

# QI-TTS: QUESTIONING INTONATION CONTROL FOR EMOTIONAL SPEECH SYNTHESIS

Haobin Tang<sup>1,2†</sup>, Xulong Zhang<sup>1†</sup>, Jianzong Wang<sup>1\*</sup>, Ning Cheng<sup>1</sup>, Jing Xiao<sup>1</sup>

<sup>1</sup>Ping An Technology (Shenzhen) Co., Ltd.

<sup>2</sup>University of Science and Technology of China

## ABSTRACT

Recent expressive text to speech (TTS) models focus on synthesizing emotional speech, but some fine-grained styles such as intonation are neglected. In this paper, we propose QI-TTS which aims to better transfer and control intonation to further deliver the speaker’s questioning intention while transferring emotion from reference speech. We propose a multi-style extractor to extract style embedding from two different levels. While the sentence level represents emotion, the final syllable level represents intonation. For fine-grained intonation control, we use relative attributes to represent intonation intensity at the syllable level. Experiments have validated the effectiveness of QI-TTS for improving intonation expressiveness in emotional speech synthesis.

**Index Terms**— Emotional speech synthesis, Intonation intensity control, Relative attribute

## 1. INTRODUCTION

Due to the rapid advancement of seq2seq modeling architecture, style transfer TTS [1, 2, 3] has become a prevalent approach for emotional speech synthesis in recent years. The approach utilizes reference audio to specify the desired speech style and its intention is to generate speech that emulates the style of the reference audio. Reference-based [1] style transfer involves unsupervised learning of a fixed-length style embedding through expressive samples, which is then utilized to model the speaking style of a reference audio. Style transfer has seen significant progress in recent times, with the emergence of numerous approaches such as Global Style Tokens (GST) [2], Variational Autoencoder (VAE) [4, 5] and their variants.

However, merely utilizing a learned sentence level style embedding in emotion transfer is insufficient for fully expressing a speaker’s attitude. The emotion embedding lacks the ability to model a combination of emotion and mutually exclusive intonation, such as “angry statement” and “angry question” accurately at the same time. By using questioning intonation, speaker can express a statement in the form of declarative question [6]. For example, “You failed the exam

this time.” to “You failed the exam this time?”. Thus, it is inadequate to model speech prosody from a single aspect and intonation is essential for intention clarification. There are three issues we would like to address: 1) The existing emotion modeling frameworks only consider the normal statement and lack of the ability to model multiple and differential prosody such as questions in each emotion; 2) The intonation expressions that exist in human language are nuanced and vary in intensity. Thus we desire to flexibly deliver questioning intonation with specific intensity; 3) limitation of the ability to disentangle prosody from other attributes like content, resulting in the quality degrade and expressiveness instability.

We propose QI-TTS, a model built upon non-autoregressive TTS system FastSpeech2 [7] with the following specific designs for above problems. We propose a multi-style extractor that extracts emotion features from sentence level while extracting intonation features from final syllable level to model intonation and emotion at the same time. By utilizing relative attribute [8], we are able to assign a relative degree of intonation strength to each audio samples in dataset via a ranking function. To minimize the information overlap between content embedding and style embedding, we utilized a gradient reversal layer (GRL) together with content predictor.

