

# BISINGER: BILINGUAL SINGING VOICE SYNTHESIS

Huali Zhou<sup>1,2\*</sup>, Yueqian Lin<sup>2\*</sup>, Yao Shi<sup>2</sup>, Peng Sun<sup>2</sup>, Ming Li<sup>1,2†</sup>

<sup>1</sup>School of Computer Science, Wuhan University, Wuhan, China

<sup>2</sup>Suzhou Municipal Key Laboratory of Multimodal Intelligent Systems,  
Duke Kunshan University, Kunshan, China

## ABSTRACT

Although Singing Voice Synthesis (SVS) has made great strides with Text-to-Speech (TTS) techniques, multilingual singing voice modeling remains relatively unexplored. This paper presents BiSinger, a bilingual pop SVS system for English and Chinese Mandarin. Current systems require separate models per language and cannot accurately represent both Chinese and English, hindering code-switch SVS. To address this gap, we design a shared representation between Chinese and English singing voices, achieved by using the CMU dictionary with mapping rules. We fuse monolingual singing datasets with open-source singing voice conversion techniques to generate bilingual singing voices while also exploring the potential use of bilingual speech data. Experiments affirm that our language-independent representation and incorporation of related datasets enable a single model with enhanced performance in English and code-switch SVS while maintaining Chinese song performance. Audio samples are available at <https://bisinger-svs.github.io>.

**Index Terms**— singing voice synthesis, bilingual singing modeling, code-switch, dataset adaptation, signal processing

## 1. INTRODUCTION

Singing voice synthesis (SVS) is becoming increasingly popular in our daily lives. It aims to create natural and expressive singing voices that match the music score. SVS is a unique type of Text-to-Speech (TTS) task because it must strictly adhere to pitch and duration limitations in the scores. Additionally, obtaining the necessary training data for SVS is more challenging due to copyright restrictions and complex annotation requirements.

The development of the SVS system is in parallel with the progressive development of the TTS system. Xiaoic-eSing [1] uses the FastSpeech [2] network to generate high-quality singing voices, while ByteSing [3] employs the autoregressive Tacotron-like [4] structure as the acoustic model. With the emergence of the end-to-end TTS model, VITS [5] and VISinger [6] were also proposed, which can effectively mitigate the two-stage mismatch problem in singing voice generation. Multiple open-source toolkits, e.g., Sinsy [7], Muskits [8], and NNSVS [9], etc., were released to boost the

development of SVS research. The diffusion probabilistic-based model, DiffSinger [10], has recently demonstrated superior performance to alternative methods. However, these systems mainly target monolingual pop songs, assuming the input lyrics are from a single language. There has been little exploration into multilingual singing voice modeling, which is essential because mixed language lyrics are prevalent in real singing. Our work seeks to address this issue by focusing on Chinese-English bilingual singing voice synthesis.

To develop a bilingual SVS model, it would be ideal to use a bilingual singing corpus with detailed music annotations. Unfortunately, the Children’s Song Dataset [11] is one of the few bilingual datasets available, consisting of Korean and English. Regrettably, there is currently no publicly available Chinese-English singing dataset, and collecting one would be time-consuming and challenging due to the need for bilingual singers and manual annotations. To overcome this challenge, on the one hand, we propose using existing monolingual singing corpora, specifically M4Singer [12] for Chinese and NUS-48E [13] for English, with a language-independent representation [14] to build our model. On the other hand, our research explores how to synthesize bilingual singing voices with the help of a bilingual speech dataset, DB-4 from Data Baker<sup>1</sup>. The highlights of our contributions are as follows:

- We study how Chinese and English singing voices can have a shared representation to learn similar pronunciations crossing languages with annotation adaption.
- Our proposed approach converts existing monolingual singing datasets with established singing voice conversion (SVC) techniques to create bilingual singing voices.
- Considering the rich resource of speech data, we also look into the possibility of developing bilingual singing voices using a bilingual speech database.

The paper is structured as follows. Section 2 covers the related works on SVS and multilingual speech synthesis, while Section 3 presents our methodology for bilingual SVS. Experiment setup and results can be found in Section 4. Finally, Section 5 provides the conclusion and future work.

## 2. RELATED WORKS

### 2.1. Singing Voice Synthesis

Singing voice synthesis (SVS) is a process that generates high-quality singing voices using music scores with lyrics for pro-

\* Equal contribution.

† Corresponding author. E-mail: [ming.li369@duke.edu](mailto:ming.li369@duke.edu)

<sup>1</sup> <https://www.data-baker.com/en>

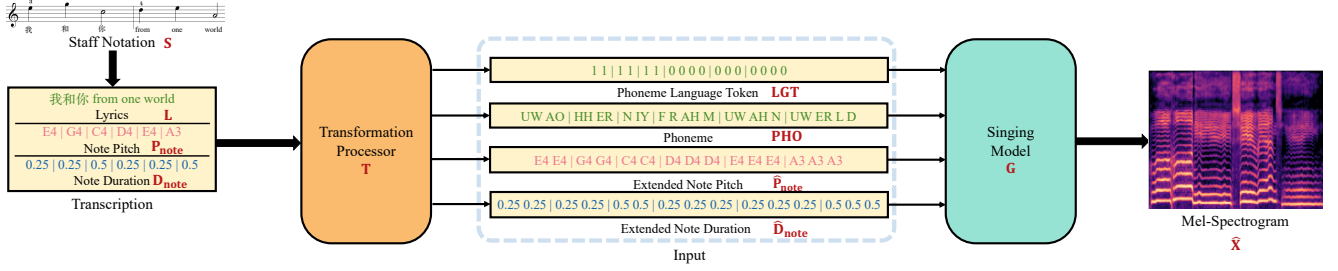


Fig. 1: System overview.

