

CONCSS: CONTRASTIVE-BASED CONTEXT COMPREHENSION FOR DIALOGUE-APPROPRIATE PROSODY IN CONVERSATIONAL SPEECH SYNTHESIS

Yayue Deng¹, Jinlong Xue¹, Yukang Jia³, Qifei Li¹, Yichen Han¹
Fengping Wang¹, Yingming Gao¹, Dengfeng Ke², Ya Li^{1,*}

¹Beijing University of Posts and Telecommunications, Beijing, China
²Beijing Language and Culture University, Beijing, China
³Perfect World Co., Ltd, Beijing, China

ABSTRACT

Conversational speech synthesis (CSS) incorporates historical dialogue as supplementary information with the aim of generating speech that has dialogue-appropriate prosody. While previous methods have already delved into enhancing context comprehension, context representation still lacks effective representation capabilities and context-sensitive discriminability. In this paper, we introduce a contrastive learning-based CSS framework, CONCSS. Within this framework, we define an innovative pretext task specific to CSS that enables the model to perform self-supervised learning on unlabeled conversational datasets to boost the model’s context understanding. Additionally, we introduce a sampling strategy for negative sample augmentation to enhance context vectors’ discriminability. This is the first attempt to integrate contrastive learning into CSS. We conduct ablation studies on different contrastive learning strategies and comprehensive experiments in comparison with prior CSS systems. Results demonstrate that the synthesized speech from our proposed method exhibits more contextually appropriate and sensitive prosody.

Index Terms— Conversational Speech Synthesis (CSS), Contrastive Learning, Context Understanding

1. INTRODUCTION

Recent advancements in speech synthesis systems [1, 2, 3, 4] have enabled the generation of high-quality speech. However, in some complex scenarios, such as human-computer interaction (HCI), these systems still fall short since they are unable to generate audio with natural and human-like prosody.

Previous psychological studies [5, 6] show that when processing conversation in real-time, our brain quickly uses various information, such as prior statements and the speaker’s identity, to help understand the current speech. Similar to human communication, several studies [7, 8, 9, 10, 11] demonstrate that incorporating historical dialogues as supplementary information into speech synthesis can improve the model’s understanding of prior statements, thereby helping to enhance the prosody of the synthesized audio. Hence, a greater interest has been shown in developing conversational speech synthesis (CSS) task which focuses on improving the model’s context-understanding capability and generating audio with context-appropriate prosody.

Guo et al. [12] first introduce a GRU-based context modeling method to extract semantic information from the history context at

the utterance level. Moreover, further studies [13, 11, 14] have verified that combining acoustic and textual context can greatly improve the naturalness of synthesized speech. Nishimura et al. [13] utilize cross-modal attention to capture both long-term linguistic and prosodic context information. Xue et al. [11] and Hu et al. [14] combine both fine-grained and coarse-grained context encoders to provide sufficient information for better context comprehension. Meanwhile, some graph-based methods [15, 16] are proposed to capture multi-scale context dependencies between different modalities.

While prior CSS frameworks have demonstrated the capability to enhance context comprehension, thereby yielding context-related prosody, the question remains: Is this output vector of the context encoder sufficiently indicative of the underlying context variations? This is a concern since these CSS approaches rely solely on jointly training the acoustic model and context encoder using the mel-reconstruction loss. Without explicit constraints, the latent vector derived from the context encoder may not possess desired characteristics like interpretability, strong representative capacity, and context-sensitive discriminability. Hence, it is imperative to include an additional constraint in order to obtain better results for context representation extraction from a jointly-trained model.