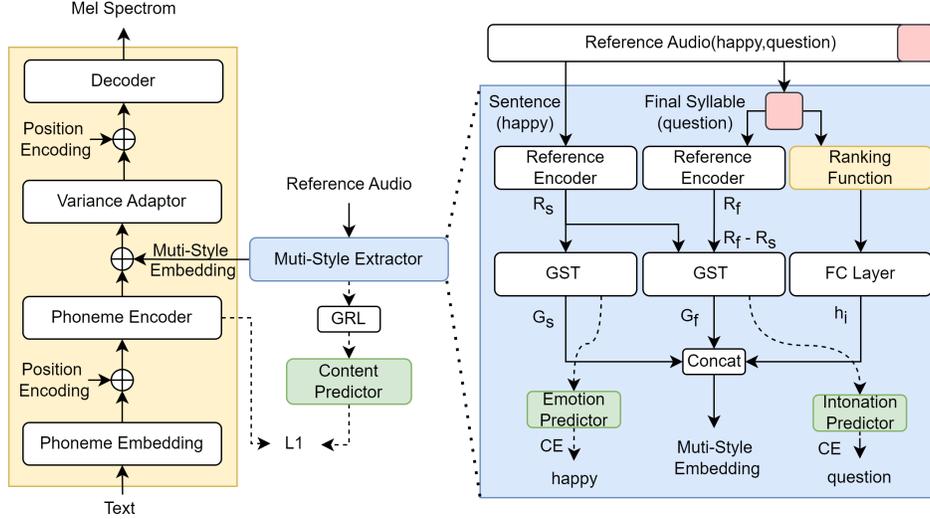
The paper presents the following contributions: 1) We jointly transfer the emotion and intonation from reference audio in an end-to-end way which further delivers the speaker’s intention; 2) QI-TTS is capable of learning the variation of intonation intensity in speech without the use of explicit labels, making it possible to control the intonation intensity effectively using either manual instructions or reference speech; 3) Experiments have validated the effectiveness of QI-TTS for expressive style transfer and intonation intensity control.

## 2. PROPOSED METHOD

As Fig. 1 depicts, the overall architecture of QI-TTS can be mainly divided into two parts based on FastSpeech 2 [7], a multi-style extractor with ranking function, and a content predictor with gradient reversal layer (GRL) [9]. Specifically, the extractor is used to extract emotion and intonation embedding at sentence and final syllable level. The ranking function outputs a relative intonation intensity and encodes it into an intonation intensity embedding which is concatenated with emo-

† Equal Contribution

\*Corresponding author: Jianzong Wang, jzwang@188.com.



**Fig. 1:** The overview architecture for QI-TTS. The red part of reference audio refers to final syllable. "GRL" denotes gradient reversal layer.  $R_s$  and  $R_f$  denote the reference embedding of the reference sentence and that of the final syllable, respectively.  $G_s$  and  $G_f$  denote emotion and intonation embedding.  $h_i$  represents intonation intensity embedding.

tion and intonation embedding to form multi-style embedding. The GRL content predictor is used to further disentangle content information from multi-style embedding.

## 2.1. Multi-Style Extractor

Linguists find that the duration and fundamental frequency in sentence-final differ from that in other positions of the utterance [10]. The duration of final syllable is 1.53 times longer than that of non-final syllable in English [11]. The final syllable's f0 valley is a significant feature of a statement. Furthermore, the absolute value of the endpoint F0 was the strongest cue in distinguishing statements from questions rather than the slope of the terminal glide [12]. Research [10] shows that the mean of final syllable in English is 0.37 seconds with a standard deviation of 0.15. So we model the last 0.52 seconds of the audio that contains the final syllable as intonation to capture the duration variance and intonation related features.

We designed a multi-style extractor, depicted in the blue section of Fig. 1, to extract the emotion and intonation information from the reference speech. The module contains reference encoders and style token layers that follow the structure proposed by [2]. Both the Mel-spectrogram of the sentence and the Mel-spectrogram of its final syllable are fed into their corresponding reference encoder.  $R_s$  and  $R_f$  denote the output sentence reference embedding and final syllable reference embedding respectively. However, the effect of final syllable reference embedding overlaps with that of sentence reference embedding, since the relationship between emotion and intonation is not completely hierarchical. To reduce such overlapping, we use  $R_f - R_s$  to represent residual embedding of final syllable level. The emotion style embedding  $G_s$

and intonation style embedding  $G_f$  are formed by passing  $R_s$  and  $R_f - R_s$  to the corresponding style token layers. These embeddings are then concatenated with a ranking embedding that represents intonation intensity to create a multi-style embedding.

