

IEEE ICASSP 2024 @ Seoul, Korea

# Extending Multilingual Speech Synthesis to 100+ Languages without Transcribed Data

Takaaki Saeki<sup>1, 2</sup>, Gary Wang<sup>1</sup>, Nobuyuki Morioka<sup>1</sup>, Isaac Elias<sup>1</sup>,  
Kyle Kastner<sup>1</sup>, Andrew Rosenberg<sup>1</sup>, Bhuvana Ramabhadran<sup>1</sup>,  
Heiga Zen<sup>1</sup>, Françoise Beaufays<sup>1</sup>, Hadar Shemtov<sup>1</sup>

<sup>1</sup>Google, <sup>2</sup>The University of Tokyo, Japan

# Overview of This Talk

**Single multilingual text-to-speech (TTS) model on 100+ languages**

Built on paired/unpaired **found** data, **without studio-quality paired data**

**Zero-supervised TTS** - w/ untranscribed found data

**20 out of 50 langs within 5% CER diff.** to ground truth

**Minimally supervised found TTS** - w/ 15min paired FLEURS

**42 out of 50 langs within 5% CER diff.** to ground truth

# Outline

**1. Background**

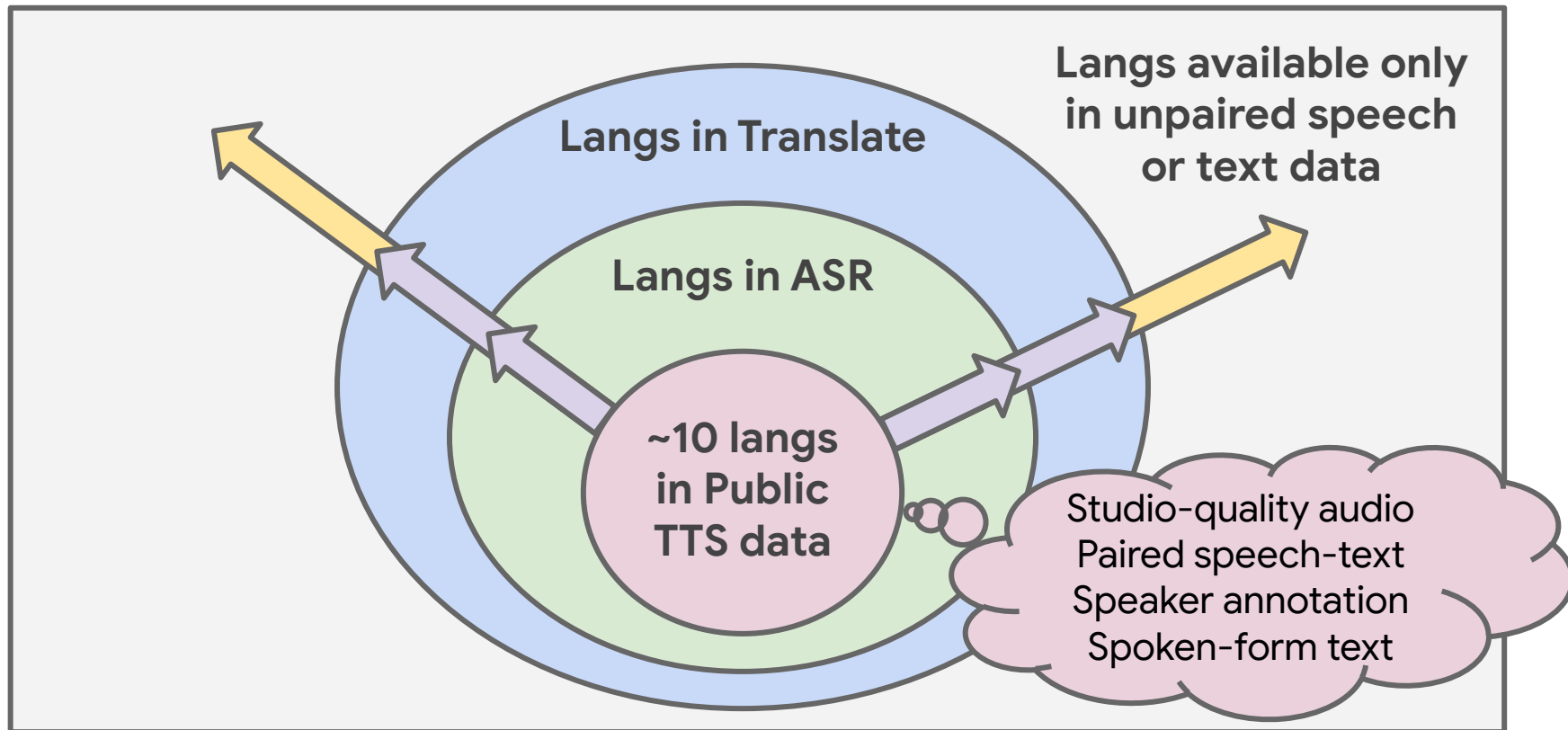
2. Method

3. Experimental Settings

4. Results

5. Summary & Future Work

# Language Extension of Speech Synthesis



# Using Untranscribed Found Data to Build TTS

Higher potential for language extension

## Typical Paired TTS data

(e.g., LJSpeech [Ito+17])

Studio recordings (MOS ~ 4.5)

Transcribed

Speaker information

Short sentences suitable for TTS

High-cost for collection

## “Found” (aka ASR) data

(e.g., FLEURS [Conneau+22])

Real-world noisy audio (MOS ~ 3.5)