Motivated by this, we propose a novel conversational speech synthesis framework CONCSS which combines contrastive learning and CSS to learn effective and context-sensitive representation. To this end, we define a novel pretext task specific to CSS as an important strategy to learn context representations using pseudo labels. By doing so, the context encoder is then forced to learn what we care about, e.g., underlying semantics and discernible context variations. To validate the effectiveness of the method, we conduct comprehensive and uniquely designed experiments. Results demonstrate that our proposed method can enhance discriminability and context sensitivity of context vectors compared with previous methods. Furthermore, the context-sensitive vectors guide downstream acoustic model to synthesize audio with more context-appropriate prosody. Our work has three main contributions:

1. The primary contribution is the novel CSS framework CONCSS for enhancing prosody, driven by contrastive learning to address the context understanding issue from a self-supervised perspective. To the best of our knowledge, we are the first to adopt contrastive learning for the CSS task.
2. To address the issue of context representation in CSS, we design an innovative pretext task specific to CSS tasks, along with a sampling strategy. This approach compels the context encoder to generate distinct representations for diverse scenarios, thereby promoting both context sensitivity and distinctiveness in prosody.
3. We comprehensively evaluate models on their ability to pro-

* Ya Li is the corresponding author.

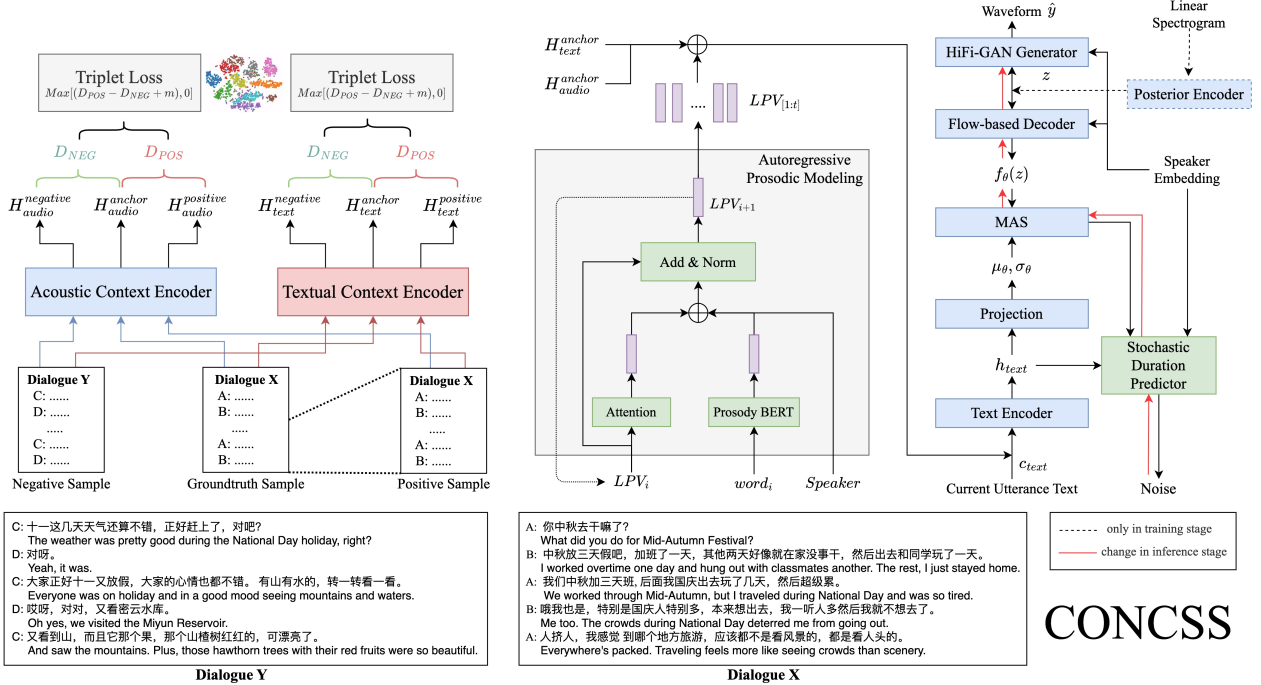


Fig. 1. Illustration of our proposed CONTRastive-based Conversational Speech Synthesis (CONCSS).

duce context-sensitive vectors and dialogue-appropriate prosody. Specifically, we compare prosodic performance between different CSS frameworks when they are exposed to diverse contexts. The rest of the paper is organized as follows. Section 2 introduces our proposed method. Experimental setup and analysis of results are shown in Section 3. Finally, Section 4 concludes the paper.

