

PRESENT: Zero-Shot Text-to-Prosody Control

Perry Lam, Huayun Zhang, Nancy F. Chen, Berrak Sisman, Dorien Herremans

Abstract—Current strategies for achieving fine-grained prosody control in speech synthesis entail extracting additional style embeddings or adopting more complex architectures. To enable zero-shot application of pretrained text-to-speech (TTS) models, we present PRESENT (PROsody Editing without Style Embeddings or New Training), which exploits explicit prosody prediction in FastSpeech2-based models by modifying the inference process directly. We apply our text-to-prosody framework to zero-shot language transfer using a JETS model exclusively trained on English LJSpeech data. We obtain character error rates (CER) of 12.8%, 18.7% and 5.9% for German, Hungarian and Spanish respectively, beating the previous state-of-the-art CER by over 2× for all three languages. Furthermore, we allow subphoneme-level control, a first in this field. To evaluate its effectiveness, we show that PRESENT can improve the prosody of questions, and use it to generate Mandarin, a tonal language where vowel pitch varies at subphoneme level. We attain 25.3% hanzi CER and 13.0% pinyin CER with the JETS model. All our code and audio samples¹ are available online.

Index Terms—speech synthesis, prosody, computational paralinguistics, zero-shot, language transfer

I. INTRODUCTION

RECENT neural text-to-speech (TTS) models have approached human-like naturalness in read speech. However, attaining similar expressiveness levels remains a challenge. A growing body of research aims to add and control speech prosody variations, progressing from digital signal processing (DSP) methods to style and emotion embeddings built into TTS architectures or even entire models to extract and transfer prosody.

On the waveform level, prosody control can be achieved through operations like time-stretching and pitch-shifting. DSP methods such as TD-PSOLA [1] and WORLD [2], despite their known artifacts, are still widely applied due to their speed and ease of use. Remarkably, they can perform as effectively as neural approaches like Controllable LPCNet [3].

In contrast, expressive TTS systems [4] allow the user to specify a style or emotion label during inference. Recent TTS models incorporate style or emotion information by extracting a reference embedding that represents the prosody or emotion from labelled audio, and adding it to the model encoder. This can be combined with a style bank for smooth style variation, such as in Global Style Tokens [5]. Further extensions include phoneme-level prosody control and hierarchical autoencoders to ensure coherence over the whole utterance [6].

Submitted on 16 Aug, 2024.

Perry Lam and Dorien Herremans are with Singapore University of Technology & Design, Singapore 487372 (e-mail: perry_lam@myemail.sutd.edu.sg and dorien_herremans@sutd.edu.sg).

Huayun Zhang and Nancy F. Chen are with the Institute of Infocomm Research, A*STAR, Singapore 138632 (e-mail: Zhang_Huayun@i2r.a-star.edu.sg and nfychen@i2r.a-star.edu.sg).

Berrak Sisman is with the University of Texas at Dallas, Richardson, TX 75080 USA (e-mail: berrak.sisman@utdallas.edu).

¹<https://github.com/iamanigeit/present> and <https://present2023.web.app/>

All of these approaches, however, require extra model components and/or further training. Therefore, to combine the simplicity of DSP methods with the naturalness of neural speech generation, we empower users to directly control prosody using the input text and inference parameters without the need for any fine-tuning or architectural modifications. We contribute significantly in the following three areas:

- **Extraction of prosodic effects from text**, such as extended duration in “A loooooong time” or the intonation variations in questions like “What was that?”. We take these prosodic parameters and modify the inference method of any TTS model with explicit duration, pitch, and energy (DPE) predictions to generate varying speech.
- **Zero-shot language transfer** with no target-language audio, relying solely on linguistic knowledge and modifying DPE to create new phonemes and speech patterns.
- **Subphoneme-level control**, achieved by subdividing phonemes and applying custom pitch and energy over the subdivisions, which helps us change long vowel intonation and model tonal languages like Mandarin.

Though our primary goal is to explore the limits of editing inference-time prosody predictions, in doing so, we achieve state-of-the-art results in zero-shot language transfer.

The rest of this paper is organized as follows: Section 2 summarizes relevant research, Section 3 describes our approach, Section 4 lists our experiment results and Section 5 concludes our paper.