Often untranscribed

No speaker information

Longer utterances or spontaneous

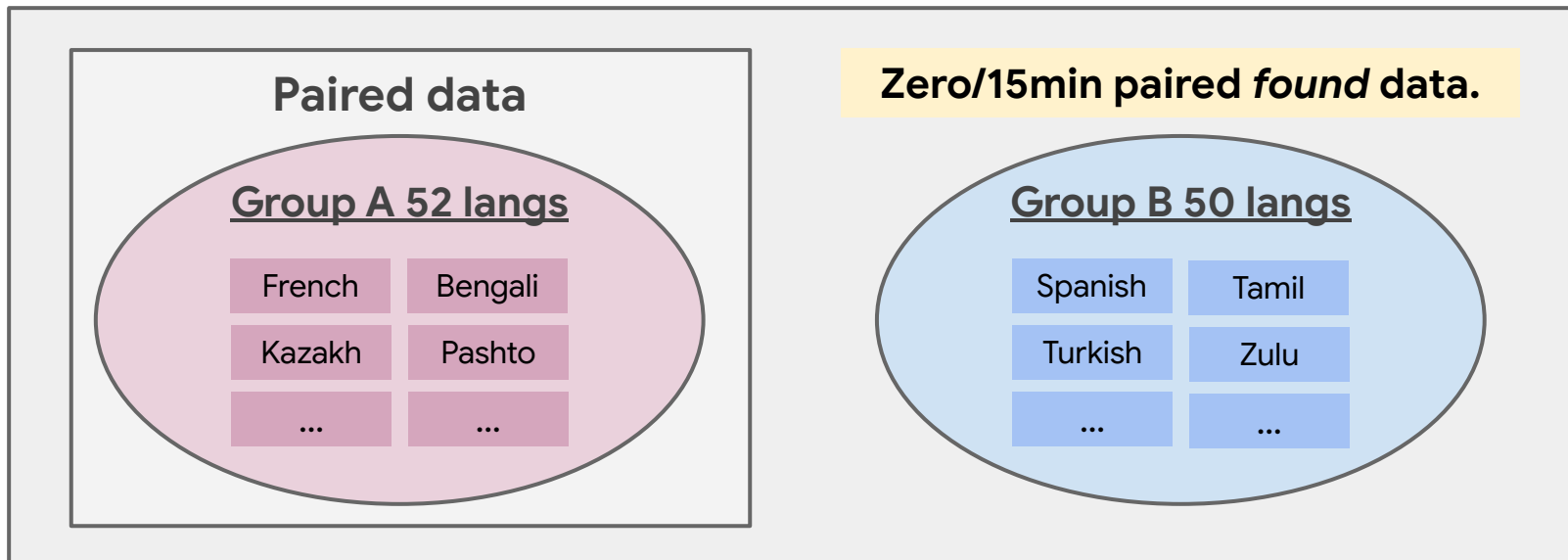
Easier to collect

# Our Concept: TTS on Paired/Unpaired Found Data

Multilingual joint semi-supervised learning on **Group A/B languages** [Chen+23].

Can we build TTS on **Group B languages** by leveraging untranscribed data?

Single TTS model on 100+ languages



# Main Contributions

Scaling a single TTS model to **100+ languages with multiple language families and writing systems.**

Showing improved capability of **zero/minimally supervised TTS.**

Robust TTS model architectures on supervised/unsupervised found data **without studio-quality paired TTS data.**

# Related Work

## **Scaling Speech Technology to 1,000+ Languages** [Pratap+23]

Building **monolingual supervised TTS** for each of 1107 languages.

**Unsupervised TTS:** w/ Unsupervised ASR [Ni+22], w/ text data [Saeki+23].

Investigated a few language families and writing systems.

**Virtuoso** (ICASSP'23): Joint multilingual semi-supervised learning for TTS

Still used full supervised ASR data to build multilingual TTS



# Outline

1. Background

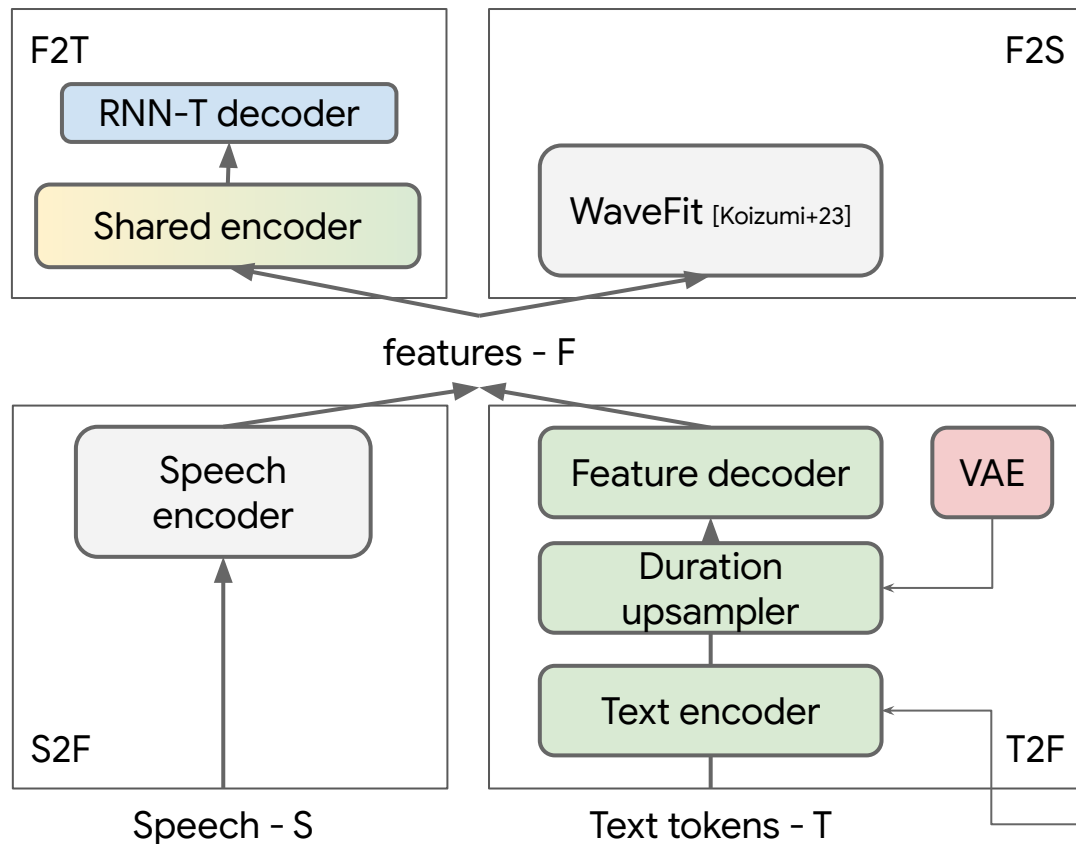
**2. Method**

3. Experimental Settings

4. Results

5. Summary & Future Work

# Overview of Model Architectures



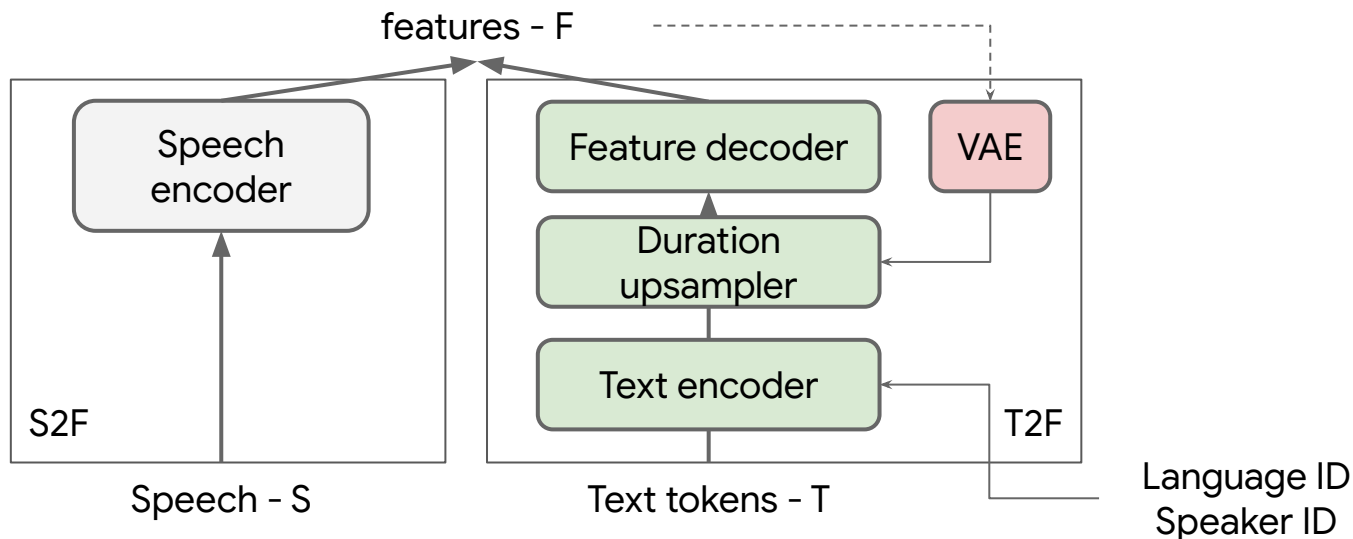
## Four components:

- \* Text-to-feature (T2F)
- \* Speech-to-feature (S2F)
- \* Feature-to-text (F2T)
- \* Feature-to-speech (F2S)

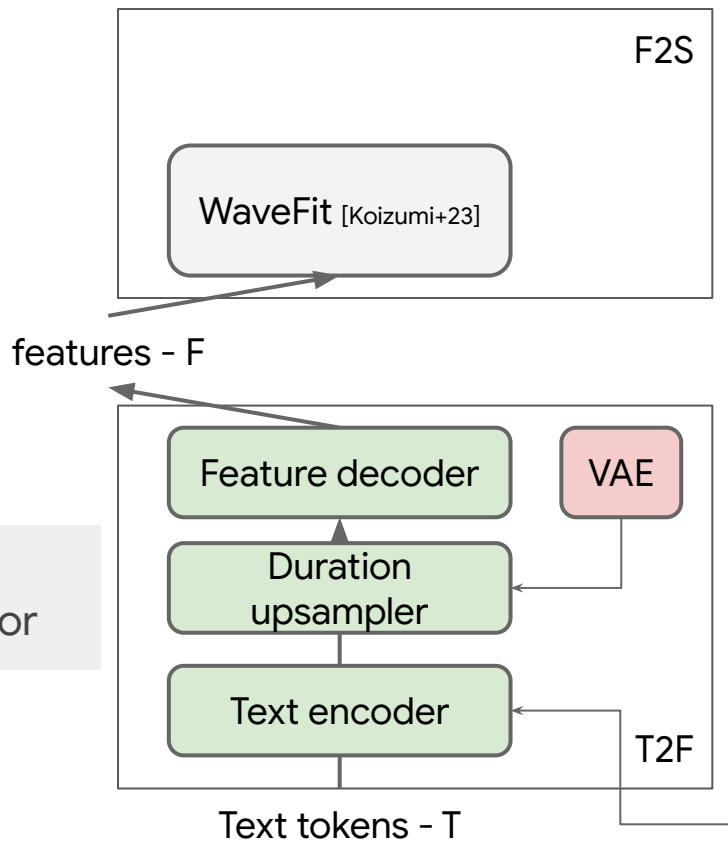
Language ID  
Speaker ID

# Text-to-Feature (T2F) Module

- Non-autoregressive model with a duration-based upsampler [Elias+20]
- **Self-supervised features (F)** from speech encoder
- **Token-level variational autoencoder (VAE)** to capture variability



# Inference Procedure



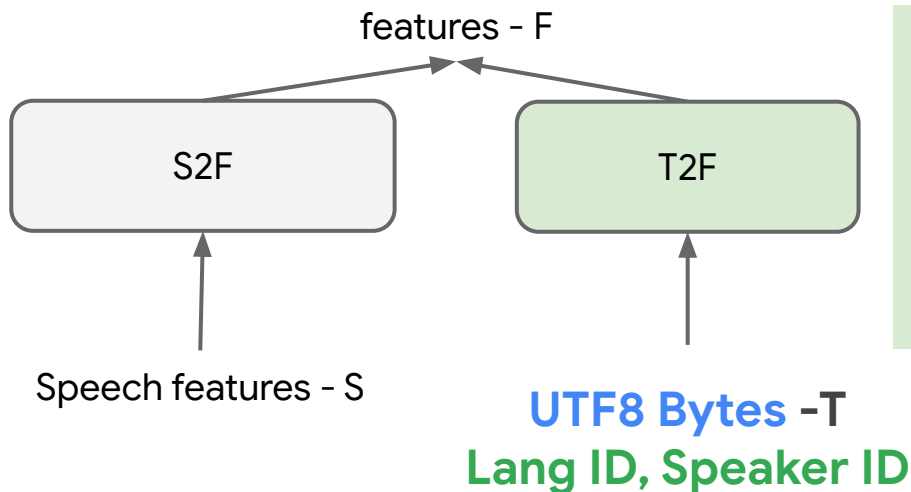
Typical cascaded non-autoregressive TTS model with T2F and F2S.

NOTE: Features are derived from self-supervised learning, instead of mel spectrograms.