nunciation and notes for prosody. Similar to the TTS task, it has evolved from the unit-selection method [15, 16], which concatenates short waveform units from the database to more advanced statistical parametric systems like Hidden Markov Model (HMM)-based systems [7, 17], and now to mainstream deep neural network (DNN)-based architectures [1, 3, 6, 8–10]. SVS typically uses a two-stage pipeline, consistent with the TTS task, where the acoustic model learns to map music score inputs to acoustic features, and the vocoder reconstructs the waveform based on the predicted acoustic feature.

Some work involves extending existing singing voice synthesis systems to accommodate other languages. In [17], they adapted the HMM-based method for Japanese SVS to work with English by developing language-independent contexts. They also investigated new syllable allocation and duplication methods, considering the distinct syllable structures of Japanese and English. This work was aided by data from a bilingual singer, which provided informative insights, but unfortunately, the available data was quite scarce [11]. With multi-singer Mandarin, English, and Cantonese singing data mined from music websites, [18] could perform cross-lingual singing voice synthesis as a byproduct of multilingual training. However, it should be noted that systems in [18] assume the inputs are in a single language. For existing systems, such as Muskits [8], multilingual synthesis requires separate models for different languages, making it incapable of synthesizing code-switch singing voices within a single model.

## 2.2. Multilingual Speech Synthesis

Multilingual speech synthesis is designed to alleviate the need for training separate models for different languages and support low-resource synthesis. In order to promote the knowledge-sharing capacity of the model for different languages, a unified input representation [19] of various languages has been sought for a long time. In [20], the Unicode bytes were used across languages owing to their language independence and fixed size. Additionally, [21] explored both a shared multilingual encoder with language embedding and a separate monolingual encoder. Another approach [22] involved using an additional network conditioned on language to generate parameters for multiple language-dependent encoders to enable cross-lingual knowledge-sharing. Cross-lingual synthesis was explored in [14], proposing common phonemic representations linked to numeric language ID codes, which

provides substantial inspiration to our research.

To preserve different languages’ characteristics while sharing knowledge with shared phoneme inputs, [23, 24] introduced stress and tone embedding. Experimental observations in [25] suggested that there is some degree of shared pronunciation across languages, which was helpful in low-resource scenarios. However, the close pronunciation between Chinese and English could lead to mutual interference, such as intonation variation and mispronunciation. To address this issue, [26] proposed an embedding strength modulator to capture the dynamic strength of language and phonology, which has also been incorporated into our work.

## 3. METHODOLOGY

### 3.1. Formulation of the Model

Fig. 1 presents an overview of our system. Given a staff notation  $S$ , our first step is to transcribe it into lyrics  $L$ , note pitch  $P_{\text{note}}$ , and note duration  $D_{\text{note}}$ . Subsequently, these transcribed elements are processed through a transformation processor  $T$  to format the data for the singing model  $G$ . The detailed transformation process is illustrated as follows:

$$\text{LGT}, \text{PHO}, \hat{P}_{\text{note}}, \hat{D}_{\text{note}} = T(L, P_{\text{note}}, D_{\text{note}}) \quad (1)$$

where  $\text{LGT}$  and  $\text{PHO}$  are the language tokens and universal phoneme representations generated according to the lyrics,  $\hat{P}_{\text{note}}$  is the phoneme’s corresponding note pitch that is extended for phoneme-level alignment, and similarly,  $\hat{D}_{\text{note}}$  is the repeated input note duration.

After obtaining phoneme-level data, we then feed it into end-to-end singing model  $G$  to predict the Mel-spectrogram  $\hat{X}$ , as shown in Eq. 2, and calculate the loss according to  $\hat{X}$ , continually updating itself throughout the training stage.

$$\hat{X} = G(\text{LGT}, \text{PHO}, \hat{P}_{\text{note}}, \hat{D}_{\text{note}}) \quad (2)$$

Once trained, the model has the capability to predict the Mel-spectrogram  $\tilde{X}$  for any input music score, such as multilingual lyrics  $L_{\text{mul}}$ , note pitch data  $P_{\text{note}}$ , and note duration data  $D_{\text{note}}$ . This prediction can be expressed by Eq. 3:

$$\tilde{X} = G(T(L_{\text{mul}}, P_{\text{note}}, D_{\text{note}})) \quad (3)$$

### 3.2. Language-independent Representation

Inspired by [14], we adopt the CMU Pronunciation Dictionary<sup>2</sup> as the shared phoneme representation for Chinese and English languages to overcome the challenge of different grapheme or phoneme sets in multilingual SVS. For Mandarin, characters are first represented as Pinyin using Pypinyin<sup>3</sup>, which can be converted to CMU phonemes by the Pinyin-to-CMU mapping table<sup>4</sup>. Likewise, each English word can be converted into CMU phonemes by referring to the mapping table for English, with examples in Table 1. When it comes to singing, the pitch is primarily determined by the score rather than tone and accent, according to [3]. Therefore, tone or stress information is not taken into consideration. Furthermore, the code-switched song shares the same BPM across linguistic boundaries.