II. RELATED WORK

Based on our main contributions, we divide the related work into the broad categories of (1) speech effect tagging, (2) zero-shot language transfer, and (3) fine-grained prosody control.

A. Speech Effects Tagging

Text-based methods for manipulating speech can be categorized into explicit and implicit forms. Explicit speech descriptors such as gender and emphasis have been integrated into the industry standard Speech Synthesis Markup Language (SSML) over the past two decades [7]. Yet, there has been relatively limited published research on SSML, even though there have been notable introductions of TTS models with new style tags, as demonstrated in [8].

Implicit methods establish connections between prosodic features and text, such that a sentence like “this is interesting!” would sound excited. Typically, this means that the text embeddings from a language model are used as input either at the subword [9] [10] or phoneme level [11] [12]. However, due to their inherent limitations in customizing prosody changes, recent projects inspired by advancements in computer vision and language processing let user input a natural-language style prompt like “sighing with helpless feeling” to generate prosodic output [13].

B. Zero-Shot Language Transfer

While multilingual TTS models have existed for some time, they rely on large multilingual corpora, which disadvantages lower-resourced languages. Transfer learning [14] [15] and data balancing [16] techniques have been employed, but these still require at least some audio data. With only International Phonetic Alphabet (IPA) transcriptions in the target language, [17] proposed using IPA phonological features to extend existing models on unseen phonemes, whereas two very recent large models have proposed zero-shot TTS with only text data available in the target language.

The first model, VALL-E X [18], uses AudioLM [19] codec codes as acoustic tokens in place of mel spectrograms as intermediate features, and treats the cross-lingual TTS model as a massive language model (LLM) that can be trained with self-supervised masking. Given a speech sample in the source language, plus source and target language phoneme sequences, it extracts the source acoustic tokens from the speech sample and the LM predicts the target acoustic tokens. Since the acoustic tokens contain speaker, recording conditions, and subphoneme information, the decoder can reconstruct the waveform for the target language in the source speaker’s voice.

The second model, ZM-Text-TTS [20], also uses masked multilingual training, but on IPA / byte tokens and raw text. The pretraining results in a language-aware embedding layer that is fed to a conventional multilingual TTS system for training with seen languages, and the model can accept IPA / byte tokens for unseen languages during inference. Nevertheless, VALL-E X is not publicly available, and ZM-Text-TTS does not account for prosody in language transfer.

C. Fine-grained Prosody Control

As utterance-level styles are now commonplace, research has shifted to controlling prosody at the phoneme level. Since acceptable prosodies are obtained by learning and sampling from a variational latent space, hierarchical variational auto-encoders (VAEs) [21] can achieve fine prosodic gradations,

down to the syllable, phone or even frame level [6].

Alternatively, others use phone-level DPE for interpretable prosody control. This was the approach of earlier research [22], but to improve output naturalness, [23] and [24] used k-means clustering on duration and pitch for each speaker, and kept the resulting centroids as discrete prosody tokens. This allows the tokens to be substituted at inference time to customize prosody, while decoding with a prosody attention module ensures information flows to the output. Meanwhile, since the advent of explicit DPE models like FastSpeech2 [25], models like [26] and [27] have extra modules attached that accept emotional dimensions (valence, arousal, dominance) that feed into phone-level DPE predictors, allowing for continuous emotion control.

III. PROPOSED METHOD

PRESENT offers a versatile approach to (1) extract inference parameters and (2) integrate them with explicit DPE predictions to generate variations in pronunciation and prosody, all without requiring additional modules or fine-tuning. The specific method of parameter extraction and integration is adaptable to the task at hand and can be customized by the user. Fig 1 illustrates this process with examples for English text-to-prosody and English-to-Mandarin language transfer.

