2022.7.15 Shimin Koza





# How do Machines Speak? Progress and Challenges of Speech Synthesis

Erica Cooper クーパー・エリカ

# Synthesized speech in everyday life

#### Virtual assistants



#### Home devices





GPS devices



# Entertainment



## Speech synthesis can break down barriers

Screen readers for vision impairments



Communication devices for vocal disabilities



https://en.wikipedia.org/wiki/Stephen\_Hawking

#### Language learning apps



Speech-to-speech translation



## History of speech synthesis



https://commons.wikimedia.org/wiki/File:Euphoniafaber.jpg

#### Joseph Faber's Euphonia 1845



https://www.youtube.com/watch?v=0rAyrmm7vv0

Bell Labs Voder 1939

# History of speech synthesis

- **1970s:** Rule-based synthesis to map linguistic units to audio
- **1980s:** Sample-based synthesis to join recordings of the smallest speech units
- **1990s:** Sample-based synthesis using longer recordings
- **2000s:** Machine learning based synthesis learning how to map text into sound from data
- **2010s:** Machine learning based synthesis using neural networks







Text normalization

- Requres real-world knowledge
- Riverside Dr. -> Riverside Drive
- Dr. Smith -> Doctor Smith
- How to expand numbers and abbreviations
  - NASA, SAT
  - **128GB**

Phonetization

• Requres linguistic knowledge

riverside drive rivərsaid draiv IPA R IH V ER S AY D D R AY V ARPAbet

- Lookup dictionary, handcrafted or learned rules
- Can also include information about stress, timing, and other information in speech



Learn a mapping between the linguistic representation created from the text, and an acoustic representation of the audio.

The mapping is learned from many examples of pairs of matching text and audio.



Sometimes an intermediate representation of audio is used.

In that case, we finally have to convert from the intermediate representation into audio.



In the past, these were all separatelydeveloped components.





Modern synthesis methods use one model for some or all of these components and learn multiple steps together.

- Why personalized speech synthesis?
  - Assistive technology for the speaking impaired
  - Narration for social media videos and presentations
  - Quick editing and correction of spoken audio

Welcome to my podcast. Today's episode is a new story.

ախոսիտ պիտարտարտ

Hello and welcome to my podcast. Today's episode is an interesting story.

allocallocallocallocallocal

How can I make a synthesized voice that sounds like me?

#### Train a single speaker model



How can I make a synthesized voice that sounds like me?

Train a single speaker model



How can I make a synthesized voice that sounds like me?

Train a single speaker model





How can I make a synthesized voice that sounds like me?



How can I make a synthesized voice that sounds like me?



Human-sounding speech synthesis requires dozens of hours of professionally-recorded speech.



https://www.curiousspeckle.net/2013/05/13/silent-room-room-for-silence/ Bell Labs anechoic chamber, 1947 Photo credit: Eric Schaal

- There are more than 7,000 languages spoken in the world !
- Popular commercial systems typically only support ~30 languages
- Apple MacinTalk: "over 30 languages"
- Amazon Polly: "31 languages"
- Google: "32 languages and variants"



Source: https://en.wikipedia.org/wiki/Linguistic\_diversity\_index

- There are more than 7,000 languages spoken in the world !
- Popular commercial systems typically only support ~30 languages
- Apple MacinTalk: "over 30 languages"
- Amazon Polly: "31 languages"
- Google: "32 languages and variants"

Can we build a speech synthesizer for a new language even if we don't have large amounts of highquality speech recordings?



Source: https://en.wikipedia.org/wiki/Linguistic\_diversity\_index

Adapt a multi-speaker model



Adapt a multi-language model



# Deepfakes vs. personalized speech synthesis

- Same technology, different purposes
- Personalized speech synthesis:
  - For your own personal use
  - To assist your everyday life
  - To create your own content
  - For entertainment

#### • Deepfakes:

- Created without the person's knowledge
- Intended to deceive
- A type of fraud
- Spread misinformation in society

#### TECHNOLOGY

#### An artificial-intelligence first: Voice-mimicking software reportedly used in a major theft

By <u>Drew Harwell</u>

September 4, 2019 at 6:27 p.m. EDT

https://www.washingtonpost.com/technology/2019/09/04/an-artificialintelligence-first-voice-mimicking-software-reportedly-used-major-theft/



https://www.ic3.gov/Media/Y2022/PSA220628

#### TECH · SCAMS

#### Beware: Phone scammers are using this new sci-fi tool to fleece victims

#### BY JENNIFER ALSEVER

May 5, 2021 8:30 AM GMT+9

https://fortune.com/2021/05/04/voice-cloning-fraud-ai-deepfakes-phone-scams/

• Synthesized speech can fool humans



• Synthesized speech can fool machines



• Synthesized speech can be detected





Synthesized speech













#### Speech Synthesis: open research problems

- Spoofing and deepfake detection
  - Humans and machines perceive speech differently and make different kinds of errors
  - Speech synthesis is always improving, so detectors need to keep up
  - Detection of new, unseen synthesizers
- Speech synthesis for low-resource languages
  - How to make it as easy as possible to build a speech synthesizer for a new language
- Expressive speech synthesis
  - How to synthesize in styles other than neutral
  - How to make synthesized speech that is entertaining to listen to
- Faster and more lightweight speech synthesis
  - Neural network models require large computational resources