**Table 1:** Unit-to-CMU phoneme mapping examples.

Chinese Pinyin	CMU Phonemes	English Word	CMU Phonemes
rang	R AE NG	cat	K AE T
wo	W AO	fan	F AE N
nuan	N UW AE N	song	S AO NG
yang	Y AE NG	total	T OW T AH L
zhui	JH UW IY	story	S T AO R IY

It is worth noting that the pronunciation of the same phoneme can differ between Chinese and English. Due to this, we use language identification tokens to preserve language-dependent characteristics while sharing unified phonemes. The tokens ‘0’ and ‘1’ represent the language ID for each phoneme, with ‘0’ signifying English and ‘1’ indicating Chinese. Take the lyrics “我和你 from one world” as an example; two token sequences are obtained as language-independent representations. The first is the phoneme sequence “UW AO HH ER N IY F R AH M UW AH N UW ER L D,” and the second is the corresponding language token sequence “1 1 1 1 1 0 0 0 0 0 0 0 0 0,” which has the same length as the former.

### 3.3. Language-style-infused Encoder

We modify the encoder in DiffSinger [10] to suit our multilingual singing voice modeling with both speech and singing data, as illustrated in Fig. 2.

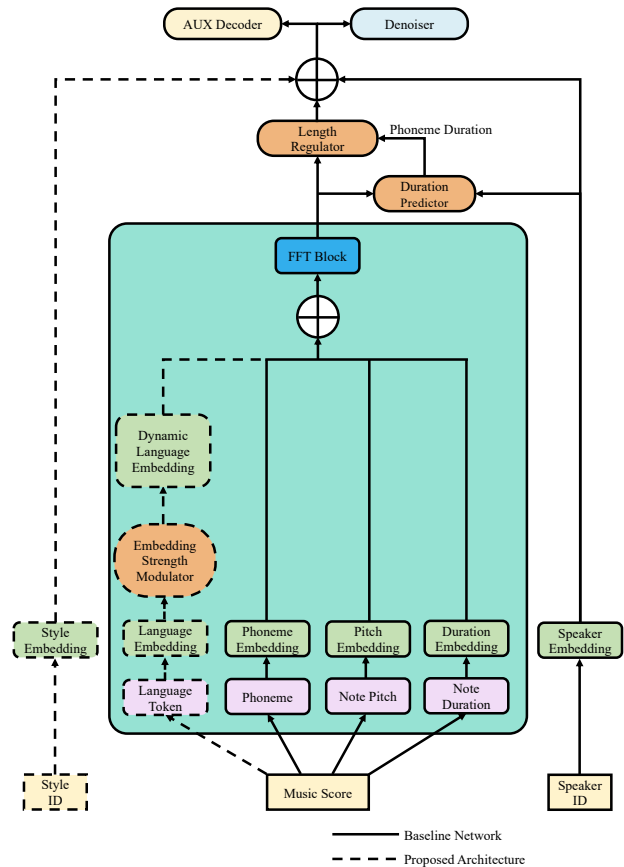
To deal with the multilingual issue, a learnable language embedding layer is introduced to convert the numeric language ID tokens to a 256-dimensional language embedding sequence, then encoded with the original phoneme embedding, note embedding and note duration embedding. Additionally, we employ an Embedding Strength Modulator (ESM) [26], a fusion of multi-head attention scheme and feed-forward network, to capture the dynamic strength of phonology and language, as shown in Fig. 2. When utilizing the ESM module, the feed-forward Transformer (FFT) block employs the dynamic language embedding that incorporates the encoded phoneme embedding instead of the static data directly calculated based on the language ID.

<sup>2</sup> <http://www.speech.cs.cmu.edu/cgi-bin/cmudict>

<sup>3</sup> <https://github.com/mozillazg/python-Pinyin>

<sup>4</sup> <https://github.com/kaldi-asr/kaldi/blob/master/egs/hkust/s5/conf/pinyin2cmu>

A style identification token for each utterance is introduced to address the training imbalance between speech and singing data and avoid the synthesized singing voices from sounding excessively smooth and fast like speech. This token consists of the numbers ‘0’, ‘1’, and ‘2’, with ‘0’ indicating speech, ‘1’ indicating singing, and ‘2’ indicating pseudo-singing obtained through the pitch shift method illustrated in Sec 3.4.3. As shown in Fig. 2, a style embedding layer is designed, converting the numeric style ID token to a 256-dimensional style embedding. Lastly, we merge the style embedding, speaker embedding, and the result of the encoded musical score to form the input for the auxiliary decoder or the denoiser.



**Fig. 2:** Language-style-infused encoder.

### 3.4. Adaptation of Datasets

As mentioned in Section 2, the scarcity of well-annotated singing data poses crucial challenges for SVS, especially multilingual ones. We propose three methods to address this challenge: phoneme-level musical annotation adaption, timbre conversion, and pseudo-singing through pitch shift.