- \* Predicted durations
- \* Latents sampled from prior

Language ID  
Speaker ID

# Input Representations for Cross-Lingual Transfer

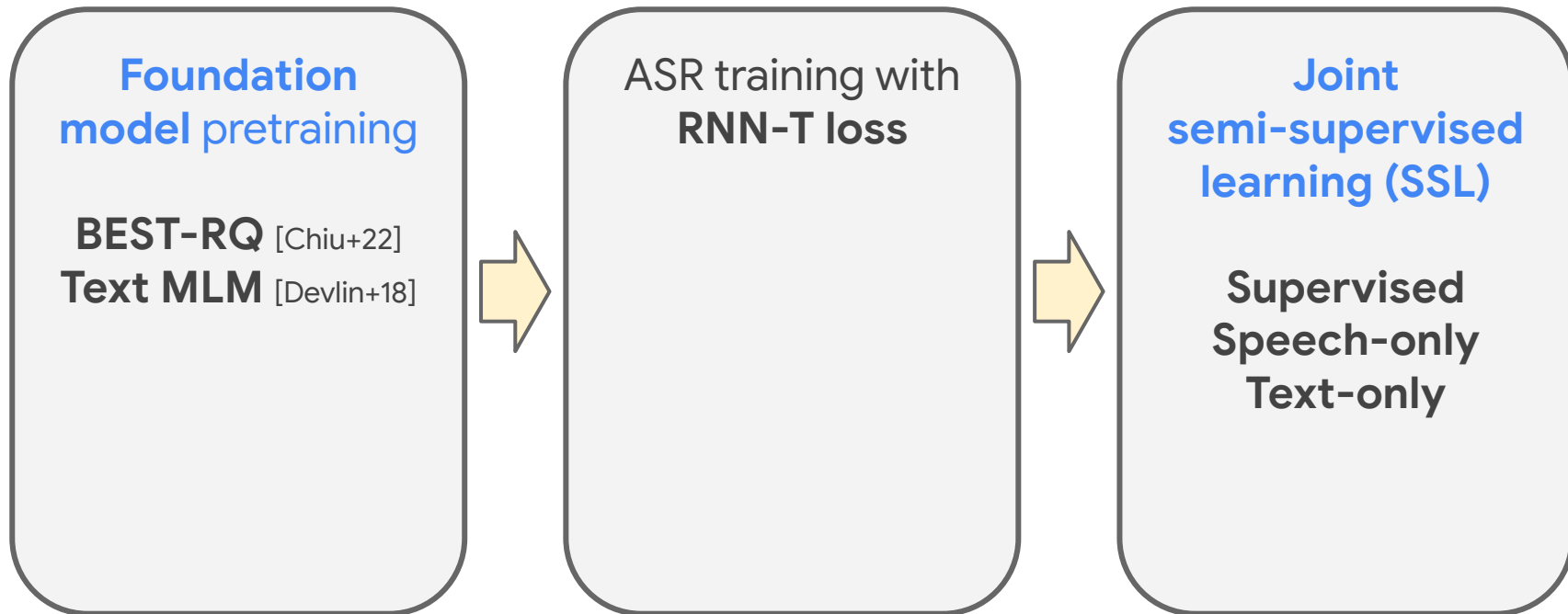


\* **UTF-8 Byte tokens** as input representation, enabling cross script transfer & robustness to *unseen* graphemes [He+21]

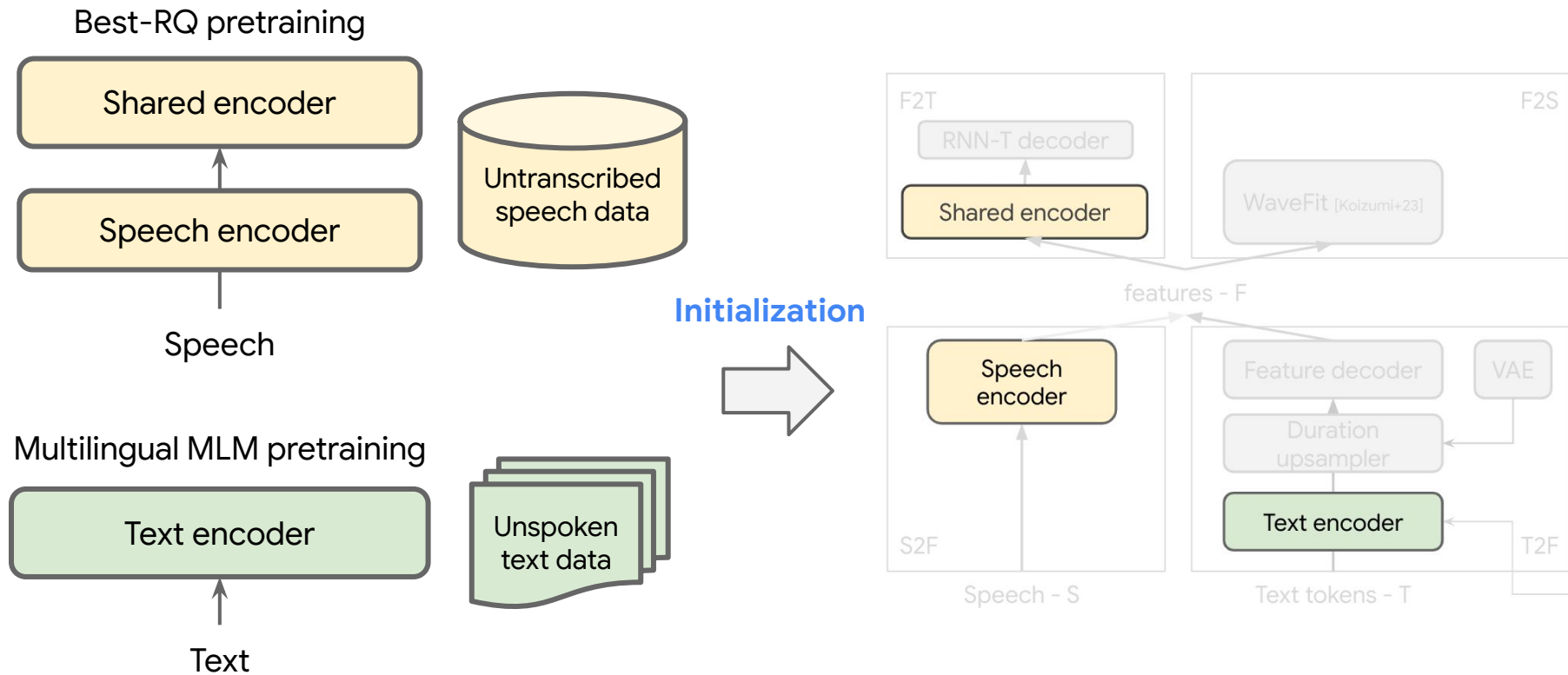
\* **Classifier Free Guidance** [Ho+22]

Randomly replace 10% IDs with `<unk_id>`, promoting cross speaker/language transfer by breaking conditional dependence