A. Obtaining Prosodic Effects from Text

We preprocess the input text to capture common dialogue features such as CAPS or *asterisks* for emphasis, repeaaaaated letters or ti~~ldes for long phonemes, and special characters like underscores and carets or questions for tone modification. As TTS systems usually rely on phoneme input, we align the text to phonemes so that DPE changes can be applied at the right positions. While grapheme-to-phoneme (G2P) systems are widely available, G2P alignment systems are outdated and lack Python implementations. Thus, we develop our own aligner combining the ‘‘Phonetic Alignment’’ and ‘‘IP Alignment’’ approaches in [28]. We begin

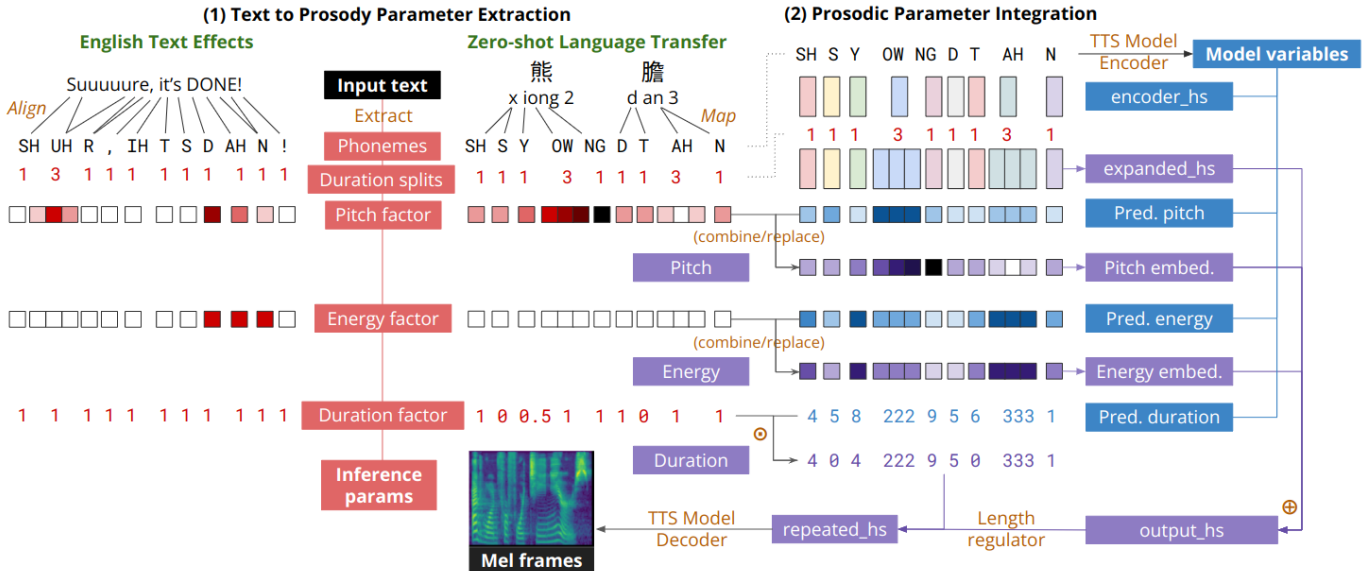


Fig. 1. Overview of inference modification in PRESENT. Green items are tasks, black items are input/outputs, Red items represent PRESENT outputs, blue items are the outputs from the FS2-based model, purple are the combined outputs, and gold are operations.

with a list of allowed grapheme-phoneme mappings (e.g. $ch \rightarrow \{CH, K, SH\}$) and search for a constraint-satisfying alignment, or return the least-cost path by dynamic programming if an alignment cannot be found. For example, given the word *where* and pronunciation $W EH R$, the aligner returns $wh \rightarrow W$, $e \rightarrow EH$, $r \rightarrow R$, $e \rightarrow \emptyset$. With an invalid pair like *whence* and $W Z EH T$, the aligner chooses the minimal cost of disallowed mappings: $\emptyset \rightarrow Z$, $n \rightarrow T$, $c \rightarrow \emptyset$. The aligner also allows us to detect possibly wrong dictionary entries, such as the **G** of $EEG \rightarrow IY IY G IY$ in the CMU Pronouncing Dictionary, and correct them ($IY IY JH IY$).

B. Generating Parameters for Zero-Shot Language Transfer

Manipulating DPE directly allows us to create phonemes and intonation patterns not found in our model’s language. For instance, we can approximate the German and Spanish /x/ with ARPabet [HH K HH] and durations [1, 0, 1] to velarize [HH] without producing a distinct [K]. Since ZM-Text-TTS tested on German (de), Hungarian (hu), and Spanish (es), we describe the main ideas for adapting American English models to these languages in Table I, and also include Mandarin (cmn) for the next section. The full conversion tables are in our code.