#### 3.4.1. Phoneme-level Musical Annotation Adaptation

To fully utilize the existing datasets’ annotation information, we intend to adjust the original annotations to fit our CMU phoneme set. The note and note durations should be considered thoughtfully for singing voice datasets. We propose a

straightforward approach that involves splitting each phoneme in the Pinyin set into its corresponding CMU phonemes; for instance, ‘ang’ will be split into ‘AE NG.’ For the musical annotations, based on the mapping relationship, we duplicate the note and note duration from the original annotation for the corresponding CMU phonemes. Next, we divide the original phoneme duration equally among the split CMU phonemes and assign the average duration to each CMU phoneme.

Moreover, we noticed that simply averaging phoneme duration may pose issues since vowels in singing tend to be longer than consonants. The average strategy can negatively affect the naturalness of synthesized singing voices. To handle this, we present distributing phoneme duration proportionally based on the Montreal Forced Aligner (MFA) [27] coarse alignment. Specifically, we utilize MFA<sup>5</sup> to align with the CMU dictionary and obtain phoneme duration. However, the alignment results are not accurate enough when it comes to the long duration of vowels and portamento. We observed that the phoneme duration closely matches M4Singer’s original annotations for the initial components of Pinyin. The ratio between consonants and vowels appears reasonable for the final components. Bearing these observations in mind, we retain the same phoneme duration for Pinyin initials when converting them to CMU-based phonemes. For the Pinyin finals, we first convert them into CMU phonemes and then proportionally distribute the original duration based on the ratio obtained in the MFA alignment result. Lastly, we calculate each mapped CMU phoneme’s note and note duration based on adapted phoneme durations and original note annotations. A concrete example can be found in Table 2.

**Table 2:** Three annotation types.

original		average		proportional	
phs	ph_dur	phs	ph_dur	phs	ph_dur
c	0.18	T	0.09	T	0.036
uen	0.245	S	0.09	S	0.144
uen	0.295	UW	0.0817	UW	0.18
		AH	0.0817	AH	0.065
		N	0.0817	AH	0.115
		UW	0.0983	N	0.18
		AH	0.0983		
		N	0.0983		

### 3.4.2. Timbre Conversion Method

While the shared CMU phonemes can help our model learn English pronunciation from Chinese data to a degree, some phonemes, such as /TH/, /Y/, /IH/, /DH/, /V/, and /OY/, are absent in Chinese. Substituting these missing phonemes with phonetically similar ones results in less natural-sounding voices, reinforcing the necessity of English data for a comprehensive and accurate representation of English phoneme pronunciation.

Due to the limited scale of existing English singing voice datasets and the need for consistent timbre to accurately capture how phonemes are pronounced in both English and Chinese, we have come up with a new method that uses advanced

SVC techniques, so-vits-svc<sup>6</sup>, an open-source repository for robust singing voice conversion. Our approach involves the conversion of original singing voices from a small English singing dataset to match all the different timbres in a larger Chinese singing dataset, significantly expanding the English singing voice dataset while creating a smooth transition between both languages.

Table 3 illustrates our pitch adjustment methodology when converting between different voice types, namely: Bass, Baritone, Tenor, Alto, and Soprano. Our goal is to preserve the natural voice range of the target singer and reduce conversion artifacts. To achieve this, we adjust the pitch in semitones. For instance, conversion from Bass to Tenor incurs a pitch shift of +8 semitones, increasing the pitch of the original Bass voice to match the typical Tenor range.

**Table 3:** Pitch adjustments reference (in semitones).

From	To				
	Bass	Baritone	Tenor	Alto	Soprano
Bass	0	+4	+8	+12	+12
Baritone	-4	0	+4	+8	+8
Tenor	-8	-4	0	+4	+8
Alto	-12	-8	-4	0	+4
Soprano	-12	-8	-8	-4	0

### 3.4.3. Musical Annotation and Pitch Shift

With easy access to existing large-scale, high-quality speech corpora, we plan to explore using a bilingual speech corpus for the bilingual SVS task. We treat speech data as plain singing data, although it lacks some singing characteristics. Unfortunately, there is no musical annotation for speech corpus, which is essential for SVS. Likewise, we utilize MFA for speech alignment, with CMU phoneme durations and word boundaries as results. We then extract F0 using the open-source Parselmouth<sup>7</sup>, which is then averaged according to word boundaries to get musical notes. We can also calculate phoneme duration and note duration from the alignment results.

To ensure the synthesized singing voice is as rhythmic as possible, we adjust the pitch of speech data with the WORLD vocoder [28] to achieve data augmentation and prevent the synthesized voice from being too stable and lacking musical expressiveness. To clarify, we start by defining ten frequently used melody frequencies. With the help of WORLD, we extract pitch contour (F0), aperiodic spectral envelope (AP), and harmonic spectral envelope (SP) for each speech audio. Next, we randomly choose one melody and replace the original F0 with it while keeping the F0 contour length unchanged. We then use WORLD to synthesize a new pseudo-singing voice based on the shifted F0 and the original AP and SP. Lastly, we update the note information in the musical annotation to match the corresponding melody. This way, we obtain pseudo-singing voices rich in pitch variation from the original stationary speech data.

<sup>5</sup> We train separate MFA models for M4Singer, Chinese speech data in DB-4, and English speech data in DB-4, using complete datasets of 29.77h, 11.84h and 5.68h, respectively.

