

# FastSpeech 2: A Text-to-Speech AI Model

Luke Buchanan, Seth Wertz

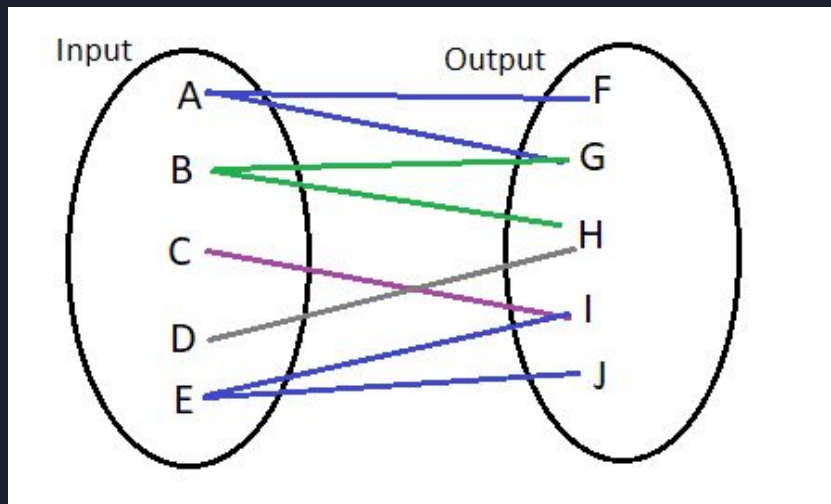


# Text-to-Speech Significance

- Increased Time Spent Online
- Accessibility
  - Visual
  - Cognitive
  - Age
- Efficiency

# Challenges

- One-to-Many Mapping Problem
  - Example: Enough, through, cough, bought, drought, dough
- Human Speech Variance





# FastSpeech 1

- Microsoft and Zhejiang University 2019
- Phoneme Encoding
- The Student-Teacher Model
- Autoregression

# FastSpeech 2 Model Overview

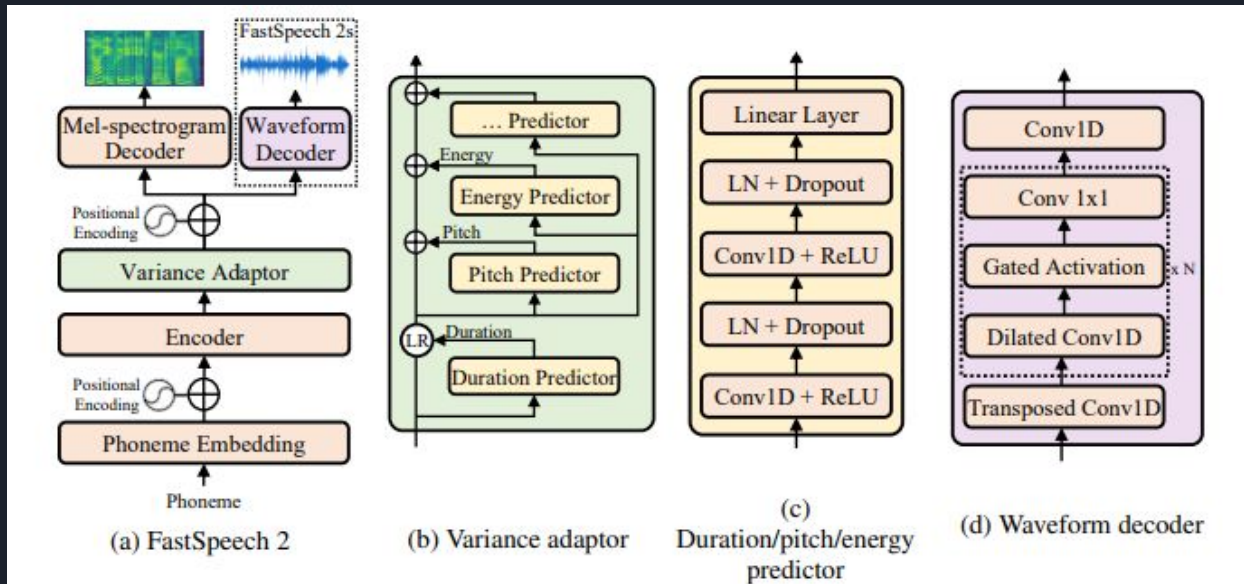


Figure 1: The overall architecture for FastSpeech 2 and 2s. LR in subfigure (b) denotes the length regulator proposed in FastSpeech. LN in subfigure (c) denotes layer normalization.



# FastSpeech 2 Model Overview

- Encoder based on self-attention transformer layers
  - Uses phoneme embeddings
- Variance adaptor
  - Predicts duration, pitch, and energy of phonemes
  - Each predictor uses CNN structure
- Convolutional layer decoder



# FastSpeech 2s

- Fully end-to-end TTS module
  - Doesn't generate intermediate spectrogram
- Generates speech waveform directly from text
- Higher quality speech than FastSpeech 1
- Faster voice synthesis than FastSpeech 2

# Results

- Trained on 13,000 English voice samples
- Mean Opinion Score (MOS)
- Ground Truth (GT)

Method	MOS
<i>GT</i>	$4.30 \pm 0.07$
<i>GT (Mel + PWG)</i>	$3.92 \pm 0.08$
<i>Tacotron 2 (Shen et al., 2018) (Mel + PWG)</i>	$3.70 \pm 0.08$
<i>Transformer TTS (Li et al., 2019) (Mel + PWG)</i>	$3.72 \pm 0.07$
<i>FastSpeech (Ren et al., 2019) (Mel + PWG)</i>	$3.68 \pm 0.09$
<i>FastSpeech 2 (Mel + PWG)</i>	$3.83 \pm 0.08$
<i>FastSpeech 2s</i>	$3.71 \pm 0.09$

(a) The MOS with 95% confidence intervals.

Method	CMOS
<i>FastSpeech 2</i>	0.000
<i>FastSpeech</i>	-0.885
<i>Transformer TTS</i>	-0.235

(b) CMOS comparison.

Table 1: Audio quality comparison.





# Model Speedup

- Faster to Train, Slower Inference

Method	Training Time (h)	Inference Speed (RTF)	Inference Speedup
<i>Transformer TTS (Li et al., 2019)</i>	38.64	$9.32 \times 10^{-1}$	/
<i>FastSpeech (Ren et al., 2019)</i>	53.12	$1.92 \times 10^{-2}$	48.5×
<i>FastSpeech 2</i>	<b>17.02</b>	$1.95 \times 10^{-2}$	47.8×
<i>FastSpeech 2s</i>	92.18	<b><math>1.80 \times 10^{-2}</math></b>	<b>51.8×</b>



# Demo

<https://cmchien.ttic.edu/FastSpeech2/>

<https://github.com/ming024/FastSpeech2>