2. PROPOSED METHOD

To strengthen context comprehension capability of the context encoder and develop context-sensitive context representation without using labeled data, we propose a novel self-supervised contrastive-based framework where we enhance the prosody of speech from context-related to context-sensitive.

2.1. Problem Definition and Framework Overview

Assume an input sequence of N utterances $[(u_1, p_1), (u_2, p_2), \dots, (u_N, p_N)]$ in a conversation where each utterance u_i is spoken by speaker p_i . In the CSS task, considering a context length of i , the context encoder takes historical utterances $[(u_1, p_1), (u_2, p_2), \dots, (u_i, p_i)]$ as input and generates context vector h_i . Then, context vector h_i is integrated into the downstream speech synthesis framework, generating audio with dialogue-appropriate prosody.

As illustrated in Figure 1, the proposed framework comprises the following four enhancements to the latest speech synthesis model VITS [4]: 1) Leveraging an innovative pretext task, context-dependent pseudo-labels are created to direct the model toward obtaining desired attributes (e.g., context-sensitive discriminability and interpretability); 2) We incorporate an acoustic and textual context encoder, maintaining structure consistent with [11]; 3) We employ triplet loss [17] coupled with a hard negative sampling strategy to amplify the model’s context awareness and comprehension, thereby endowing the synthesized speech with context-sensitive distinctiveness and context-appropriate prosody; 4) We utilize an autoregressive prosodic modeling (APM) module with a pre-trained prosodic language model [18].

2.2. Context-aware Contrastive Learning

2.2.1. Pretask Definition

We hypothesize that a context encoder with adept context understanding capabilities can detect variations in prior statements and, furthermore, exhibit appropriate and distinguish representations to different contexts.

Hence, the context-based pretext task is defined as follows: Positive sample h_i^p is the output of context encoder given by the same historical dialogue with the groundtruth sample but has different context lengths, whereas negative sample h_i^n is given by different dialogues with non-overlapping context backgrounds.

Thus we want context vectors can satisfy the following criteria:

$$D(h_i, h_i^p) < D(h_i, h_i^n) \quad (1)$$

where similarity measurement $D(\cdot)$ is defined as squared Euclidean distance in the context representation space.

2.2.2. Sampling Strategy

Furthermore, previous work [17] has demonstrated that hard negative samples provide the most significant assistance to update the gradients in the optimization process. In order to select hard negative samples, we consider that negative samples should be derived from two sources: 1) from dialogues with the same speakers but in entirely disparate contexts (intra-speaker classes); 2) from dialogues involving different speakers with non-overlapping context backgrounds (inter-speaker classes).

In our case, intra-speaker classes are hard negative samples because of the potential similarities for context vectors of the same speaker. Although the context variations are totally different, the context spoken by the same speakers may lead to similar speaker-related prosody. Hence, we further employ intra-speaker classes as a hard negative sampling strategy to enhance contrastive learning.

CONCSS

Dialogue Y
 C: 十一这几天天气还算不错，正好赶上了，对吧？
 The weather was pretty good during the National Day holiday, right?
 D: 对呀。
 Yeah, it was.
 C: 大家正好十一又放假，大家的心情也都不错。有山有水的，转一转看一看。
 Everyone was on holiday and in a good mood seeing mountains and waters.
 D: 哎呀，对呀，又看密云水库。
 Oh yes, we visited the Miyun Reservoir.
 C: 又看到山，而且它那个景，那个山楂树红红的，可漂亮了。
 And saw the mountains. Plus, those hawthorn trees with their red fruits were so beautiful.