<sup>6</sup> <https://github.com/svc-develop-team/so-vits-svc>

<sup>7</sup> <https://github.com/YannickJadoul/Parselmouth>

## 4. EXPERIMENTS

### 4.1. Datasets

We use three datasets for experiments, namely, M4Singer [12], NUS-48E [13] and DB-4<sup>1</sup>, as summarized in Table 4. Singing and speech audios are recorded at 44.1kHz, and 48kHz, respectively, with 16-bit quantization, and we downsample them to 24kHz in our experiments.

**Table 4:** Datasets used in our experiments.

Dataset	Language	Voice parts	Duration (hour)	# speakers		# utterances	
				M	F	Singing	Speech
M4Singer [12]	CN	S, A, T, B <sub>1</sub> <sup>8</sup>	29.77	10	10	20942	–
NUS-48E [13]	EN	S, A, T, B <sub>1</sub> , B <sub>2</sub> <sup>8</sup>	1.91	6	6	1262	–
DB-4 <sup>1</sup>	CN	–	11.84	0	1	–	10000
	EN	–	5.68	0	1	–	5000

To be detailed, M4Singer [12] are aligned on the Pinyin-phoneme-level, which cannot capture English pronunciation. An example from M4Singer is provided in Table 5. Phonemes with corresponding musical notes and note durations serve as input to the acoustic model’s front end, carrying phoneme durations for accurate phoneme prediction.

**Table 5:** Phoneme-level annotation in M4Singer.

phs	is_slur	ph_dur	notes	notes_dur
r	0	0.28	56	0.63
ang	0	0.35	56	0.63
uo	0	0.23	56	0.23
n	0	0.32	58	0.51
uan	0	0.19	58	0.51
iang	0	0.0905	61	0.0905
iang	1	0.1636	60	0.1636
iang	0	0.9	58	0.9

### 4.2. Experimental Setup

To verify the proposed methods, we conduct experiments based on the DiffSinger model [10] and settings [12], using the pre-trained HiFi-GAN [29] provided by M4Singer [12] as vocoder. There are two types of model structures and three corpora, and systems to compare are described in Table 6.

- **Model 1:** Original DiffSinger model.
- **Model 2:** DiffSinger with Language-style-infused encoder illustrated in Sec 3.3.
- **Corpus 1:** M4Singer with averaged CMU-phoneme-based annotation described in Sec 3.4.1, totally 29.77h.
- **Corpus 2:** M4Singer with proportional CMU-phoneme-based annotation described in Sec 3.4.1, totally 29.77h.
- **Corpus 3:** As described in Sec 3.4.2, we use so-vits-svc<sup>6</sup> to convert the singing voices in the NUS-48E dataset to the target 20 singers in the M4Singer dataset, thereby each singer owns both Mandarin and English singing data, with a total of  $1.91h \times 20 = 38.2h$ .
- **Corpus 4:** As stated in Sec 3.4.3, we get speech DB-4 and pseudo-singing DB-4,  $11.84h + 5.68h = 17.52h$  each.

**Table 6:** System description.

System	System 1	System 2	System 3	System 4	System 5
<b>Model</b>	Model 1	Model 1	Model 2	Model 2	Model 2
<b>Corpora</b>	Corpus 1	Corpus 2	Corpus 2, 3	Corpus 2, 4	Corpus 2, 3, 4

<sup>8</sup> B<sub>1</sub>, the abbreviation of Bass; B<sub>2</sub>, the abbreviation of Baritone

### 4.3. Evaluation

We create 25 test cases for each language scenario, cumulating a total of 75 cases. To compare the systems<sup>9</sup>, we conduct both objective and subjective evaluations with four target speakers, each of whom represents a different voice part. We recruit a qualified singer to record clean singing voices for each test case, then converted the recording to the four target timbres using so-vits-svc<sup>6</sup> as ground truth.

For the objective evaluation, following [8], we utilize four metrics, including Melcepstrum distortion (MCD), F0 root mean square error (F0\_RMSE), voice/unvoiced error rate (VUV\_E), and semitone accuracy (SA). We employ Whisper [30] to compute the word error rate (WER), which serves as our multilingual evaluation criterion. We use a ResNet101-based speaker verification model, trained with the VoxCeleb2 development set and achieving an Equal Error Rate (EER) of 0.44% on the Vox-O test set, to extract speaker embeddings and calculate cosine similarity (SIM) for the synthesized singing, reaching around 0.6 for the three voice parts: Tenor, Alto, and Soprano, which objectively exhibits the speaker’s identity well hold<sup>10</sup>. For the subjective evaluation, 14 native Chinese speakers with work-proficient English are recruited to score the entire set of generated samples, ranging from 1 to 5 with 0.5 increments. The results are shown in Table 7.