## 2.2. Modelling Intonation Intensity

We aim to conduct fine grained control of questioning intonation but the intensity labels are not readily available. A ranking-based method called relative attributes [8] is used for unsupervised intensity modelling. We regard the questioning intensity as a speech attribute which can be depicted by learned relative attributes. The questioning intensity of a sentence should be zero since it lacks any questioning intonation variation. Therefore, we consider questioning intensity to be a relative difference between statement and question. Assuming we have a training set  $T = X_t$ , and  $X_t$  denotes the  $t^{th}$  training sample's acoustic features.  $A$  and  $B$  are the question and statement set respectively. We aim to learn the following ranking function:

$$f(X_t) = WX_t \quad (1)$$

where  $W$  is a weighting matrix that we need to learn. Considering the intonation intensity of question should always higher than that of statement we have to satisfy the following constraints:

$$\begin{aligned} \forall (X_a \in A \text{ and } X_b \in B) : WX_a > WX_b \\ \forall (X_a, X_b) \in A \text{ or } (X_a, X_b) \in B : WX_a = WX_b \end{aligned} \quad (2)$$

To estimate the weighting matrix  $W$ , we solve the following problem [13]:

$$\begin{aligned} & \min_W \left( \frac{1}{2} \|W\|_2^2 + C \left( \sum \xi_{a,b}^2 + \sum \gamma_{a,b}^2 \right) \right) \\ \text{s.t. } & W(X_a - X_b) \geq 1 - \xi_{a,b}; \forall (X_a \in A \text{ and } X_b \in B) \\ & |W(X_a - X_b)| \leq \gamma_{a,b}; \forall (X_a, X_b) \in A \text{ or } (X_a, X_b) \in B \\ & \xi_{a,b} \geq 0; \gamma_{a,b} \geq 0 \end{aligned} \quad (3)$$

where  $C$  is the trade-off between the margin and the size of slack variables  $\xi_{a,b}^2$  and  $\gamma_{a,b}^2$ . Once the relative ranking function  $f(x)$  is trained, a normalized relative attribute in range  $[0, 1]$  can be calculated for a speech sample. This attribute is subsequently fed to a FC layer, which yields intensity embedding  $h_i$ . Fine-grained intonation control can be achieved during the inference stage through the use of reference speech or manual instructions. Specifically, the intensity can be predicted by analyzing the reference audio or assigned a value manually within the interval of  $[0, 1]$ .

### 2.3. Prediction Tasks

The multi-style extractor unsupervisedly learns multi-level styles from reference audio. We add emotion and intonation prediction tasks to force each level module to pay more attention to learning the corresponding style. The cross-entropy (CE) loss is used for emotion predictor. We use the weighted cross entropy function as intonation loss function because of the sparse question labels:

$$L = -\hat{y}_1 \log y_1 - \sigma \hat{y}_2 \log y_2 \quad (4)$$

where  $[\hat{y}_1, \hat{y}_2]$  is the probability of the ‘‘statement’’ and ‘‘question’’ categories respectively. One-hot encoding of ground-truth label is represented by  $[y_1, y_2]$ . We adjust  $\sigma$  to ensure the data balance in the training.

An adversarial content predictor network inspired by Mask-And-Predict (MAP) [14] is designed to disentangle overlapped content information in multi-style embedding. The network is composed of a gradient reverse layer and a content predictor [15]. The predictor is trained to predict content representation as accurate as possible by minimizing the loss:  $L_{content} = \|\hat{c} - c\|_1$ , where the content embedding generated by the phoneme encoder and the output of the adversarial content predictor are denoted as  $c$  and  $\hat{c}$ , respectively. The gradient is reversed before backward propagated to the multi-style extractor to minimize the content information contained in the multi-style embedding.

## 3. EXPERIMENT

### 3.1. Training Setting

To train our model, we use the english part of ESD dataset [16], which consists of five emotions spoken by 10 native English

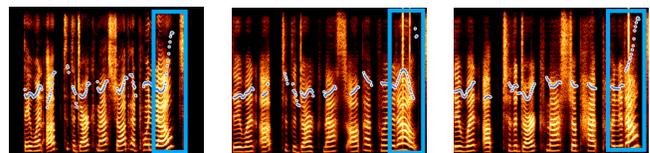
(5 male and 5 female): *Neutral, Sad, Happy, Angry, and Surprise*. We adopt the same data partition as provided in ESD. Importantly, we add ‘‘statement’’ and ‘‘question’’ labels with the help of K-means method similar to Into-TTS [17] and there are 310 questions for each speaker on average.