Dialogue X
 A: 你中秋去干嘛了？
 What did you do for Mid-Autumn Festival?
 B: 中秋放三天假吧，加班了一天，其他两天好像就在家没事干，然后出去和同学玩了一天。
 I worked overtime one day and hung out with classmates another. The rest, I just stayed home.
 A: 我们中秋加三天班，后面我国庆出去玩了几次，然后超级累。
 We worked through Mid-Autumn, but I traveled during National Day and was so tired.
 B: 哦我也是，特别是国庆人特别多，本来想出去，我一听人多然后我就不想去了。
 Me too. The crowds during National Day deterred me from going out.
 A: 人挤人，我感觉到哪个地方旅游，应该都不是看风景的，都是看人头的。
 Everywhere's packed. Traveling feels more like seeing crowds than scenery.

2.2.3. Multi-modal Context Comprehension with Triplet Loss

To achieve context-sensitive discriminability of context vector, we directly maximize the similarity between positive pairs and minimize the similarity of negative pairs via a triplet loss [17], which can be expressed as:

$$L(h_i^a, h_i^p, h_i^n) = \max\{D(h_i^a, h_i^p) - D(h_i^a, h_i^n) + m, 0\} \quad (2)$$

where the margin parameter m imposing the distance between negative samples and positive samples should be larger than m . The anchor h_i^a , positive h_i^p , and negative samples h_i^n have been generated as described above for textual and acoustic modalities.

For each modality, the losses are averaged over each loss element in the batch and used to update the context encoder.

$$\begin{cases} \mathcal{L}_{text}^k = L(H_{text}^a, H_{text}^p, H_{text}^n) \\ \mathcal{L}_{audio}^k = L(H_{audio}^a, H_{audio}^p, H_{audio}^n) \end{cases} \quad (3)$$

$$\mathcal{L}_{contra} = \frac{1}{N} \sum_{k=1}^N (\mathcal{L}_{text}^k + \mathcal{L}_{audio}^k) \quad (4)$$

where H_{audio} and H_{text} represent acoustic and textual context vectors, respectively, as shown in Figure 1. Minimization of this loss encourages context encoder following:

$$\|h_i^a - h_i^p\|_2^2 + m < \|h_i^a - h_i^n\|_2^2, \quad \forall (h_i^a, h_i^p, h_i^n) \in \mathcal{T} \quad (5)$$

where \mathcal{T} contains all possible triplets in the training set.

Consequently, the context encoder yields quality representations of underlying context variations which are used later for transferring knowledge to CSS tasks.

2.3. Autoregressive Prosodic Modeling

Inspired by [19], we adopt an autoregressive prosodic modeling (APM) module to promote fluent and natural prosody. This module ensures that, while generating the current latent prosody vector LPV_{i+1} , it takes into account both word-level prosody information and preceding latent prosody vectors (LPV_s). The APM is substituted by an attention mechanism [20] coupled with a pre-trained prosody BERT architecture.

3. EXPERIMENTS

3.1. Experimental Setting

We conduct experiments on the open-source Chinese conversational speech corpus¹ which consists of 10 hours of transcribed Mandarin conversational speech spoken by 30 speakers on certain topics. To better utilize and train the data, we use the ffmpeg toolkit to split continuous dialogues into distinct audio clips, removing non-lexical noises. These clips are then sequenced by utterance order, and their text is converted to phonemes using an open-source tool². The final processed data totals around 9.2 hours.

For the backbone of CSS, we adhere to the vanilla training setups and implementation of VITS. We first pre-train the backbone on the Biaobei Chinese TTS dataset³ for an initial 5k steps, then the whole CONCSS framework is trained on the Chinese conversation speech corpus for 20k steps with a batch size of 16 to achieve satisfactory results. The Prosody BERT [18] is finetuned on the prosodic annotation of Biaobei data to learn word-level prosody information.

