

*The 32nd International Joint Conference on
Artificial Intelligence (IJCAI 2023)*



IJCAI/2023 MACAO

Learning to Speak from Text: Zero-Shot Multilingual Text-to-Speech with Unsupervised Text Pretraining

Takaaki Saeki¹, Soumi Maiti², Xinjian Li², Shinji Watanabe²,
Shinnosuke Takamichi¹, Hiroshi Saruwatari¹



¹The University of Tokyo, Japan
²Carnegie Mellon University, USA



TTS is widely used in AI voice user interfaces.

Recent neural TTS [Kim+21] achieves **human-like natural speech.**



Only open to limited number of resource-rich languages.

- Requiring **paired speech-text data** with studio recording audio.
- **Hard to collect enough data** for low-resource languages.



Low-resource TTS approaches to expand number of languages.

- Adapting multilingual model to low-resource language [Lee+18] [He+21]
- Using untranscribed speech for training [Zhang+20] [Ni+21]
- Joint semi-supervised learning with different types of data [Saeki+23]

Previous approaches heavily rely on speech recordings.

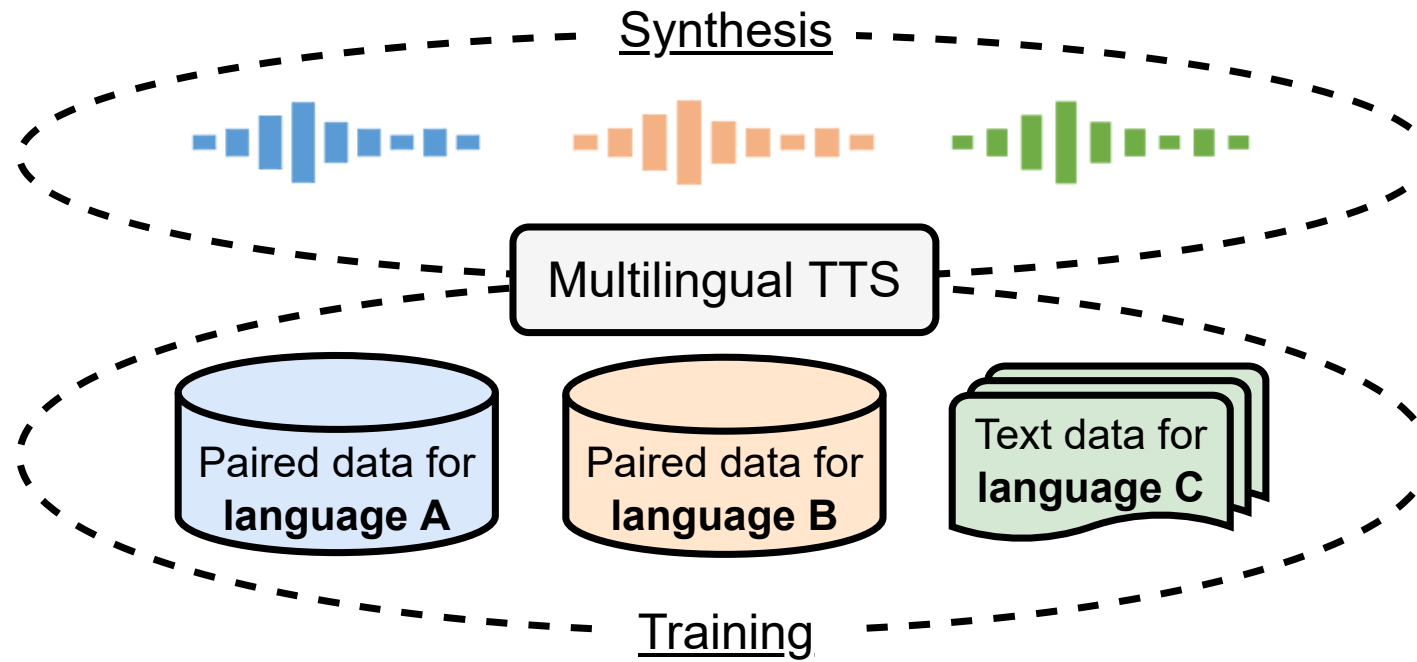
➤ **Often challenging to collect training data for target languages.**

Text data is much easier to collect than paired speech-text data.