TABLE I

MAPPING RULES TO ARPABET. IPA SYMBOLS IN /SLASHES/, PINYIN (ANGLE BRACKETS), ARPABET [SQUARE BRACKETS]. D = DURATION FACTOR, P = PITCH CHANGE, E = ENERGY CHANGE, || = WORD SEP.

All	Use combinations, possibly with zero duration, for phonemes that don’t exist in English (for /β λ ɲ ç x y o/ etc). Examples: German /œ/ → [W EH] D=[0,1] and /ç/ → [H SH S] D=[0,1,0] Hungarian /y/ → [UH Y] D=[0,1] and /j/ → [G Y] D=[0.7,0] Spanish /t/ → [R HH R] D=[1,0,1] and /β/ → [B V] D=[0,1] Mandarin ⟨zh⟩ → [T SH] D=[1,0] and ⟨x⟩ → [SH S] D=[1,0]
de	Shorten long vowels corresponding to German short vowels: /a/ → [AA] D=0.5 Make schwa clearer and prevent merging into next phoneme: /ə/ → [AX] D=1.5 E=+1 and [AH] → [AH ,] D=[1,0]
hu	Shorten long vowels corresponding to Hungarian short vowels: /u/ → [UW] D=0.5 Reduce phoneme lengths as Hungarian has faster speaking speed: /b/ → [B] D=0.7 Double consonants for long consonants: /k:/ → [K K] D=[0.7,0.7]
es	Reduce phoneme lengths as Spanish has faster speaking speed: /t/ → [T] D=0.7 Shorten long vowels corresponding to Spanish short vowels: /o/ → [OW] D=0.4 Insert semivowels or very short pauses between consecutive vowels to avoid diphthongization, and raise the stressed vowel: /o, i/ → [OW W IY] D=[0.4,0.4,0.7] P=[0,0,+1], E=[0,0,+0.5] Double plosives between vowels to prevent devoicing: /apa/ → [AA P P AA] D=[0.7, 0, 0.7, 0.7]
cmn	Define conversion for initial and rimes instead of individual phonemes, as they don’t always combine sequentially: ⟨i⟩ → [IY ,] D=[1,0] but ⟨in⟩ → [IH IY N] D=[1,0,1] Use 0-duration voiceless phone to stop voicing vs [HH] for aspiration: ⟨g⟩ → [G K] D=[1,0] and ⟨k⟩ → [K HH] D=[1,0,5] Syllables that are pronounced differently from initial + rime mapping have their own rules: ⟨i⟩ before ⟨z(h),c(h),s(h)⟩ → [Z UH] D=[0.5,0.7] but ⟨chi⟩ → [CH HH R R] D=[1,0.5,1,1] Add pause before characters that start with a vowel: ⟨ai⟩ → [ʔai/ → [, AY] D=[0.2,1] Set glides to half duration: ⟨iu⟩ → [Y OW] D=[0.5,1]

C. Subphoneme-level Control and Tonal Languages

To achieve tone contour effects like the rising-falling pitch in “Suuuuure!”, the phoneme must be split and each subphoneme assigned separate pitches. Thus, we repeat the encoder output h for the divided phoneme (boxes for OW and AH of `encoder_hs` in Fig 1). This differs from simply repeating the phonemes for inference, as that would generate varying h and possibly make the phoneme pronounced multiple times.

One evident use case for pitch effects pertains to questions. While humans can clearly perceive a question via prosody, TTS systems still lack proper intonation. For English question prosody, [29] enabled users to choose from a range of pretrained prosody templates. However, to maintain simplicity and avoid adding the complexity of language models, we follow the prosodic analysis of [30], applying a low-to-high accent on the locus of interrogation and the final word of the question to convey question intonation.

Another critical test of our subphoneme tone contour approach is its ability to model tonal languages. After applying the phoneme changes in Table I, we split each vowel nucleus into subphonemes and assign their pitch following Mandarin tone contours in Table II. Contours on the five-point tonal scale are then normalized to pitch values between $[-2.0, +2.0]$.

TABLE II
TONE-PITCH MAPPING.