¹<https://magichub.com/datasets/mandarin-chinese-conversational-speech-corpus-multiple-devices/>

²<https://github.com/PaddlePaddle/PaddleSpeech/>

³https://www.data-baker.com/open_source.html

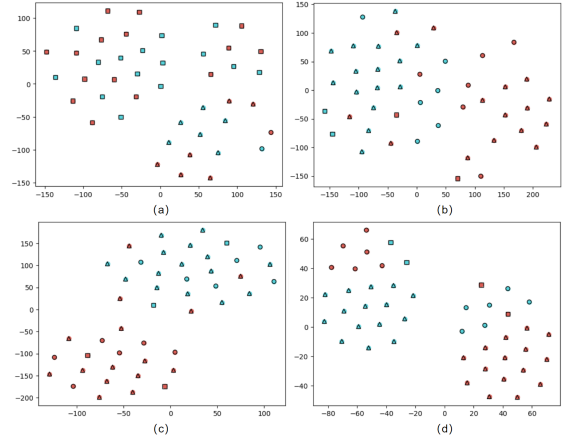


Fig. 2. T-SNE [21] visualization of context vector distribution from three separate dialogues. (a) M2CTTS; (b) S1; (c) S2; (d) S3. Identical shapes represent synthesized text from the same dialogue. Red indicates a Fake context, while blue signifies a Real context.

3.2. Compared Models and Evaluation Metrics

In our experimental evaluation, we conduct comparative and ablation studies on the following seven CSS methods:

- **GRU-based**: proposed by Guo et al. [12] to achieve GRU-based context understanding.
- **M2CTTS**: proposed by Xue et al [11] for multi-modal context understanding (M2CTTS-M6).
- **CONCSS-w/o-APM**: The proposed framework without APM module. With regard to various contrastive learning strategies, we ran ablation studies as follows:
 - **S1**: adopts basic contrastive loss proposed by Chopra et al. [22]. Besides, it does not employ the sampling strategy.
 - **S2**: utilizes triplet loss [17] without using the hard negative sampling strategy.
 - **S3**: employs triplet loss and hard negative sampling strategy, as illustrated in Section 2.2
- **CONCSS**: The proposed CSS framework, denoted as **S4**.

We assess the effectiveness of our proposed method using both subjective and objective measure metrics:

- **Naturalness MOS**: evaluates the overall naturalness of speech on a 1-5 scale, particularly evaluating whether its prosody fits historical dialogues. The greater the score, the better the naturalness.
- **CMOS**: assesses the context-aware prosodic expression. Raters score from -3 (A model is completely better) to 3 (B model is completely better) in 1-point increments. The preference between A and B is denoted as (A vs. B). A higher score indicates a stronger preference for B.
- **Mel Loss**: calculates the Mean Squared Error (MSE) between the predicted and ground-truth mel-spectrograms.
- **Log F0 RMSE and MCD**: are calculated after Dynamic Time Warping [23] to evaluate prosody-related performance of speech.

3.3. Comparison of Dialogue-Appropriate Prosody

We randomly select 20 synthesized speech for each comparison model. Subjective and objective results are compiled and presented in Table 1 and Table 2, and we have the following observations: 1) The MOS score in a multimodal setting surpasses that of a unimodal approach, which can also be supported by previous studies

Table 1. Subjective evaluation (context-appropriate prosody and naturalness) for different models.

Model	GRU-based	M2CTTS	S1	S2	S3	S4		
MOS (\uparrow)	3.396 \pm 0.107	3.438 \pm 0.104	3.528 \pm 0.097	3.708 \pm 0.108	3.838 \pm 0.110	3.967 \pm 0.120		
Models	GRU-based vs. M2CTTS	M2CTTS vs. S1	S1 vs. S2	S1 vs. S3	S2 vs. S3	S3 vs. S4	GRU-based vs. S4	M2CTTS vs. S4
CMOS (\uparrow)	0.200	0.388	0.796	0.983	0.492	0.325	1.846	1.788

Table 2. Objective evaluation metrics primarily focus on the context-sensitive prosody. The `Real` context type uses the correct context for the current synthesized sentence, whereas the `Fake` type randomly selects from unrelated dialogues.