**Table 7:** Experimental results in terms of objective and subjective metrics.

System	Language	MCD ↓	F0_RMSE ↓	VUV_E ↓	SA ↑	SIM ↑	MOS <sup>11</sup> ↑
<b>Ground Truth</b>	CN	–	–	–	–	–	3.78 ± 0.08
	EN	–	–	–	–	–	3.80 ± 0.08
	MIX	–	–	–	–	–	3.60 ± 0.09
	ALL	–	–	–	–	–	3.73 ± 0.05
<b>System 1</b>	CN	<b>9.5983</b>	<b>0.1607</b>	<b>6.53%</b>	<b>45.16%</b>	0.64	<b>3.41 ± 0.09</b>
	EN	10.1743	0.1654	<b>6.36%</b>	39.95%	0.59	2.75 ± 0.09
	MIX	10.8097	0.2312	<b>9.01%</b>	34.67%	0.65	3.06 ± 0.09
	ALL	10.1941	0.1858	<b>7.30%</b>	39.93%	0.63	3.07 ± 0.05
<b>System 2</b>	CN	9.7394	0.1810	7.17%	43.18%	0.62	3.34 ± 0.09
	EN	10.2490	0.1883	8.40%	36.25%	0.57	2.56 ± 0.09
	MIX	11.0579	0.2293	9.93%	33.51%	0.62	2.90 ± 0.09
	ALL	10.3488	0.1995	8.50%	37.64%	0.60	2.96 ± 0.05
<b>System 3</b>	CN	9.7953	0.1709	7.29%	43.38%	0.58	2.97 ± 0.09
	EN	<b>8.6386</b>	<b>0.1440</b>	6.93%	<b>45.05%</b>	0.52	3.40 ± 0.08
	MIX	10.3429	<b>0.2163</b>	9.69%	<b>38.88%</b>	0.58	<b>3.32 ± 0.08</b>
	ALL	<b>9.5923</b>	<b>0.1770</b>	7.97%	<b>42.44%</b>	0.56	3.23 ± 0.05
<b>System 4</b>	CN	9.8097	0.1711	7.00%	44.14%	0.60	3.14 ± 0.09
	EN	10.4531	0.1854	7.69%	33.11%	0.49	2.58 ± 0.09
	MIX	11.3648	0.2262	10.54%	31.95%	0.58	2.89 ± 0.08
	ALL	10.5425	0.1942	8.41%	36.40%	0.56	2.89 ± 0.05
<b>System 5</b>	CN	9.8542	0.1761	7.15%	44.28%	0.57	3.04 ± 0.09
	EN	8.7834	0.1541	7.68%	41.93%	0.49	<b>3.42 ± 0.09</b>
	MIX	<b>10.2845</b>	0.2182	10.14%	37.21%	0.59	3.24 ± 0.09
	ALL	9.6407	0.1828	8.33%	41.14%	0.55	<b>3.24 ± 0.05</b>

Based on the WER (Fig. 3) and MOS results, it is evident that **System 1** and **System 2** perform poorly on English songs. With the integration of Corpus 3, **System 3** and **System 5** have shown significant improvement, with MOS scores of 3.40 and 3.42, respectively, compared to **System 2**’s score of 2.56. Comparing the results from **System 3, 4, 5**, using Corpus 4 has effectively alleviated the decline in the Chinese synthesis effect. At the same time, the real English singing data in Corpus 3 has assisted the model in effectively learning English

<sup>9</sup> Please note that when encountering missing phonemes, we replace them with phonetically similar ones for **System 1** and **System 2**.

<sup>10</sup>The highly-precise verification model results in a relatively strict score.

<sup>11</sup>With 95% confidence interval.

from Corpus 4, consisting of speech data and pseudo-singing data.

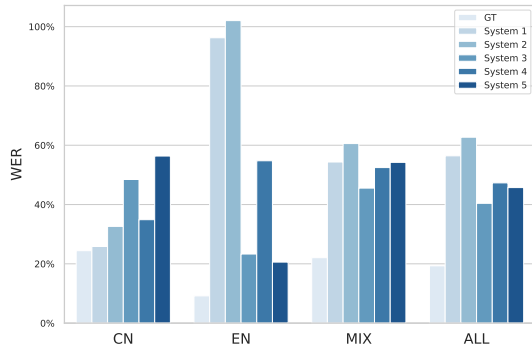


Fig. 3: WER results.

To compare these systems more effectively, we conducted a preference test, asking the subjects to choose the best system in pronunciation (Fig. 4) for Chinese and English songs. We found that **System 1** performed best in Chinese; however, its understanding of Chinese phoneme duration proportions could not be transferred to English synthesis. With the proportionally distributed phoneme duration, the synthesized effect of Chinese is still maintained, as evidenced by the results of **System 2, 3, 4, 5**. Furthermore, the proportionally obtained duration proved more suitable for English tasks, compared to **System 1**. It is worth noting that even though **System 2**, trained exclusively with Chinese data, could generate English singing voices pleasing to the human ear, the large model, Whisper [30], classified such pronunciation effects as Chinese based on its transcription results and WER calculations. However, after adding Corpus 3 and Corpus 4, the English WER results of **System 3, 4, 5** have dropped significantly, indicating that real English data is necessary to learn English phonemes’ pronunciation better. Particularly, according to preference results, **System 3** and **System 5** have gained substantial benefits from incorporating real English singing data.

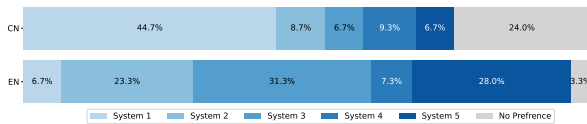


Fig. 4: Preference results in terms of pronunciation.