Tone	1	2	3	4	5
Contour	55	24	212	52	-
Pitch	+2, +2	-1, +1	-1, -2, -1	+2, -1	0

The initial and coda (if any) take the start and end pitch of the tone contour. Table III demonstrates one example of how pinyin ⟨tian2⟩ maps to [T HH Y EH N]. Pitch transitions are smoothed across syllables to avoid abrupt pitch changes.

TABLE III
EXAMPLE OF ARPA-PINYIN MAPPING.

ARPA	T	HH	Y		EH		N
Duration	×1	×0.5	×0.5		×1		×1
Subphonemes	T	HH	Y	EH	EH	EH	N
Pitch	-1	-1	-1	-0.33	0.33	+1	+1

As Mandarin is a syllable-timed language, we keep the ×1 duration constant, with the neutral tone at half duration. Finally, we leverage `pywordseg` to segment Mandarin text and introduce brief pauses between words for better enunciation.

IV. EXPERIMENTS

We conducted our experiments using the ESPnet [31] toolkit for reproducibility. Our source model was the publicly released English-only single-speaker JETS [32] model pretrained on LJSpeech, known for achieving state-of-the-art naturalness.

A. Zero-Shot Language Transfer

We first evaluate the ability of the PRESENT to produce intelligible speech in other languages. To generate audio samples, we extract raw text from the CSS10 dataset and phonemize them with `espeak-ng`, then perform phoneme conversion and DPE editing. As a baseline, we employ ZM-Text-TTS [20], the only open-source zero-shot language transfer system. As ZM-Text-TTS was trained and evaluated on the CSS10 [33] datasets’ European languages subset, we compare

PRESENT on the 3 languages on which ZM-Text-TTS has done zero-shot TTS: German, Hungarian, and Spanish.

For fair comparison, we follow them in evaluating character error rate (CER) by running generated audio through Whisper’s [34] multilingual speech recognition. We use the large-v2 model from SYSTRAN’s faster-whisper on default settings for its speed and robustness. CER is computed from the length-normalized Levenshtein distance between Whisper transcripts and ground truth, ignoring punctuation and whitespace.

Since ZM-Text-TTS pretrains on multilingual text before training on text-audio pairs, there are two settings in their evaluation: text-unseen (where the target language text is not available for pretraining) and text-seen (where the text is available, but there is no paired audio). Naturally, the text-unseen case leads to higher CER.

TABLE IV
CER COMPARISON. EURO. LANG. = EUROPEAN LANGUAGES. GROUND TRUTH = RAW CSS10 AUDIO. RANGE GIVEN FOR SPANISH ZM-TEXT-TTS IS FROM USING EITHER PHONEMES OR BYTES AS INPUT.

Target	(Source Languages) Model	Text	CER
German	(6 Euro. langs.) ZM-Text-TTS	Unseen	38.75
		Seen	28.01
	(English only) PRESENT	Unseen	12.82
		Ground Truth	–
Hungarian	(6 Euro. langs.) ZM-Text-TTS	Unseen	52.62
		Seen	50.11
	(English only) PRESENT	Unseen	18.73
		Ground Truth	–
Spanish	(7 Euro. langs.) ZM-Text-TTS	Unseen	44.75 – 64.07
		Seen	11.69 – 18.27
	(English only) PRESENT	Unseen	5.92
		Ground Truth	–

PRESENT reduces the previous state-of-the-art CER by over $2\times$ for each language, even with a single off-the-shelf English-only model to generate them with no further training. In fact, our CER is close to the ZM-Text-TTS multilingual model trained *with* target audio (German 9.76, Hungarian 9.11 and Spanish 5.32). This shows that phoneme conversion followed by prosody manipulation is critical for zero-shot language transfer.

B. Subphoneme-Level Control

For question prosody, we took the first 10 dialogues from the DailyTalk dataset [35] and extracted the first single-sentence question from each of them, making 10 questions in total. We report the Mean Opinion Score (MOS) for ground truth audio from DailyTalk, unaccented JETS-generated audio, and PRESENT-accented audio. Experiments were conducted with PsyToolkit [36] [37] and 15 responses were received.

TABLE V
MOS FOR QUESTION PROSODY.