- **No need to collect studio recording audio.**
- **No need of preprocessing** to align speech and text.
- **Free from sensitive speaker-related information.**

Goal: Building TTS for languages **with only textual resources.**

Open-up TTS systems to much more languages.



Strong zero-shot cross-lingual transferability of multilingual BERT [Devlin+19] in natural language processing tasks [Pires+19].

We investigate **cross-lingual transfer of multilingual LM for TTS**.

Strong zero-shot cross-lingual transferability of multilingual BERT [Devlin+19] in natural language processing tasks [Pires+19].

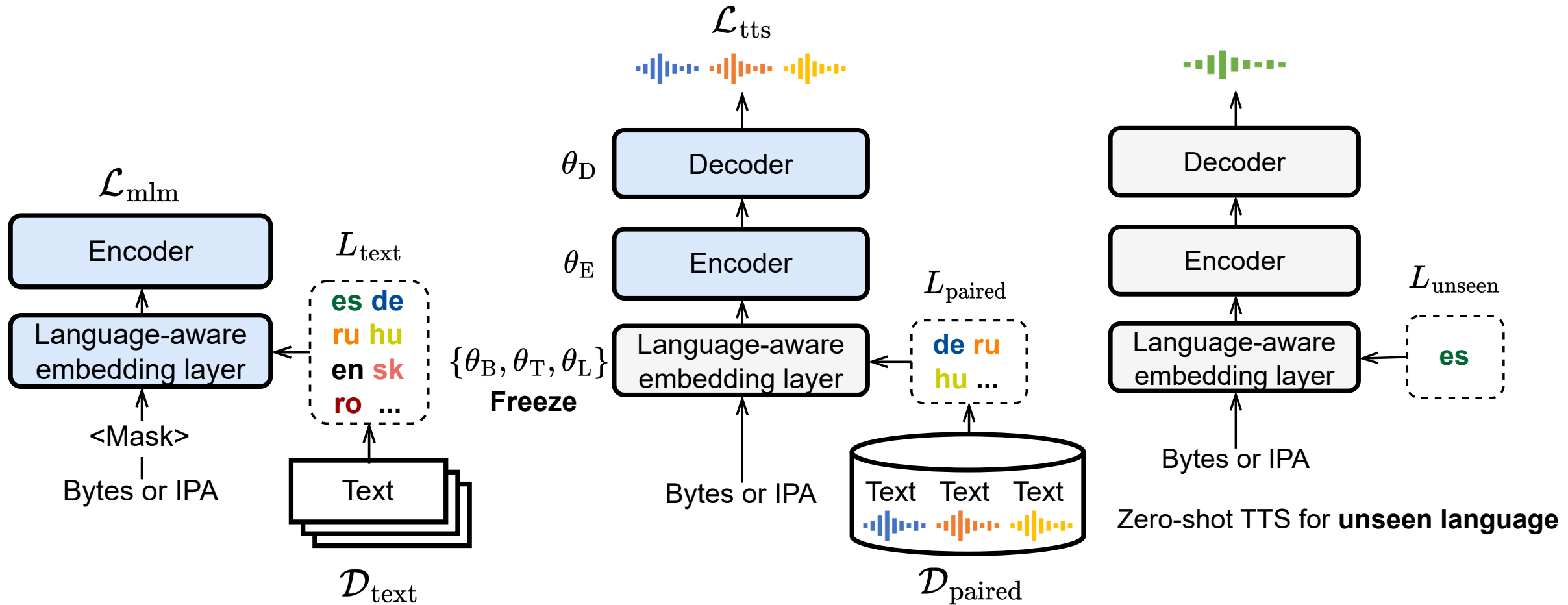
We investigate **cross-lingual transfer of multilingual LM for TTS**.

Contributions:

- ❑ **Propose zero-shot TTS from text data, achieving high intelligibility.**
- ❑ **Improve multilingual TTS without per-language pronunciation knowledge.**
- ❑ **Conducted comprehensive ablation studies.**

- Background
- Method
- Experiment
- Conclusions

Transformer-based multilingual TTS model using text-only and paired data.