#### 4.4. Discussion on Phoneme Substitution

Based on the results from **System 1** and **System 2**, it is evident that shared CMU phonemes can partially assist our model in learning English pronunciation from Chinese data. Most of these phonemes are present in both languages. However, certain English phonemes like /TH/, /Y/, /IH/, /DH/, /V/, and /OY/ are absent in Chinese data. During the inference process, these are substituted with similarly pronounced Chinese phonemes. However, this approach may only sometimes work since some phonemes lack similar pronunciations in Chinese. In these situations, the pronunciation of each phoneme remains isolated and cannot be connected coherently as a word. Frequent substitutions may also result in synthesized voices that sound

unclear and similar to Chinese. Consider the lyrics “I’m in love with the shape of you” as an example (Fig. 5); the system with phoneme substitution performs worse than other systems. The corresponding audio example is also available online<sup>12</sup>.

Expected: AY M IH N L AH V W IH DH TH SH EY P AH V Y UW  
Substituted: AY M AY N L AH W W AY Z S SH EY P AH W IY UW

Fig. 5: Substitution example.

#### 4.5. Ablation Study

We further conduct ablation studies to confirm the effectiveness of the methods proposed in our BiSinger model, which include: 1) using a timbre conversion approach for the monolingual singing dataset, 2) applying pitch shift to the speech dataset, and 3) introducing the Language-style-infused encoder. For clarity, we undertook an additional MOS assessment, where we randomly selected and evaluated 18 distinct test cases based on a single target singer. The results, shown in Table 8, confirm the effectiveness of our outlined strategy for multilingual SVS. When compared to **ASM 1** and **ASM 2**, **ASM 3** demonstrates that timbre conversion and the introduction of language ID tokens can improve performance. Moreover, the DB-4 dataset exhibits better performance with the application of pitch shift as opposed to without it.

Table 8: MOS scores for ablation study.

System	SVC	LGT	NUS	DB-4	pitch-shift	MOS <sup>11</sup> ↑
<b>ASM 1</b>	-	-	✓	-	-	3.83 ± 0.10
<b>ASM 2</b>	✓	-	✓	-	-	3.87 ± 0.10
<b>ASM 3</b>	✓	✓	✓	-	-	3.96 ± 0.09
<b>ASM 4</b>	-	-	-	✓	-	3.83 ± 0.10
<b>ASM 5</b>	-	✓	-	✓	-	3.78 ± 0.10
<b>ASM 6</b>	-	✓	-	✓	✓	3.87 ± 0.10

## 5. CONCLUSION

In this paper, we have developed BiSinger, a system that can synthesize singing voices in both Chinese and English. To overcome language barriers, we adopted a language-independent representation, transitioning from a Pinyin-based annotation to a CMU-based one. Furthermore, we have enhanced the model’s performance in English and code-switch SVS by incorporating a monolingual singing dataset using advanced SVC techniques. In addition, we have explored the use of a bilingual speech dataset to facilitate multilingual SVS. Experimental results demonstrate that our methods can generate multilingual singing voices and enhance English and code-switch SVS while maintaining performance in Chinese songs. In our future work, we will continue to study multilingual and code-switched singing voice synthesis. And it would be promising work to find other solid evaluation metrics for singing voices.

#### Acknowledgement

This research is funded by the Kunshan Municipal Government Research Funding under the project “Deep Learning based Singing Voice Synthesis for Kun Opera”.