We follow an a public version of relative attribute [8] to train the ranking function for questioning intonation intensity. In practice, final syllable’s acoustic features of statements and questions are used for calculating intensity. The acoustic features are extracted by the openSMILE [18]. We train QI-TTS for 350k iterations with 16 batch size. In our experiments, we employ HiFi-Gan [19] as the vocoder for synthesizing waveforms from the generated Mel-spectrograms.

### 3.2. Result Evaluation

Our comparison involves the audio samples generated by QI-TTS and the systems listed below: 1) GT, Mel-spectrogram of reference audio; 2) GT mel + Vocoder, speech samples generated by HiFi-GAN using ground truth Mel-spectrogram; 3) MultiEmo FS2 [20], which adds the emotion d-vectors to Fastspeech 2. ‘‘FS2’’ denotes Fastspeech 2; 4) Styler [21], a speech synthesis model that employs speech decomposition to represent the style and achieve expressiveness.

For subject evaluation, we assess the quality and expressiveness of the synthesized audios via mean opinion score (MOS) and similarity mean opinion score (SMOS) tests. Subjects are asked to rate 5 question and 5 statement speech in each emotion on a scale from 1 to 5 with 1 point intervals. Besides, the accuracy of intonation perception is taken to evaluate the intonation transfer. Participants are tasked with assessing whether the audio sample is a question or a statement. We calculate the objective matrixes (*i.e.*, mel cepstral distortion (MCD) [22], F0 Frame Error (FFE), and Duration MSE) for objective evaluation. MCD evaluates the spectrum similarity while FFE is the proportion of frames with error F0. MSE between the ground-truth duration and predicted one is computed for duration. As depicts in Table 1, QI-TTS is capable of producing speech samples that more closely resemble the intonation of the reference audio with no hinder to emotion transfer, clearly reflecting the correct pitch and duration.



(a) Ground Truth (b) w/o final syllable (c) QI-TTS

**Fig. 2:** Visualizations of a declarative question’s Mel-spectrograms in ablation study.

**Table 1:** Subjective and objective evaluation results.

Model	MOS $\uparrow$	SMOS $\uparrow$	Intonation $\uparrow$	MCD $\downarrow$	FFE $\downarrow$	Duration MSE $\downarrow$
GT	4.47 $\pm$ 0.08	/	/	/	/	/
GTmel + Vocoder	4.40 $\pm$ 0.09	4.47 $\pm$ 0.10	99.2%	2.40	0.07	0.031
MutiEmo FS2 [20]	3.81 $\pm$ 0.08	3.85 $\pm$ 0.08	81.6%	<b>3.15</b>	0.43	0.144
Styler [21]	3.76 $\pm$ 0.08	3.97 $\pm$ 0.08	85.9%	5.57	0.41	0.149
QI-TTS	<b>3.84 <math>\pm</math> 0.10</b>	<b>4.01 <math>\pm</math> 0.08</b>	<b>95.2%</b>	4.89	<b>0.39</b>	<b>0.141</b>

### 3.3. Ablation Study

The effectiveness of techniques employed in QI-TTS, such as final syllable level style, residual style embedding, and predictors, were demonstrated through ablation studies. To compare the expressiveness of the synthesized speeches, CSMOS (comparative similarity mean opinion score) results are presented in Table 2. Fig. 2 shows that the model without final syllable level performs wrong intonation while the QI-TTS performs well indicating the importance of modeling intonation style representation in final syllable level. The absence of intonation related function results in more severe decline in question than in statement which demonstrates the efficiency in modeling questioning intonation. Moreover, removing the residual method results in -0.08 CSMOS suggesting that minimizing interference and overlap between emotion and intonation styles is necessary. The degradation of audio similarity also demonstrates the efficacy of the predictors in accurately modeling specific emotions and intonations while reducing the impact of textual content.