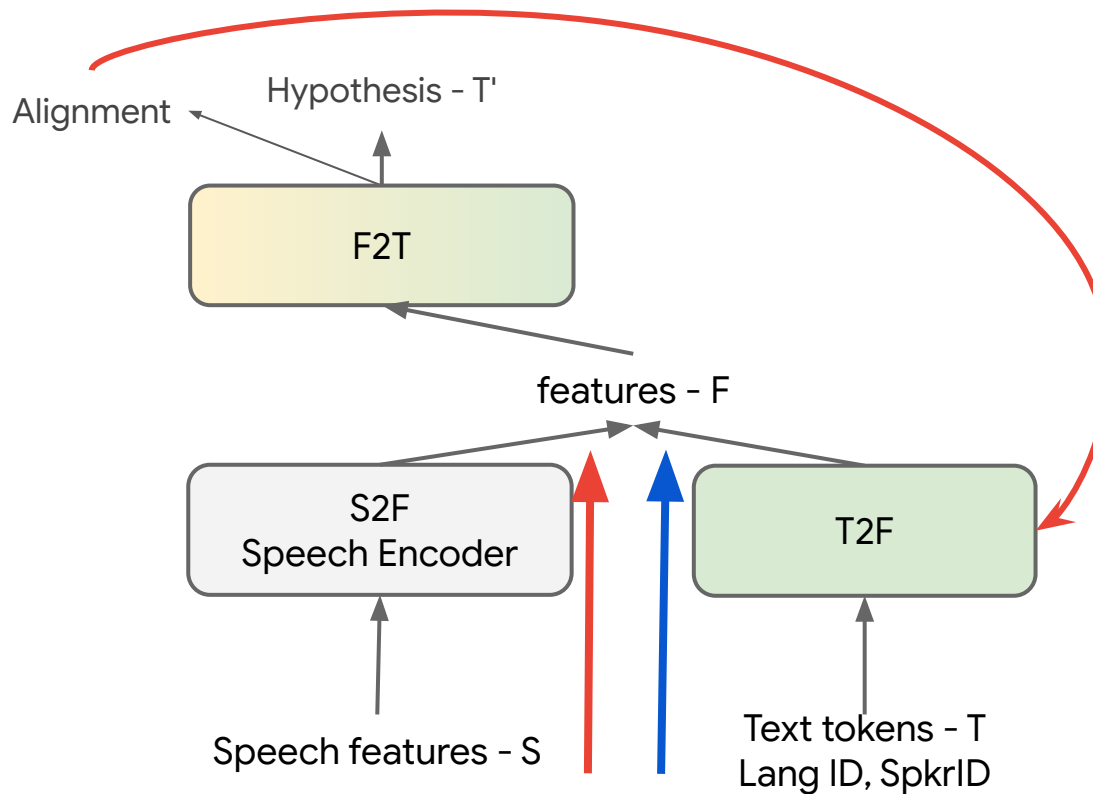
# Curriculum Training Process



# Foundation Model Pretraining



# Supervised Learning with Paired *Found* Data



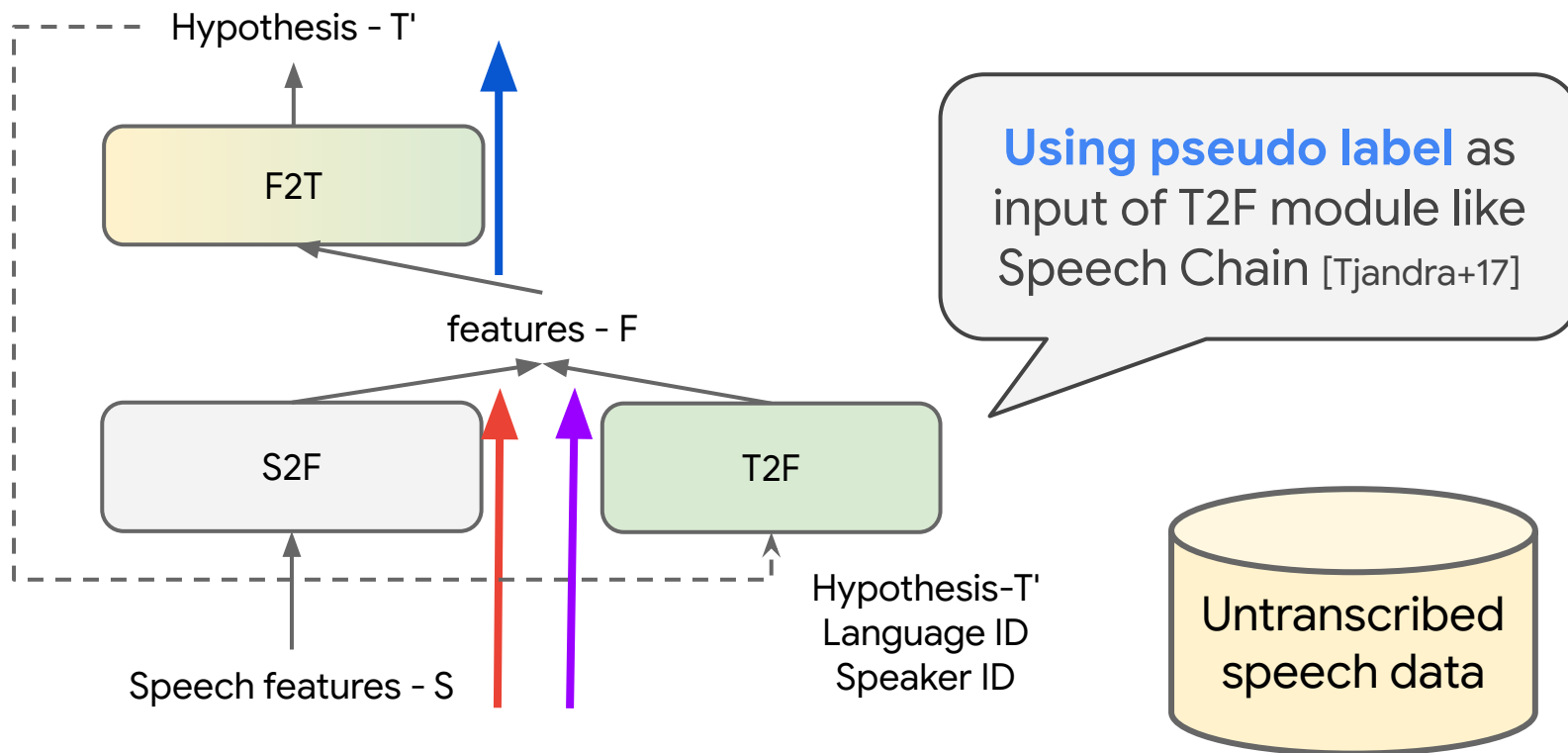
\* F2T module to get RNN-T alignment is jointly optimized.

\* RNN-T alignment is used for upsampling and duration model training in T2F module.

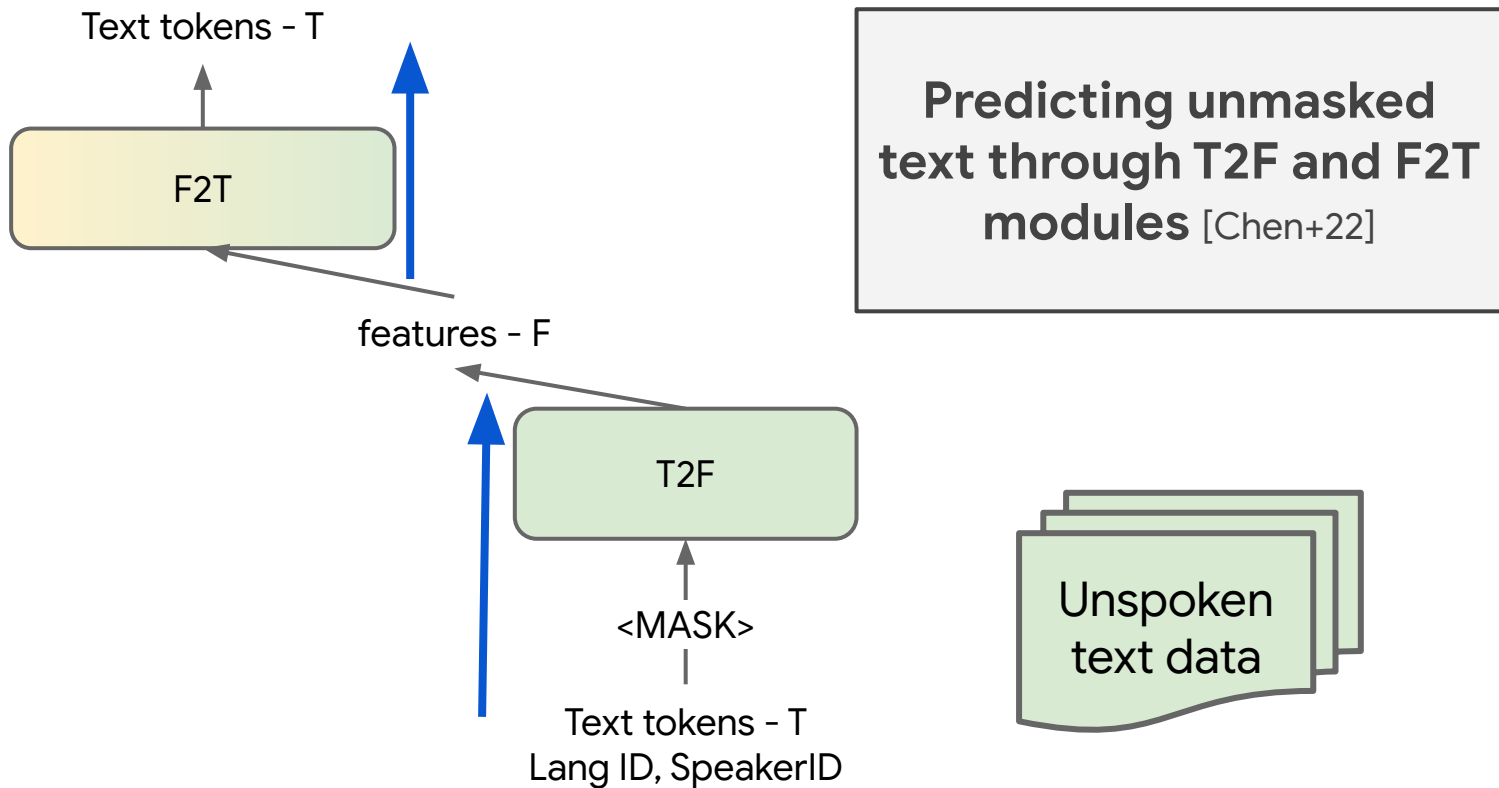
\* T2F module is trained to **predict self-supervised features**.



