol. abs/2312.15821, 2023. [Online]. Available: <https://doi.org/10.48550/arXiv.2312.15821>
- [7] M. Kawamura, R. Yamamoto, Y. Shirahata, T. Hasumi, and K. Tachibana, "Libritts-p: A corpus with speaking style and speaker identity prompts for text-to-speech and style captioning," in *Proc. Interspeech 2024*, Sep. 2024.
- [8] D. Lyth and S. King, "Natural language guidance of high-fidelity text-to-speech with synthetic annotations," *CoRR*, vol. abs/2402.01912, 2024. [Online]. Available: <https://doi.org/10.48550/arXiv.2402.01912>
- [9] Z. Jiang, J. Liu, Y. Ren, J. He, Z. Ye, S. Ji, Q. Yang, C. Zhang, P. Wei, C. Wang, X. Yin, Z. MA, and Z. Zhao, "Mega-TTS 2: Boosting prompting mechanisms for zero-shot speech synthesis," in *The Twelfth International Conference on Learning Representations*, 2024. [Online]. Available: <https://openreview.net/forum?id=mvMI3N4AvD>
- [10] A. Baby, A. L. Thomas, N. Nishanthi, T. Consortium *et al.*, "Resources for indian languages," in *Proceedings of Text, Speech and Dialogue*, 2016.
- [11] P. Srinivasa Varadhan, A. Sankar, G. Raju, and M. M. Khapra, "Rasa: Building expressive speech synthesis systems for indian languages in low-resource settings," in *Interspeech 2024*, 2024, pp. 1830–1834.
- [12] A. Singh, A. Nagireddi, D. G. J. Bandekar, R. R. S. Badiger, S. Udupa, P. K. Ghosh, H. A. Murthy, P. Kumar, K. Tokuda, M. Hasegawa-Johnson, and P. Olbrich, "Limmits'24: Multi-speaker, multi-lingual indic tts with voice cloning," in *2024 IEEE International Conference on Acoustics, Speech, and Signal Processing Workshops (ICASSPW)*, 2024, pp. 61–62.
- [13] A. Sankar, S. Anand, P. S. Varadhan, S. Thomas, M. Singal, S. Kumar, D. Mehendale, A. Krishana, G. Raju, and M. M. Khapra, "Indicvoices-r: Unlocking a massive multilingual multi-speaker speech corpus for scaling indian TTS," in *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024. [Online]. Available: <https://openreview.net/forum?id=3qH8q02x0n>
- [14] M. Morrison, C. Hsieh, N. Pruyne, and B. Pardo, "Cross-domain neural pitch and periodicity estimation," in *arXiv preprint arXiv:2301.12258*, 2023.
- [15] M. Lavechin, M. Métais, H. Titeux, A. Boissonnet, J. Copet, M. Rivière, E. Bergelson, A. Cristia, E. Dupoux, and H. Bredin, "Brouhaha: multi-task training for voice activity detection, speech-to-noise ratio, and C50 room acoustics estimation," *ASRU*, 2023.
- [16] H. Bredin, R. Yin, J. M. Coria, G. Gelly, P. Korshunov, M. Lavechin, D. Fustes, H. Titeux, W. Bouaziz, and M.-P. Gill, "pyannote.audio: neural building blocks for speaker diarization," in *ICASSP 2020, IEEE International Conference on Acoustics, Speech, and Signal Processing*, Barcelona, Spain, May 2020.
- [17] J. P. Gala, P. A. Chitale, R. AK, V. Gumma, S. Doddapaneni, A. K. M., J. A. Nawale, A. Sujatha, R. Puduppully, V. Raghavan, P. Kumar, M. M. Khapra, R. Dabre, and A. Kunchukuttan, "Indictrans2: Towards high-quality and accessible machine translation models for all 22 scheduled indian languages," *Trans. Mach. Learn. Res.*, vol. 2023, 2023. [Online]. Available: <https://openreview.net/forum?id=vfT4YuzAYA>
- [18] T. Javed, J. Nawale, E. George, S. Joshi, K. Bhogale, D. Mehendale *et al.*, "IndicVoices: Towards building an inclusive multilingual speech dataset for Indian languages," in *Findings of the Association for Computational Linguistics: ACL 2024*, L.-W. Ku, A. Martins, and V. Srikumar, Eds. Bangkok, Thailand: Association for Computational Linguistics, Aug. 2024, pp. 10 740–10 782. [Online]. Available: <https://aclanthology.org/2024.findings-acl.639>
- [19] W. Wang, Y. Song, and S. Jha, "Globe: A high-quality english corpus with global accents for zero-shot speaker adaptive text-to-speech," in *Interspeech 2024*, 2024, pp. 1365–1369.
- [20] N. Mundra, A. N. K. Khandavally, R. Dabre, R. Puduppully, A. Kunchukuttan, and M. M. Khapra, "An empirical comparison of vocabulary expansion and initialization approaches for language models," in *Proceedings of the 28th Conference on Computational Natural Language Learning*, L. Barak and M. Alikhani, Eds. Miami, FL, USA: Association for Computational Linguistics, Nov. 2024, pp. 84–104. [Online]. Available: <https://aclanthology.org/2024.conll-1.8>
- [21] V. Noroozi, S. Majumdar, A. Kumar, J. Balam, and B. Ginsburg, "Stateful conformer with cache-based inference for streaming automatic speech recognition," in *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2024, Seoul, Republic of Korea, April 14-19, 2024*. IEEE, 2024, pp. 12 041–12 045. [Online]. Available: <https://doi.org/10.1109/ICASSP48485.2024.10446861>
- [22] P. Manocha and A. Kumar, "Speech quality assessment through mos using non-matching references," in *Interspeech 2022*, 2022, pp. 654–658.