Ground Truth	JETS	PRESENT
4.46	3.73	3.92

We then evaluate the ability of the JETS model to produce intelligible Mandarin speech by synthesizing speech based on the AISHELL-3 test set transcripts. As a baseline, we take the IPA multilingual model (pretrained on 7 European languages)

model from ZM-Text-TTS, convert pinyin transcripts to IPA, and use the best-approximation phoneme when a Mandarin phoneme does not exist in the pretrained IPA symbol set. We then input the the ground truth, PRESENT, and ZM-Text-TTS audio into the state-of-the-art Paraformer automatic speech recognition (ASR) framework [38], and ensure transcriptions only contain Chinese by masking decoder outputs to $-\infty$ for alphanumeric tokens. Due to hallucination issues in individual models, we use a mixture-of-experts consisting of the aishell2-vocab5212 and paraformer-large-vocab8404 models.

The CER for transcriptions are computed at both Hanzi level and pinyin level in Table VI. As Mandarin has many homophones, pinyin CER is a better measure of intelligibility; romanization also makes it comparable to European-language CER where a change like *la* to *le* is 50% CER, not 100%. Thus, we romanize Mandarin transcripts with `pypinyin` and include pinyin CER with tone counting as one character (i.e. `ping1` versus `ping2` would be 20% CER).

We also measured MOS (naturalness only) by selecting 15 utterances from the AISHELL-3 test set that made sense as full standalone sentences without jargon, rare words, or names. We skipped testing ZM-Text-TTS and PRESENT without tones or duration control, since they did not produce intelligible results. As before, the survey was created with PsyToolkit and 10 responses were received from Mandarin speakers.

TABLE VI
ENGLISH-TO-MANDARIN LANGUAGE TRANSFER RESULTS.

	% Hanzi CER	% Pinyin CER	MOS
Ground Truth	1.2	0.9	4.65
PRESENT	25.3	13.0	2.18
– w/o tones	59.5	33.8	1.92
– w/o tones/duration	105.4	63.9	–
ZM-Text-TTS	95.0	71.7	–

The dramatic CER reductions from ZM-Text-TTS to PRESENT with phoneme conversion, duration and tones applied in succession demonstrates the effectiveness of our subphoneme-level DPE control. Using Latin orthography, the English-to-Mandarin language transfer CER is equal to the average of English-to-{German, Hungarian, Spanish}, and even outperforms all previous baselines on those languages. Still, MOS testing reveals the naturalness limitation of PRESENT-generated audio for human listeners due to the strong American accent, despite some improvement via tone contouring.

V. CONCLUSIONS

We have introduced PRESENT, a novel approach that explores the limits of using only DPE predictions in a single-speaker English-only JETS model, without any additional embeddings or training. Our technique allows us to create prosodic effects from text and synthesize speech in unseen languages. Our zero-shot language transfer far outstrips previous state-of-the-art for European languages. Furthermore, the phoneme conversion and tone contour techniques we develop could enable direct accented speech generation (as the results are all American-accented), or TTS for hundreds of tonal minority languages within the Mainland Southeast Asian linguistic area that are only recorded in phonetic transcriptions.