Method	Set	Type	Mel Loss (\downarrow)	Log F0 RMSE (\downarrow)	MCD (\downarrow)
GRU-based		Real	3.599	0.2949 \pm 0.1192	5.3590
		Fake	3.683	0.3001 \pm 0.1164	5.3781
M2CTTS		Real	3.579	0.2936 \pm 0.1014	5.3236
		Fake	3.596	0.3036 \pm 0.1277	5.3882
CONCSSS	S1	Real	3.609	0.2911 \pm 0.1099	5.3382
		Fake	3.626	0.3203 \pm 0.1093	5.4923
	S2	Real	3.556	0.2906 \pm 0.1047	5.2883
		Fake	3.638	0.3311 \pm 0.1417	5.5157
	S3	Real	3.530	0.2821 \pm 0.0960	5.2748
		Fake	3.715	0.3272 \pm 0.1455	5.6923
	S4	Real	3.525	0.2803 \pm 0.0961	5.2634
		Fake	3.649	0.3252 \pm 0.1097	5.6041

Table 3. Subjective evaluation between different context types.

Model	MOS (\uparrow)		CMOS (\uparrow)
	Real	Fake	Real vs Fake
GRU-based	3.442 \pm 0.111	3.388 \pm 0.102	0.325
M2CTTS	3.504 \pm 0.100	3.312 \pm 0.112	0.445
S1	3.638 \pm 0.091	3.250 \pm 0.116	0.492
S2	3.796 \pm 0.076	3.229 \pm 0.101	0.529
S3	3.958 \pm 0.074	3.308 \pm 0.100	0.804

[11, 15]; 2) Within contrastive-based frameworks, S4 achieves optimal performance with a MOS score of 3.967. Objective metrics also reach the best scores, as highlighted in bold in Table 2. In comparison to baseline systems (GRU-based and M2CTTS), it displays a pronounced preference with CMOS scores of 1.846 and 1.788, respectively; 3) Employing triplet loss and hard negative sampling strategy can both enhance prosody performance. Specifically, after utilizing triplet loss, the CMOS score is 0.796 and the MCD score is 5.2883. Further incorporation of the negative sample strategy elevates the CMOS to 0.983 and results in an MCD score of 5.2748; 4) The APM module also displays an enhancement in prosody, with the MOS score increasing from 3.838 to 3.967.

To summarize, both subjective and objective evaluations confirm the efficacy of the contrastive learning approach in enhancing the dialogue-appropriate prosody and naturalness of synthesized speech. Ablation studies are conducted on the contrastive loss function, sampling strategy, and APM module, all of which positively impact the dialogue-appropriate prosody. Compared to basic contrastive loss, the utilization of triplet loss provides more efficient contrastive learning, deepening model’s context understanding and thereby facilitating the acquisition of effective context representation. Additionally, Employing the hard negative sampling strategy further boosts contrastive learning, resulting in more dialogue-appropriate prosody.

3.4. Comparison of Context-sensitive Distinctiveness

To compare discriminability of context vectors generated by different context modeling methods, we first visualize the distance of context vectors. As shown in Fig 2, we can observe that: 1) non-

contrastive framework is locally clustered based on the semantic space of current synthesized sentence (the same shape), which indicates previous frameworks primarily generate context vectors based on the current text, rather than history context; 2) both S2 and S3 methods demonstrate superior context-sensitive discriminative capabilities in comparison to the non-contrastive and S1 methods.

Furthermore, to assess the performance of context-sensitive prosody, we employ a CSS task-specific evaluation method where various CSS frameworks are exposed to diverse contexts. Specifically, given the same synthesized text, we feed compared models with distinct contexts either from the current dialogue or from other irrelevant dialogues, named `Real` and `Fake` respectively. Objective and subjective results are shown in Table 2 and Table 3, respectively. We observe that, in subjective evaluation, S3 has the top CMOS score of 0.804, indicating that S3 method exhibits the highest sensitivity to context. Besides, the MOS score in S3 shows a clearer distinction between the `Real` and `Fake` types by 0.65, compared to the 0.054 of the GRU-based method. In objective evaluation, the `Fake` and `Real` differences for Mel Loss in the GRU-based system and S3 are 0.084 and 0.185 respectively, while for MCD they are 0.0191 and 0.4175. We can infer that contrastive-based approaches exhibit greater sensitivity to different contexts than non-contrastive approaches. Both subjective and objective metrics demonstrate that the sampling strategy and triplet loss aid in enhancing the model’s context comprehension. This, in turn, results in more discriminative context vectors for diverse context inputs, subsequently influencing the prosodic expression in downstream speech generation. Audio samples are available on the project page ⁴.

