

TTS (Text to speech)

SSR / VR on next page

- TTS (Text to speech) \Rightarrow input is text \rightarrow output is spoken audio
 - o goal: to convert text into natural human-like speech

Steps:

- 1) Text analysis/normalization
 - Convert num, abbreviations etc into spoken form

Ex 12.50 2
Twelve dollar fifty cents

- 2) Phoneme conversion

- translate words into their phonetic pronunciation

- 3) Prosody Prediction

- predict Pitch, Rhythm, and stress patterns to make speech sound natural (instead of robotic)

- 4) Wave form Generation

- Earlier: Concatenative TTS - stitched together pre-recorded human voice samples
- Later: parametric TTS - generated speech from acoustic parameters (less natural sounding)
- Now: Neural TTS - Neural vocoders like WaveNet generate extremely realistic waveforms

- Modern pipeline Example (Tacotron 2 + WaveNet)

- 1) Tacotron 2 - takes text, outputs a mel spectrogram (time vs Frequency plot of audio energy)

- 2) WaveNet / HiFi GAN - takes that spectrogram & generates realistic audio waveforms