**Table 2:** CSMOS comparison for ablation study.

Model	Question	Statement
QI-TTS	/	/
w/o final syllable level	-0.15	-0.09
w/o residual style	-0.08	-0.08
w/o Emotion predictor	-0.10	-0.10
w/o Intonation predictor	-0.11	-0.04
w/o GRL content predictor	-0.08	-0.09

### 3.4. Intonation Intensity Control

To evaluate the intonation intensity control, three distinct intensity values were chosen: 0.3, 0.6, and 0.9 which correspond to weak, medium, and strong intensity levels, respectively. Best-worst scaling (BWS) test is conducted to evaluate intonation and emotion perception of generated audio with different questioning intensities in each emotion category. The assessors are requested to select the speech sample that best and worst represents a specific emotion or intonation based on their perception. There is little difference between the experimental results of angry, happy, and sad, so we show angry and surprise only in Table 3.

We first evaluate the perception of the questioning intonation. According to Table 3a, the highest “Best” score for question intonation always occurs at 90% questioning intensity, and this score increases with increasing intensity. Conversely, the highest “Worst” score is always observed at the lowest percentage of questioning intonation. These findings provide

empirical support for the effectiveness of controlling intonation intensity. We further evaluate emotion perception in synthesized speech to study the effect of intonation intensity control on emotional expression in Table 3b. The “Best” score of “Surprise” slightly increases as the percentage of questioning intonation intensity increases while that of “Angry” slightly decreases. We attribute this phenomenon to the potential link between surprises and questions. Stronger question helps express surprise, but excessive transfer confuses other emotions with surprise to some extent.

**Table 3:** Best-worst scaling (BWS) test for questioning intonation and emotion perception.

(a) Perception of questioning intonation

Configuration		Best(%)	Worst(%)
Surprise	30% Question	8	79
	60% Question	11	21
	90% Question	81	0
Angry	30% Question	8	69
	60% Question	15	31
	90% Question	77	0

(b) Perception of emotion

Configuration		Best(%)	Worst(%)
Surprise	30% Question	29	39
	60% Question	34	33
	90% Question	37	28
Angry	30% Question	39	28
	60% Question	40	27
	90% Question	21	45

## 4. CONCLUSION

This paper proposes QI-TTS, a multi-style transfer model, which can better transfer emotion and intonation from reference audio and achieve intonation intensity control for expressive TTS. Experimental results demonstrate that QI-TTS performs better on expressive speech synthesis and intonation intensity control missions. In future research, we will explore the effectiveness of QI-TTS in multilingual scenarios.

## 5. ACKNOWLEDGEMENT

Supported by the Key Research and Development Program of Guangdong Province (grant No. 2021B0101400003) and corresponding author is Jianzong Wang (jzwang@188.com).