4. CONCLUSION

In this paper, we introduce CONCSS, a contrastive-based CSS framework that leverages self-supervised training to enhance context understanding. Specifically, we define a pretext task to enable the model to utilize pseudo-labels, thereby increasing context sensitivity to various scenarios. Furthermore, we propose a hard negative sampling strategy to boost context comprehension and the generation of effective context representation. Comprehensive subjective and objective evaluations demonstrate that the proposed method can enhance context comprehension and yield well-representative context vectors, enabling the generation of speech with more appropriate and context-sensitive prosody. Furthermore, we propose a CSS task-specialized subjective evaluation method to assess models’ performance in context understanding and context-sensitive distinctiveness. The effectiveness of the proposed framework is verified in comprehensive experiments.

5. ACKNOWLEDGEMENTS

The work was supported by the National Natural Science Foundation of China (NSFC) (No.62271083), the Fundamental Research Funds for the Central Universities (No.2023RC13), the open research fund of The State Key Laboratory of Multimodal Artificial Intelligence Systems (No.202200042).

⁴<https://anonymous.4open.science/w/DEMO-ICASSP2024-5A69/>

6. REFERENCES

- [1] Kai Shen, Zeqian Ju, Xu Tan, Yanqing Liu, Yichong Leng, Lei He, Tao Qin, Sheng Zhao, and Jiang Bian, “Naturalspeech 2: Latent diffusion models are natural and zero-shot speech and singing synthesizers,” *arXiv preprint arXiv:2304.09116*, 2023.
- [2] Jonathan Shen, Ruoming Pang, Ron J. Weiss, Mike Schuster, Navdeep Jaitly, Zongheng Yang, Zhifeng Chen, Yu Zhang, Yuxuan Wang, RJ-Skerrv Ryan, Rif A. Saurous, Yannis Agiomyriannakis, and Yonghui Wu, “Natural TTS synthesis by conditioning wavenet on MEL spectrogram predictions,” in *ICASSP*. 2018, pp. 4779–4783, IEEE.
- [3] Yi Ren, Chenxu Hu, Xu Tan, Tao Qin, Sheng Zhao, Zhou Zhao, and Tie-Yan Liu, “Fastspeech 2: Fast and high-quality end-to-end text to speech,” in *ICLR*. 2021, OpenReview.net.
- [4] Jaehyeon Kim, Jungil Kong, and Juhee Son, “Conditional variational autoencoder with adversarial learning for end-to-end text-to-speech,” in *ICML*. 2021, vol. 139 of *Proceedings of Machine Learning Research*, pp. 5530–5540, PMLR.
- [5] Greg J Stephens, Lauren J Silbert, and Uri Hasson, “Speaker-listener neural coupling underlies successful communication,” *Proceedings of the National Academy of Sciences*, vol. 107, no. 32, pp. 14425–14430, 2010.
- [6] Bingjiang Lyu, Hun S Choi, William D Marslen-Wilson, Alex Clarke, Billi Randall, and Lorraine K Tyler, “Neural dynamics of semantic composition,” *Proceedings of the National Academy of Sciences*, vol. 116, no. 42, pp. 21318–21327, 2019.
- [7] Pilar Oplustil Gallegos, Johannah O’Mahony, and Simon King, “Comparing acoustic and textual representations of previous linguistic context for improving text-to-speech,” in *The 11th ISCA Speech Synthesis Workshop (SSW11)*. ISCA, 2021, pp. 205–210.