REFERENCES

- [1] E. Moulines and F. Charpentier, "Pitch-synchronous waveform processing techniques for text-to-speech synthesis using diphones," *Speech communication*, vol. 9, no. 5-6, pp. 453-467, 1990.
- [2] M. Morise, F. Yokomori, and K. Ozawa, "World: a vocoder-based high-quality speech synthesis system for real-time applications," *IEICE TRANSACTIONS on Information and Systems*, vol. 99, no. 7, pp. 1877-1884, 2016.
- [3] M. Morrison, Z. Jin, N. J. Bryan, J.-P. Caceres, and B. Pardo, "Neural pitch-shifting and time-stretching with controllable lpcnet," *arXiv preprint arXiv:2110.02360*, 2021.
- [4] J. Pitrelli, R. Bakis, E. Eide, R. Fernandez, W. Hamza, and M. Picheny, "The ibm expressive text-to-speech synthesis system for american english," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 14, no. 4, pp. 1099-1108, 2006.
- [5] Y. Wang, D. Stanton, Y. Zhang, R.-S. Ryan, E. Battenberg, J. Shor, Y. Xiao, Y. Jia, F. Ren, and R. A. Saurous, "Style tokens: Unsupervised style modeling, control and transfer in end-to-end speech synthesis," in *International conference on machine learning*. PMLR, 2018, pp. 5180-5189.
- [6] T. Kenter, V. Wan, C.-A. Chan, R. Clark, and J. Vit, "Chive: Varying prosody in speech synthesis with a linguistically driven dynamic hierarchical conditional variational network," in *International Conference on Machine Learning*. PMLR, 2019, pp. 3331-3340.
- [7] Z. W. Shuang and D. Burnett, "Speech synthesis markup language (SSML) version 1.1," W3C, W3C Recommendation, Sep. 2010, <https://www.w3.org/TR/2010/REC-speech-synthesis11-20100907/>.
- [8] M. Kim, S. J. Cheon, B. J. Choi, J. J. Kim, and N. S. Kim, "Expressive Text-to-Speech Using Style Tag," in *Proc. Interspeech 2021*, 2021, pp. 4663-4667.
- [9] S. Ammar Abbas, T. Merritt, A. Moinet, S. Karlapati, E. Muszynska, S. Slangen, E. Gatti, and T. Drugman, "Expressive, Variable, and Controllable Duration Modelling in TTS," in *Proc. Interspeech 2022*, 2022, pp. 4546-4550.
- [10] L. Chen, Y. Deng, X. Wang, F. K. Soong, and L. He, "Speech bert embedding for improving prosody in neural tts," in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 6563-6567.
- [11] Y. A. Li, C. Han, X. Jiang, and N. Mesgarani, "Phoneme-level bert for enhanced prosody of text-to-speech with grapheme predictions," in *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2023, pp. 1-5.
- [12] Y. Jia, H. Zen, J. Shen, Y. Zhang, and Y. Wu, "PnG BERT: Augmented BERT on Phonemes and Graphemes for Neural TTS," in *Proc. Interspeech 2021*, 2021, pp. 151-155.
- [13] D. Yang, S. Liu, R. Huang, G. Lei, C. Weng, H. Meng, and D. Yu, "Instructtts: Modelling expressive tts in discrete latent space with natural language style prompt," *arXiv preprint arXiv:2301.13662*, 2023.
- [14] T. Nekvinda and O. Dušek, "One Model, Many Languages: Meta-Learning for Multilingual Text-to-Speech," in *Proc. Interspeech 2020*, 2020, pp. 2972-2976.
- [15] K. Azizah, M. Adriani, and W. Jatmiko, "Hierarchical transfer learning for multilingual, multi-speaker, and style transfer dnn-based tts on low-resource languages," *IEEE Access*, vol. 8, pp. 179 798-179 812, 2020.
- [16] J. Yang and L. He, "Towards Universal Text-to-Speech," in *Proc. Interspeech 2020*, 2020, pp. 3171-3175.
- [17] M. Staib, T. H. Teh, A. Torresquintero, D. S. R. Mohan, L. Foglianti, R. Lenain, and J. Gao, "Phonological Features for 0-Shot Multilingual Speech Synthesis," in *Proc. Interspeech 2020*, 2020, pp. 2942-2946.
- [18] Z. Zhang, L. Zhou, C. Wang, S. Chen, Y. Wu, S. Liu, Z. Chen, Y. Liu, H. Wang, J. Li *et al.*, "Speak foreign languages with your own voice: Cross-lingual neural codec language modeling," *arXiv preprint arXiv:2303.03926*, 2023.
- [19] Z. Borsos, R. Marinier, D. Vincent, E. Kharitonov, O. Pietquin, M. Sharifi, D. Roblek, O. Teboul, D. Grangier, M. Tagliasacchi *et al.*, "Audiolm: a language modeling approach to audio generation," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2023.