# Training with Speech-Only Data: *Pseudo-Labeling*



# Training with Text-Only Data: *Aligned MLM*



# Outline

1. Background

2. Method

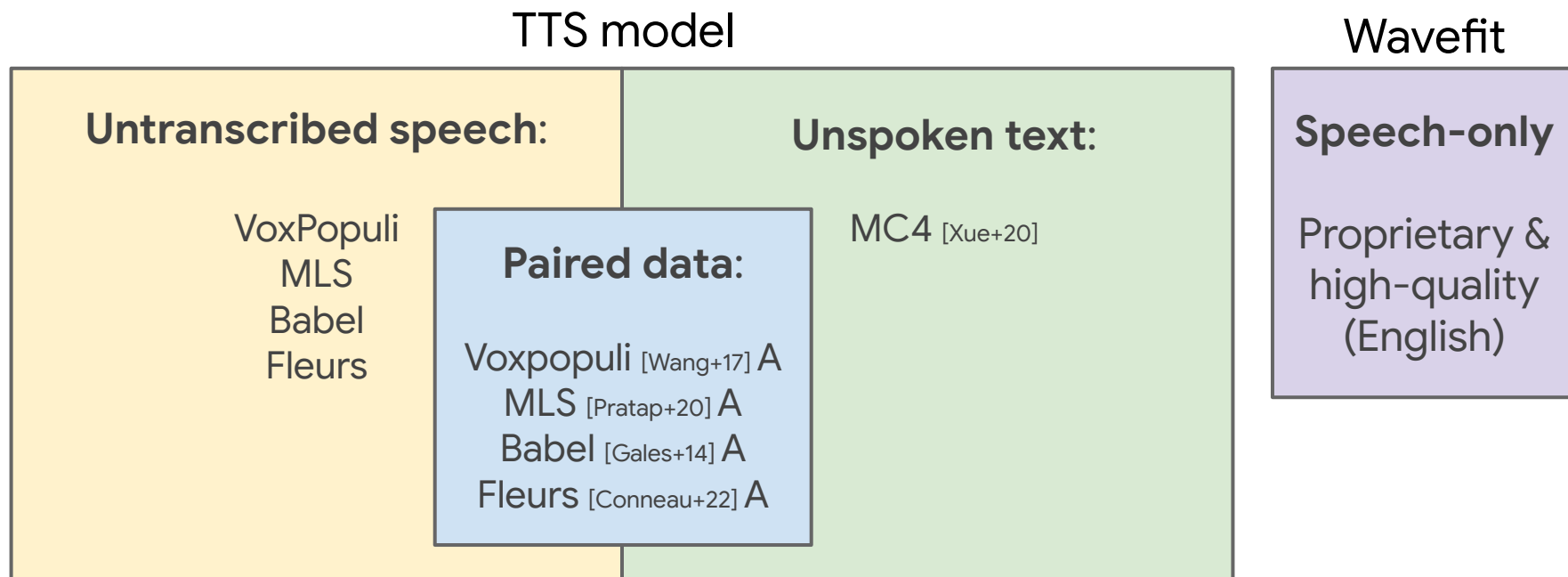
**3. Experimental Settings**

4. Results

5. Summary & Future Work

# Training Dataset

**We do not use any paired TTS data** (studio-recording audio & text).



\* Proprietary dataset with 56 langs was used for speech encoder pretraining.

# Settings

## Tokenization

**Grapheme:** 4096 vocab.  
Sentence-piece tokens

**Bytes:** UTF-8 bytes with  
256 tokens

## Training schemes

**Baseline** (w/o joint SSL)

**Proposed** (w/ joint SSL)

## Data Condition

**Zero** (No supervised)

**15m** (15min. Fleurs)

**Supervised** (MLS, Voxp,  
Babel, Fleurs)

# Data Setting for Zero

**Speech-only data:** Voxpopuli, MLS, Babel, Fleurs

**Text-only data:** MC4

**Paired ASR data:**

Voxpopuli, MLS, Babel, Fleurs

**Group A**

52 languages

French

Bengali

Kazakh

Pashto

...

...