- [8] Shun Lei, Yixuan Zhou, Liyang Chen, Jiankun Hu, Zhiyong Wu, Shiyin Kang, and Helen Meng, “Towards multi-scale speaking style modelling with hierarchical context information for mandarin speech synthesis,” *arXiv preprint arXiv:2204.02743*, 2022.
- [9] Vivek Kumar Rangarajan Sridhar, Ann K. Syrdal, Alistair Conkie, and Srinivas Bangalore, “Enriching text-to-speech synthesis using automatic dialog act tags,” in *INTERSPEECH*. 2011, pp. 317–320, ISCA.
- [10] Keon Lee, Kyumin Park, and Daeyoung Kim, “Dailytalk: Spoken dialogue dataset for conversational text-to-speech,” *arXiv preprint arXiv:2207.01063*, 2022.
- [11] Jinlong Xue, Yayue Deng, Fengping Wang, Ya Li, Yingming Gao, Jianhua Tao, Jianqing Sun, and Jiaen Liang, “M2-ctts: End-to-end multi-scale multi-modal conversational text-to-speech synthesis,” in *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2023, pp. 1–5.
- [12] Haohan Guo, Shaofei Zhang, Frank K Soong, Lei He, and Lei Xie, “Conversational end-to-end tts for voice agents,” in *2021 IEEE Spoken Language Technology Workshop (SLT)*. IEEE, 2021, pp. 403–409.
- [13] Yuto Nishimura, Yuki Saito, Shinnosuke Takamichi, Kentaro Tachibana, and Hiroshi Saruwatari, “Acoustic modeling for end-to-end empathetic dialogue speech synthesis using linguistic and prosodic contexts of dialogue history,” *arXiv preprint arXiv:2206.08039*, 2022.
- [14] Yifan Hu, Rui Liu, Guanglai Gao, and Haizhou Li, “Fctalker: Fine and coarse grained context modeling for expressive conversational speech synthesis,” *CoRR*, vol. abs/2210.15360, 2022.
- [15] Jingbei Li, Yi Meng, Chenyi Li, Zhiyong Wu, Helen Meng, Chao Weng, and Dan Su, “Enhancing speaking styles in conversational text-to-speech synthesis with graph-based multi-modal context modeling,” in *ICASSP*. 2022, pp. 7917–7921, IEEE.
- [16] Jingbei Li, Yi Meng, Xixin Wu, Zhiyong Wu, Jia Jia, Helen Meng, Qiao Tian, Yuping Wang, and Yuxuan Wang, “Inferring speaking styles from multi-modal conversational context by multi-scale relational graph convolutional networks,” in *ACM Multimedia*. 2022, pp. 5811–5820, ACM.
- [17] Florian Schroff, Dmitry Kalenichenko, and James Philbin, “Facenet: A unified embedding for face recognition and clustering,” in *CVPR*. 2015, pp. 815–823, IEEE Computer Society.
- [18] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova, “BERT: pre-training of deep bidirectional transformers for language understanding,” in *NAACL-HLT (1)*. 2019, pp. 4171–4186, Association for Computational Linguistics.
- [19] Yi Ren, Ming Lei, Zhiying Huang, Shiliang Zhang, Qian Chen, Zhijie Yan, and Zhou Zhao, “Prosospeech: Enhancing prosody with quantized vector pre-training in text-to-speech,” in *ICASSP*. 2022, pp. 7577–7581, IEEE.
- [20] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin, “Attention is all you need,” in *NIPS*, 2017, pp. 5998–6008.
- [21] Laurens Van der Maaten and Geoffrey Hinton, “Visualizing data using t-sne.,” *Journal of machine learning research*, vol. 9, no. 11, 2008.
- [22] Sumit Chopra, Raia Hadsell, and Yann LeCun, “Learning a similarity metric discriminatively, with application to face verification,” in *CVPR (1)*. 2005, pp. 539–546, IEEE Computer Society.
- [23] Meinard Müller, “Dynamic time warping,” *Information retrieval for music and motion*, pp. 69–84, 2007.