(a) Unsupervised multilingual text pretraining

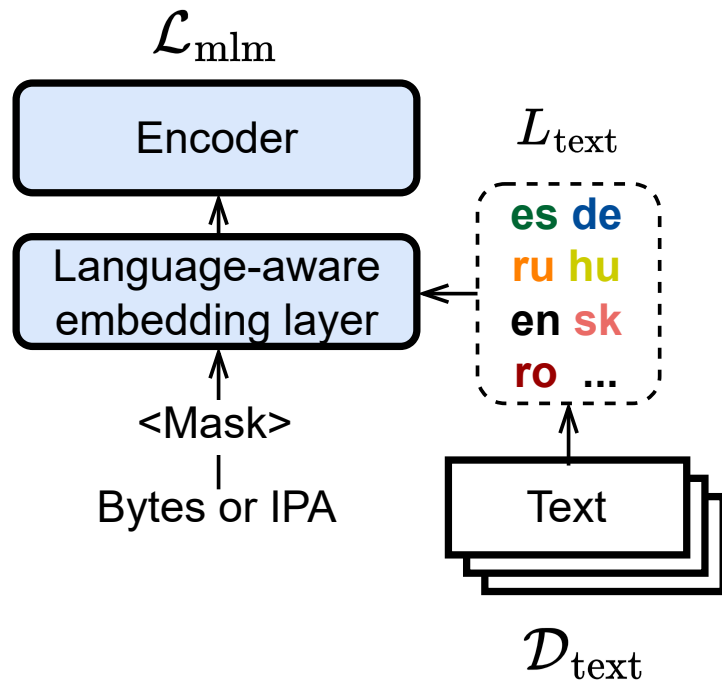
(b) Supervised learning with paired data

(c) Inference

MLM pretraining with multilingual text including many languages

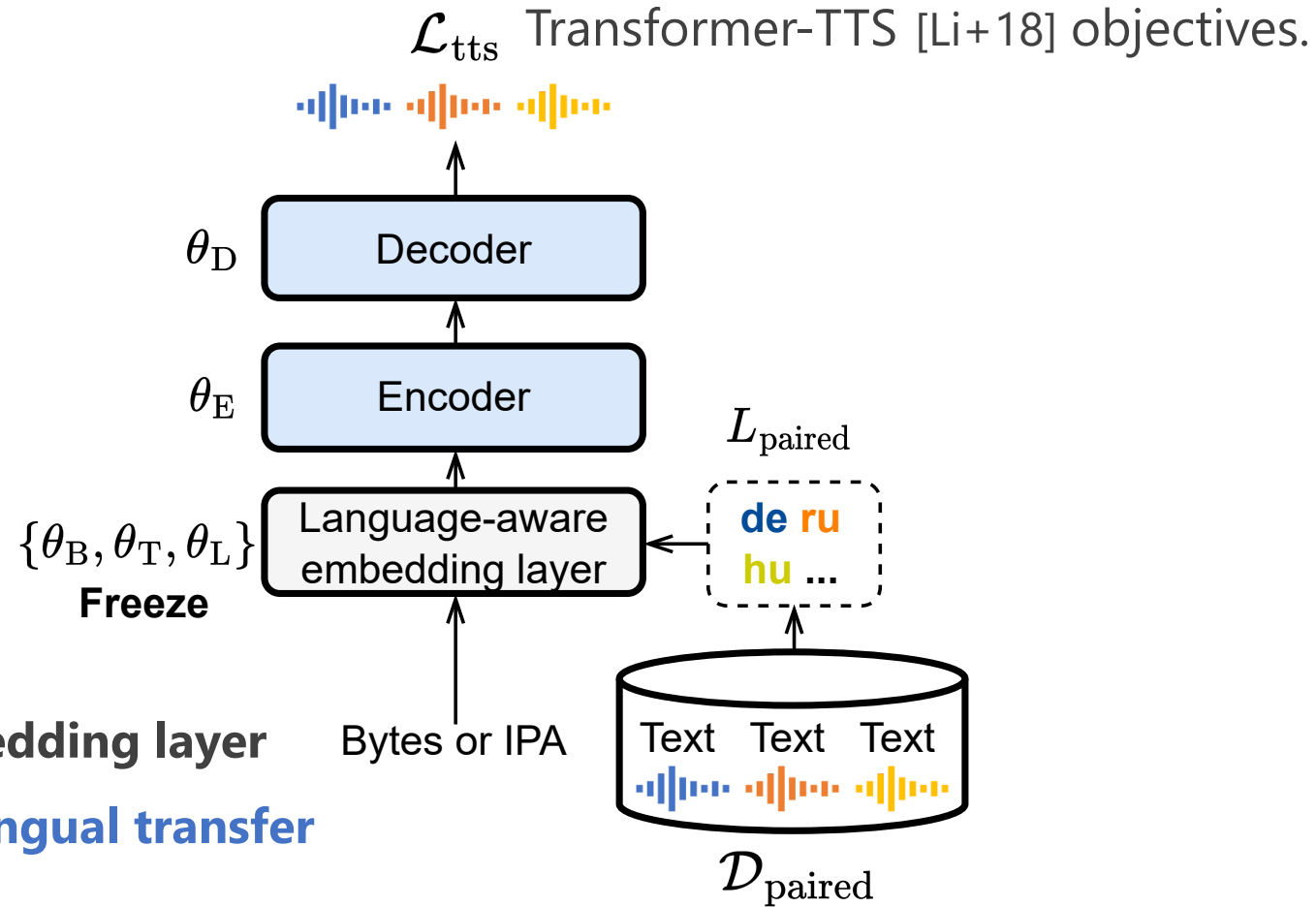
- ❑ Language-agnostic tokens (Bytes and IPA*) for TTS
- ❑ Language-aware embedding layer to inject language info.