<sup>12</sup><https://bisinger-svs.github.io>

## 6. REFERENCES

- [1] P. Lu, J. Wu, J. Luan, X. Tan, and L. Zhou, “Xiaoice-singer: A high-quality and integrated singing voice synthesis system,” *arXiv preprint arXiv:2006.06261*, 2020.
- [2] Y. Ren, Y. Ruan, X. Tan, T. Qin, S. Zhao, Z. Zhao, and T.-Y. Liu, “Fastspeech: Fast, robust and controllable text to speech,” *Advances in neural information processing systems*, vol. 32, 2019.
- [3] Y. Gu, X. Yin, Y. Rao, Y. Wan, B. Tang, Y. Zhang, J. Chen, Y. Wang, and Z. Ma, “Bytesinger: A chinese singing voice synthesis system using duration allocated encoder-decoder acoustic models and wavernn vocoders,” in *2021 12th International Symposium on Chinese Spoken Language Processing (ISCSLP)*. IEEE, 2021, pp. 1–5.
- [4] Y. Wang, R. Skerry-Ryan, D. Stanton, Y. Wu, R. J. Weiss, N. Jaitly, Z. Yang, Y. Xiao, Z. Chen, S. Bengio, et al., “Tacotron: Towards end-to-end speech synthesis,” *arXiv preprint arXiv:1703.10135*, 2017.
- [5] J. Kim, J. Kong, and J. Son, “Conditional variational autoencoder with adversarial learning for end-to-end text-to-speech,” in *International Conference on Machine Learning*. PMLR, 2021, pp. 5530–5540.
- [6] Y. Zhang, J. Cong, H. Xue, L. Xie, P. Zhu, and M. Bi, “Visinger: Variational inference with adversarial learning for end-to-end singing voice synthesis,” in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022, pp. 7237–7241.
- [7] K. Oura, A. Mase, T. Yamada, S. Muto, Y. Nankaku, and K. Tokuda, “Recent development of the hmm-based singing voice synthesis system—sinsy,” in *Seventh ISCA Workshop on Speech Synthesis*, 2010.
- [8] J. Shi, S. Guo, T. Qian, N. Huo, T. Hayashi, Y. Wu, F. Xu, X. Chang, H. Li, P. Wu, et al., “Muskits: an end-to-end music processing toolkit for singing voice synthesis,” *arXiv preprint arXiv:2205.04029*, 2022.
- [9] R. Yamamoto, R. Yoneyama, and T. Toda, “Nnsvs: A neural network-based singing voice synthesis toolkit,” in *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2023, pp. 1–5.
- [10] J. Liu, C. Li, Y. Ren, F. Chen, and Z. Zhao, “Diff-singer: Singing voice synthesis via shallow diffusion mechanism,” in *Proceedings of the AAAI conference on artificial intelligence*, 2022, vol. 36, pp. 11020–11028.
- [11] S. Choi, W. Kim, S. Park, S. Yong, and J. Nam, “Children’s song dataset for singing voice research,” in *International Society for Music Information Retrieval Conference (ISMIR)*, 2020.
- [12] L. Zhang, R. Li, S. Wang, L. Deng, J. Liu, Y. Ren, J. He, R. Huang, J. Zhu, X. Chen, et al., “M4singer: A multi-style, multi-singer and musical score provided mandarin singing corpus,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 6914–6926, 2022.
- [13] Z. Duan, H. Fang, B. Li, K. C. Sim, and Y. Wang, “The nus sung and spoken lyrics corpus: A quantitative comparison of singing and speech,” in *2013 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference*. IEEE, 2013, pp. 1–9.
- [14] Z. Cai, Y. Yang, and M. Li, “Cross-lingual multi-speaker speech synthesis with limited bilingual training data,” *Computer Speech & Language*, vol. 77, pp. 101427, 2023.
- [15] H. Kenmochi and H. Ohshita, “Vocaloid-commercial singing synthesizer based on sample concatenation,” in *Interspeech*, 2007, vol. 2007, pp. 4009–4010.
- [16] J. Bonada, M. Umberto Morist, and M. Blaauw, “Expressive singing synthesis based on unit selection for the singing synthesis challenge 2016,” *Morgan N, editor. Interspeech 2016; 2016 Sep 8-12; San Francisco, CA.[place unknown]: ISCA; 2016. p. 1230-4.*, 2016.
- [17] K. Nakamura, K. Oura, Y. Nankaku, and K. Tokuda, “Hmm-based singing voice synthesis and its application to japanese and english,” in *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2014, pp. 265–269.
- [18] Y. Ren, X. Tan, T. Qin, J. Luan, Z. Zhao, and T.-Y. Liu, “Deepsinger: Singing voice synthesis with data mined from the web,” in *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2020, pp. 1979–1989.
- [19] Y. Zhang, R. J. Weiss, H. Zen, Y. Wu, Z. Chen, R. Skerry-Ryan, Y. Jia, A. Rosenberg, and B. Ramabhadran, “Learning to speak fluently in a foreign language: Multilingual speech synthesis and cross-language voice cloning,” *arXiv preprint arXiv:1907.04448*, 2019.
- [20] B. Li, Y. Zhang, T. Sainath, Y. Wu, and W. Chan, “Bytes are all you need: End-to-end multilingual speech recognition and synthesis with bytes,” in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 5621–5625.
- [21] Y. Cao, X. Wu, S. Liu, J. Yu, X. Li, Z. Wu, X. Liu, and H. Meng, “End-to-end code-switched tts with mix of monolingual recordings,” in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 6935–6939.
- [22] T. Nekvinda and O. Dušek, “One model, many languages: Meta-learning for multilingual text-to-speech,” *arXiv preprint arXiv:2008.00768*, 2020.

- [23] Z. Liu and B. Mak, “Cross-lingual multi-speaker text-to-speech synthesis for voice cloning without using parallel corpus for unseen speakers,” *arXiv preprint arXiv:1911.11601*, 2019.
- [24] Z. Liu and B. Mak, “Multi-lingual multi-speaker text-to-speech synthesis for voice cloning with online speaker enrollment.,” in *Interspeech*, 2020, pp. 2932–2936.
- [25] Y. Lee, S. Shon, and T. Kim, “Learning pronunciation from a foreign language in speech synthesis networks,” *arXiv preprint arXiv:1811.09364*, 2018.
- [26] F. Yang, J. Luan, and Y. Wang, “Improve bilingual tts using dynamic language and phonology embedding,” *arXiv preprint arXiv:2212.03435*, 2022.
- [27] M. McAuliffe, M. Socolof, S. Mihuc, M. Wagner, and M. Sonderegger, “Montreal forced aligner: Trainable text-speech alignment using kaldi.,” in *Interspeech*, 2017, vol. 2017, pp. 498–502.
- [28] 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.
- [29] J. Kong, J. Kim, and J. Bae, “Hifi-gan: Generative adversarial networks for efficient and high fidelity speech synthesis,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 17022–17033, 2020.
- [30] A. Radford, J. W. Kim, T. Xu, G. Brockman, C. McLevey, and I. Sutskever, “Robust speech recognition via large-scale weak supervision,” in *International Conference on Machine Learning*. PMLR, 2023, pp. 28492–28518.