## 6. REFERENCES

- [1] RJ Skerry-Ryan, Eric Battenberg, Ying Xiao, Yuxuan Wang, Daisy Stanton, Joel Shor, Ron Weiss, Rob Clark, and Rif A Saurous, “Towards end-to-end prosody transfer for expressive speech synthesis with tacotron,” in *ICML*, 2018, pp. 4693–4702.
- [2] Yuxuan Wang, Daisy Stanton, Yu Zhang, RJ-Skerry Ryan, Eric Battenberg, Joel Shor, Ying Xiao, Ye Jia, Fei Ren, and Rif A Saurous, “Style tokens: Unsupervised style modeling, control and transfer in end-to-end speech synthesis,” in *ICML*, 2018, pp. 5180–5189.
- [3] Yanyao Bian, Changbin Chen, Yongguo Kang, and Zhenglin Pan, “Multi-reference tacotron by intercross training for style disentangling, transfer and control in speech synthesis,” *arXiv preprint arXiv:1904.02373*, 2019.
- [4] Alexander H Liu, Tao Tu, Hung-yi Lee, and Lin-shan Lee, “Towards unsupervised speech recognition and synthesis with quantized speech representation learning,” in *ICASSP. IEEE*, 2020, pp. 7259–7263.
- [5] Diederik P Kingma and Max Welling, “Auto-encoding variational bayes,” *arXiv preprint arXiv:1312.6114*, 2013.
- [6] Qibing Bai, Tom Ko, and Yu Zhang, “A study of modeling rising intonation in cantonese neural speech synthesis,” in *Interspeech*, 2022, pp. 501–505.
- [7] 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 *ICML*, 2020.
- [8] Devi Parikh and Kristen Grauman, “Relative attributes,” in *2011 International Conference on Computer Vision (ICCV)*. IEEE, 2011, pp. 503–510.
- [9] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky, “Domain-adversarial training of neural networks,” *The journal of machine learning research*, vol. 17, no. 1, pp. 2096–2030, 2016.
- [10] Rochele Berkovits, “Duration and fundamental frequency in sentence-final intonation,” *Journal of Phonetics*, vol. 12, no. 3, pp. 255–265, 1984.
- [11] Pierre Delattre, “A comparison of syllable length conditioning among languages,” *International Review of Applied Linguistics*, vol. 4, pp. 183–198, 1966.
- [12] Wojciech Majewski and Richard Blasdell, “Influence of fundamental frequency cues on the perception of some synthetic intonation contours,” *The Journal of the Acoustical Society of America*, vol. 45, no. 2, pp. 450–457, 1969.
- [13] Olivier Chapelle, “Training a support vector machine in the primal,” *Neural computation*, vol. 19, no. 5, pp. 1155–1178, 2007.
- [14] Jie Wang, Jingbei Li, Xintao Zhao, Zhiyong Wu, Shiyin Kang, and Helen Meng, “Adversarially learning disentangled speech representations for robust multi-factor voice conversion,” in *Interspeech*, 2021, pp. 846–850.
- [15] Andy T Liu, Shu-wen Yang, Po-Han Chi, Po-chun Hsu, and Hung-yi Lee, “Mockingjay: Unsupervised speech representation learning with deep bidirectional transformer encoders,” in *ICASSP. IEEE*, 2020, pp. 6419–6423.
- [16] Kun Zhou, Berrak Sisman, Rui Liu, and Haizhou Li, “Seen and unseen emotional style transfer for voice conversion with a new emotional speech dataset,” in *ICASSP. IEEE*, 2021, pp. 920–924.
- [17] Jihwan Lee, Joun Yeop Lee, Heejin Choi, Seongkyu Mun, Sangjun Park, and Chanwoo Kim, “Into-tts: Intonation template based prosody control system,” *arXiv preprint arXiv:2204.01271*, 2022.
- [18] Florian Eyben, Martin Wöllmer, and Björn Schuller, “Opensmile: the munich versatile and fast open-source audio feature extractor,” in *Proceedings of the 18th ACM international conference on Multimedia*, 2010, pp. 1459–1462.
- [19] Jungil Kong, Jaehyeon Kim, and Jaekyoung Bae, “Hifigan: Generative adversarial networks for efficient and high fidelity speech synthesis,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 17022–17033, 2020.
- [20] Chenye Cui, Yi Ren, Jinglin Liu, Feiyang Chen, Rongjie Huang, Ming Lei, and Zhou Zhao, “Emovie: A mandarin emotion speech dataset with a simple emotional text-to-speech model,” in *Interspeech*, 2021, pp. 2766–2770.
- [21] Keon Lee, Kyumin Park, and Daeyoung Kim, “Styler: Style factor modeling with rapidity and robustness via speech decomposition for expressive and controllable neural text to speech,” in *Interspeech*, 2021, pp. 3431–3435.
- [22] Robert Kubichek, “Mel-cepstral distance measure for objective speech quality assessment,” in *Proceedings of IEEE pacific rim conference on communications computers and signal processing*. IEEE, 1993, vol. 1, pp. 125–128.