*IPA: Using International Phonetic Alphabet Symbols



(a) Unsupervised multilingual text pretraining

Supervised learning with paired data including a few languages.



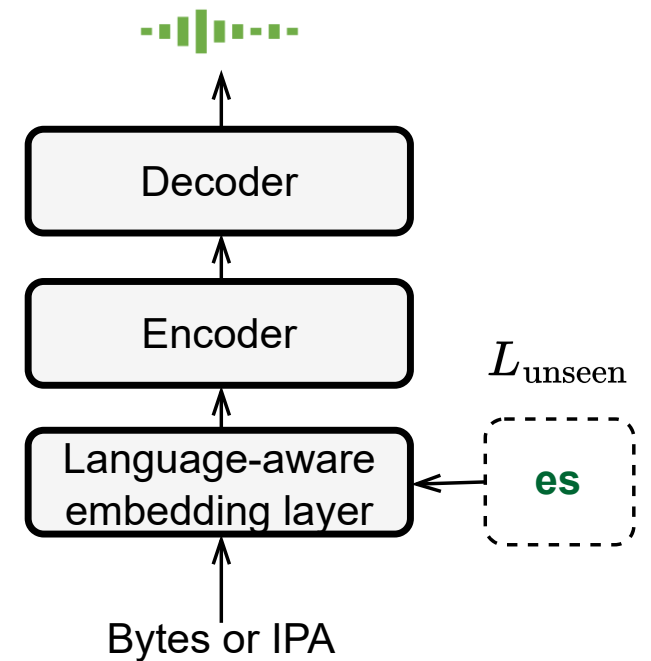
Frozen language-aware embedding layer

Facilitating zero-shot cross-lingual transfer

(b) Supervised learning with paired data

Zero-shot TTS by using *unseen* language IDs

Using language IDs only included in text data, not in paired data.



Zero-shot TTS for **unseen language**

(c) Inference

- Background
- Related Work
- Method
- Experiment
- Conclusions

Text data: 19 European languages.

Paired data: 7 European languages.

Languages	Code	Text-only data	Paired data	
			Text	Audio
<i>Seen languages for evaluation L_{seen}</i>				
German	de	359MB	0.73MB	16.13h
French	fr	372MB	0.94MB	19.15h
Dutch	nl	336MB	0.75MB	14.10h
Finnish	fi	308MB	0.47MB	21.36h
Hungarian	hu	104MB	0.51MB	10.53h
Russian	ru	4.9MB	1.5MB	10.00h
Greek	el	0.39MB	0.39MB	4.13h
<i>Unseen language for evaluation L_{unseen}</i>				
Spanish	es	345MB	0.0MB (1.2MB)	0.00h (23.81h)

Chose **Spanish** as an unseen language for main evaluation.