**No paired data.**

**Group B**

50 languages

Spanish

Tamil

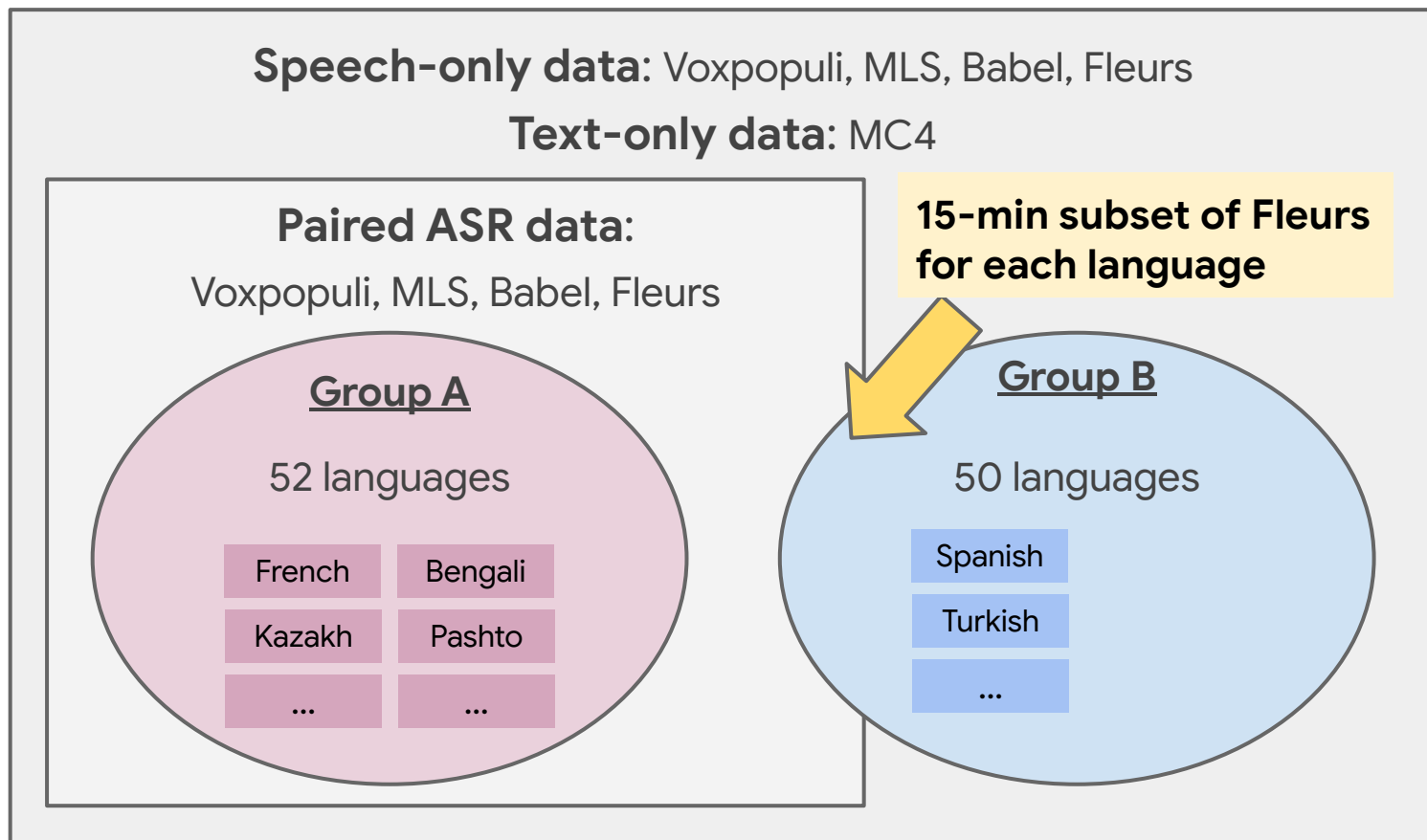
Turkish

Zulu

...

...

# Data Setting for 15m



# Data Setting for *Supervised*

**Speech-only data:** Voxpopuli, MLS, Babel, Fleurs

**Text-only data:** MC4

**Paired ASR data:**

Voxpopuli, MLS, Babel, Fleurs

**Group A**

52 languages

French

Bengali

Kazakh

Pashto

...

...

**Group B**

50 languages

Spanish

Tamil

Turkish

Zulu

...

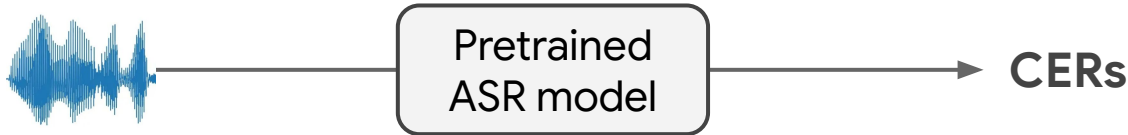
...



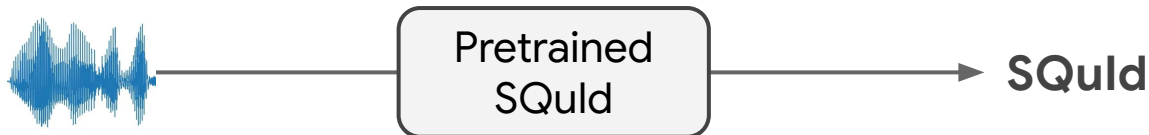
# Evaluation Metrics

1. Subjective **mean opinion score (MOS)** tests for **naturalness**  
4 languages with different language families.
2. **Character error rates (CER)** for **intelligibility** for all 50 Group B languages.

Synthetic speech



3. **SQuld** [Sellam+23]: **Automatic MOS prediction**, for all 50 Group B languages.



# Outline

1. Background

2. Method

3. Experimental Settings

**4. Results**

5. Summary & Future Work

# Main Results for All Languages

MOS: Avg. of **4** langs  
CER and SQuld: Avg. of **50** langs

	MOS	CER (%)	SQuld
<i>Groundtruth</i>	3.67	6.55	3.64
<i>Supervised</i>	3.21	6.39	3.88
<i>Zero (Baseline)</i>	2.48	28.28	3.84
<i>Zero (Proposed)</i>	2.53	23.44	3.77
<i>15m (Baseline)</i>	2.93	11.17	3.91
<i>15m (Proposed)</i>	3.18	7.33	3.88

# Main Results for All Languages

	MOS		SQuid
<i>Groundtruth</i>	3.67	6.55	3.64
<i>Supervised</i>	3.21	6.39	3.88
<i>Zero (Baseline)</i>	2.48	28.28	3.84
<i>Zero (Proposed)</i>	2.53	23.44	3.77
<i>15m (Baseline)</i>	2.93		3.91
<i>15m (Proposed)</i>	3.18		3.88

Due to noisy  
real-world  
found data

**3.18 MOS with  
15-min paired  
found data**

# Main Results for All Languages

	MOS	CER (%)	SQuid
<i>Groundtruth</i>	3.67	6.55	3.64
<i>Supervised</i>	3.21	6.39	3.88
<i>Zero (Baseline)</i>	2.48	28.28	3.84
<i>Zero (Proposed)</i>	2.53	23.44	3.77
<i>15m (Baseline)</i>	2.93	11.17	
<i>15m (Proposed)</i>	3.18	7.33	