- [20] T. Saeki, S. Maiti, X. Li, S. Watanabe, S. Takamichi, and H. Saruwatari, "Learning to speak from text: Zero-shot multilingual text-to-speech with unsupervised text pretraining," in *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence, IJCAI 2023, 19th-25th August 2023, Macao, SAR, China*. ijcai.org, 2023, pp. 5179-5187. [Online]. Available: <https://doi.org/10.24963/ijcai.2023/575>
- [21] G. Sun, Y. Zhang, R. J. Weiss, Y. Cao, H. Zen, and Y. Wu, "Fully-hierarchical fine-grained prosody modeling for interpretable speech synthesis," in *ICASSP 2020-2020 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, 2020, pp. 6264-6268.
- [22] Y. Lee and T. Kim, "Robust and fine-grained prosody control of end-to-end speech synthesis," in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 5911-5915.
- [23] P. Tsiakoulis, "Improved prosodic clustering for multispeaker and speaker-independent phoneme-level prosody control," in *Speech and Computer: 23rd International Conference, SPECOM 2021, St. Petersburg, Russia, September 27-30, 2021, Proceedings*, vol. 12997. Springer Nature, 2021, p. 112.
- [24] N. Ellinas, M. Christidou, A. Vioni, J. S. Sung, A. Chalamandaris, P. Tsiakoulis, and P. Matorocostas, "Controllable speech synthesis by learning discrete phoneme-level prosodic representations," *Speech Communication*, vol. 146, pp. 22-31, 2023.
- [25] Y. Ren, C. Hu, X. Tan, T. Qin, S. Zhao, Z. Zhao, and T. Liu, "Fastspeech 2: Fast and high-quality end-to-end text to speech," in *9th International Conference on Learning Representations, ICLR, 2021*.
- [26] S. Sivaprasad, S. Kosgi, and V. Gandhi, "Emotional Prosody Control for Speech Generation," in *Proc. Interspeech 2021*, 2021, pp. 4653-4657.
- [27] S. Kosgi, S. Sivaprasad, N. Pedanekar, A. Nelakanti, and V. Gandhi, "Empathic machines: using intermediate features as levers to emulate emotions in text-to-speech systems," in *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, 2022*, pp. 336-347.
- [28] S. Jiampojamarn and G. Kondrak, "Phoneme alignment: An exploration," in *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, 2010, pp. 780-788.
- [29] J. Lee, J. Y. Lee, H. Choi, S. Mun, S. Park, J.-S. Bae, and C. Kim, "Intotts: Intonation template based prosody control system," *arXiv preprint arXiv:2204.01271*, 2022.
- [30] N. Hedberg and J. M. Sosa, "The prosody of questions in natural discourse," in *Speech Prosody 2002, International Conference, 2002*.
- [31] T. Hayashi, R. Yamamoto, K. Inoue, T. Yoshimura, S. Watanabe, T. Toda, K. Takeda, Y. Zhang, and X. Tan, "Espnet-TTS: Unified, reproducible, and integratable open source end-to-end text-to-speech toolkit," in *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 7654-7658.
- [32] D. Lim, S. Jung, and E. Kim, "JETS: Jointly Training FastSpeech2 and HiFi-GAN for End to End Text to Speech," in *Proc. Interspeech 2022*, 2022, pp. 21-25.
- [33] K. Park and T. Mulc, "CSS10: A Collection of Single Speaker Speech Datasets for 10 Languages," in *Proc. Interspeech 2019*, 2019, pp. 1566-1570.
- [34] A. Radford, J. W. Kim, T. Xu, G. Brockman, C. McLeavey, and I. Sutskever, "Robust speech recognition via large-scale weak supervision," in *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, ser. Proceedings of Machine Learning Research, A. Krause, E. Brunskill, K. Cho, B. Engelhardt, S. Sabato, and J. Scarlett, Eds., vol. 202. PMLR, 2023, pp. 28 492-28 518. [Online]. Available: <https://proceedings.mlr.press/v202/radford23a.html>
- [35] K. Lee, K. Park, and D. Kim, "Dailytalk: Spoken dialogue dataset for conversational text-to-speech," in *ICASSP 2023 - 2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2023, pp. 1-5.
- [36] G. Stoet, "Psytoolkit: A software package for programming psychological experiments using linux," *Behavior Research Methods*, vol. 42, pp. 1096-1104, 2010.
- [37] —, "Psytoolkit: A novel web-based method for running online questionnaires and reaction-time experiments," *Teaching of Psychology*, vol. 44, no. 1, pp. 24-31, 2017.
- [38] Z. Gao, S. Zhang, I. McLoughlin, and Z. Yan, "Paraformer: Fast and Accurate Parallel Transformer for Non-autoregressive End-to-End Speech Recognition," in *Proc. Interspeech 2022, 2022*, pp. 2063-2067.