* To ensure enough human evaluators.

Languages	Code	Text-only data	Paired data	
			Text	Audio
<i>Seen languages for evaluation L_{seen}</i>				
German	de	359MB	0.73MB	16.13h
French	fr	372MB	0.94MB	19.15h
Dutch	nl	336MB	0.75MB	14.10h
Finnish	fi	308MB	0.47MB	21.36h
Hungarian	hu	104MB	0.51MB	10.53h
Russian	ru	4.9MB	1.5MB	10.00h
Greek	el	0.39MB	0.39MB	4.13h
<i>Unseen language for evaluation L_{unseen}</i>				
Spanish	es	345MB	0.0MB (1.2MB)	0.00h (23.81h)

Methods

- ❑ **Baseline:** **without** unsupervised text pretraining.
- ❑ **Proposed:** **with** unsupervised text pretraining.
- ❑ **Oracle:** **Using paired data** for the target language.

Token types

- ❑ **Bytes:** Without language-specific knowledge.
- ❑ **IPA:** Using per-language pronunciation dictionary.

Evaluation metrics

- Mel Cepstral Distortion (**MCD**) [Fukada+1992]
- Character Error Rates (**CER**) computed by Whisper [Radford+22]
- Automatic Mean opinion scores (**AMOS**)
- Subjective mean opinion scores (**MOS**)

Results (Unseen Language)

19/23

Oracle: Using unseen language during training

	Spanish	
	MCD (↓)	CER (↓)
<i>Ground-truth</i>	-	2.71
<i>Oracle</i> (Bytes)	8.65	10.70
<i>Oracle</i> (IPA)	6.20	5.32
<i>Baseline</i> (IPA)	10.75	44.75
<i>Proposed</i> (Bytes)	9.05	18.27
<i>Proposed</i> (IPA)	9.44	11.69

Results (Unseen Language)

20/23

Oracle: Using unseen language during training

	Spanish	
	MCD (↓)	CER (↓)
<i>Ground-truth</i>	-	2.71
<i>Oracle</i> (Bytes)	8.65	10.70
<i>Oracle</i> (IPA)	6.20	5.32
<i>Baseline</i> (IPA)	10.75	44.75
<i>Proposed</i> (Bytes)	9.05	18.27
<i>Proposed</i> (IPA)	9.44	11.69

Proposed achieved intelligible (< 20% CER) zero-shot TTS compared with *Baseline*

Results (Unseen Language)





21/23

Oracle: Using unseen language during training

	Spanish	
	MCD (↓)	CER (↓)
<i>Ground-truth</i>	-	2.71
<i>Oracle</i> (Bytes)	8.65	10.70
<i>Oracle</i> (IPA)	6.20	5.32
<i>Baseline</i> (IPA)	10.75	44.75
<i>Proposed</i> (Bytes)	9.05	18.27
<i>Proposed</i> (IPA)	9.44	11.69

Proposed was still worse than *Oracle* (Phones) but comparable to *Oracle* (Bytes)

Se me representaba el sonido de las campanas de la iglesia, tocadas por los cuatro muchachos o por el ingrato padre.

	Spanish	
	Sample	CER (↓)
<i>Ground-truth</i>		2.71
<i>Oracle</i> (IPA)		5.32
<i>Baseline</i> (Bytes)		66.45
<i>Proposed</i> (Bytes)		18.27

Background

- ❑ Need to reduce **cost of data collection** of neural multilingual TTS

Method

- ❑ **Multilingual unsupervised text pretraining**
- ❑ Zero-shot TTS from unseen language

Results

- ❑ Achieved highly **intelligible** (CER < 12%) **zero-shot TTS**
- ❑ Observed **language dependency**

Future work

- ❑ Need to improve **naturalness and prosody**
- ❑ Need to develop a method that **works well for many languages**

Paper



More audio samples