7.33% CER with 15min paired found data

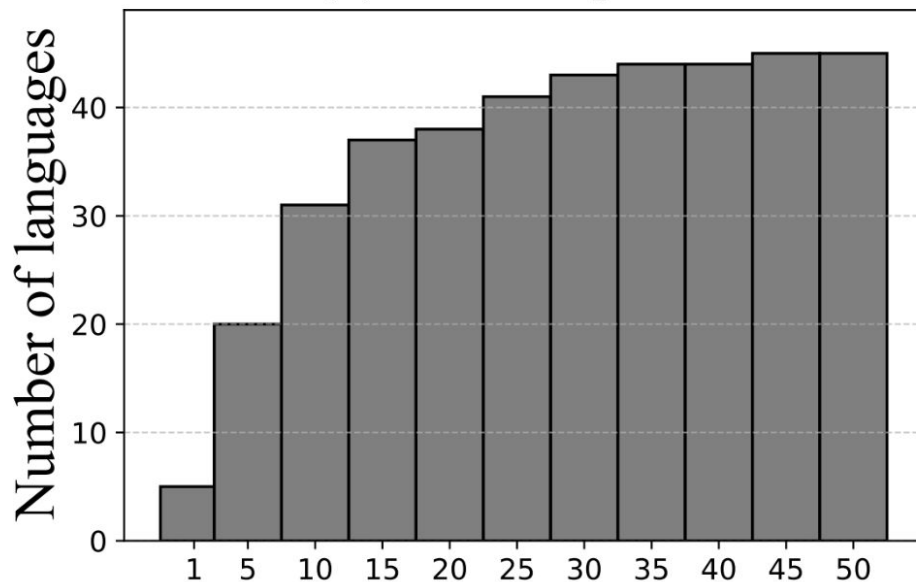
# Ablation Study of Each Method

Joint speech-text  
semi-supervised  
learning is effective

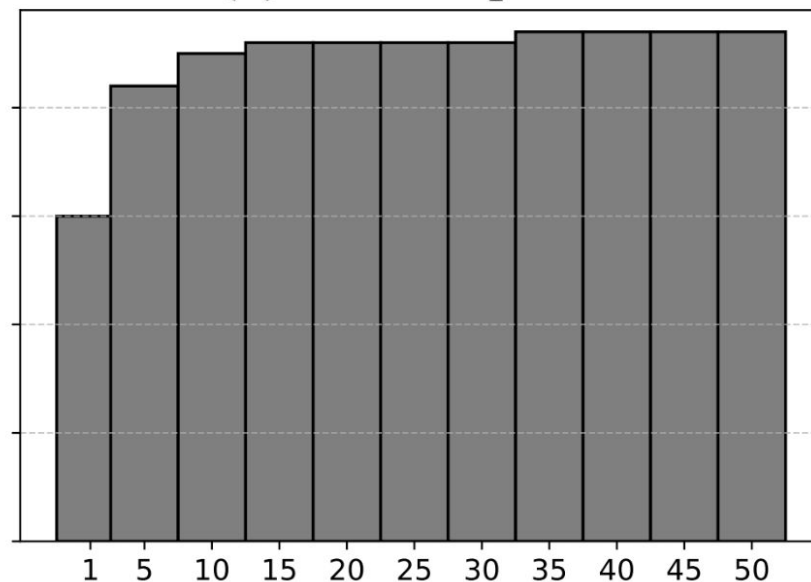
	CER (%)	
	<i>Zero</i>	<i>15m</i>
Baseline	28.28	11.17
+ Text MLM pretraining	26.13	8.47
+ Aligned text MLM	27.90	8.35
+ Pseudo labeling	<b>23.44</b>	<b>7.33</b>

# Histograms for TTS Language Extension

(a) *Zero Proposed*



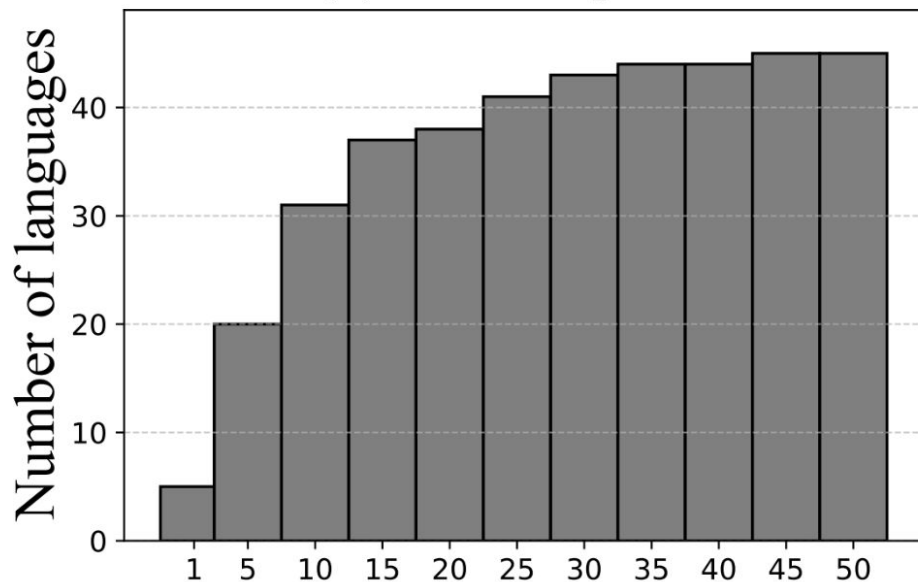
(b) *15m Proposed*



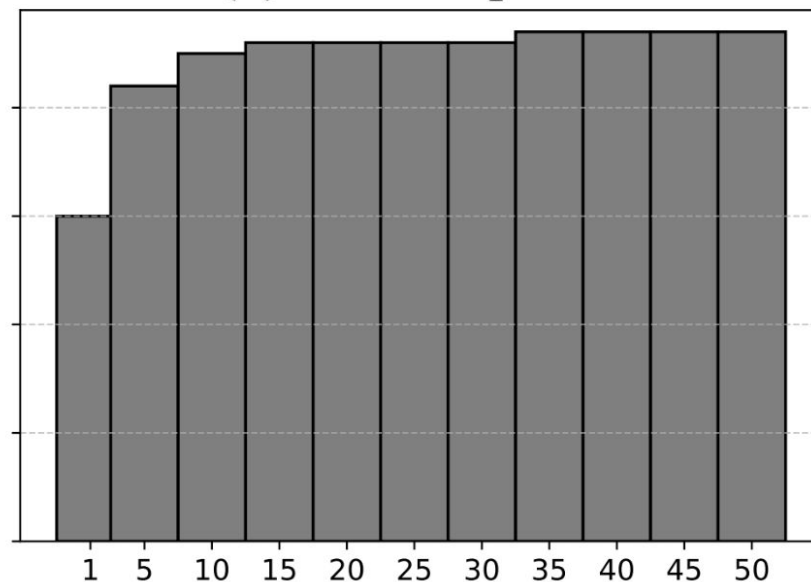
CER-diff to Groundtruth (%)

# Histograms for TTS Language Expansion

(a) *Zero Proposed*



(b) *15m Proposed*



**Zero: 20/50 langs  
within 5% CER diff.**

CER-diff to Groundtruth (%)





**15m: 42/50 langs  
within 5% CER diff.**







# Audio Samples

[https://google.github.io/tacotron/publications/extending\\_ts/](https://google.github.io/tacotron/publications/extending_ts/)

## (1) Spanish

Ground-truth	Supervised	Zero	15m
			

## (2) Zulu

Ground-truth	Supervised	Zero	15m
			

Thank you for your attention!

## Summary

**TTS language extension** with speech/text semi-supervised learning.

**Zero supervision:** 20/50 langs within 5% CER diff. to ground truth.

**Minimal supervision:** 42/50 langs within 5% CER diff. using 15-min Fleurs.

## Future work

Experiments on more languages and cross-lingual speaker transfer.