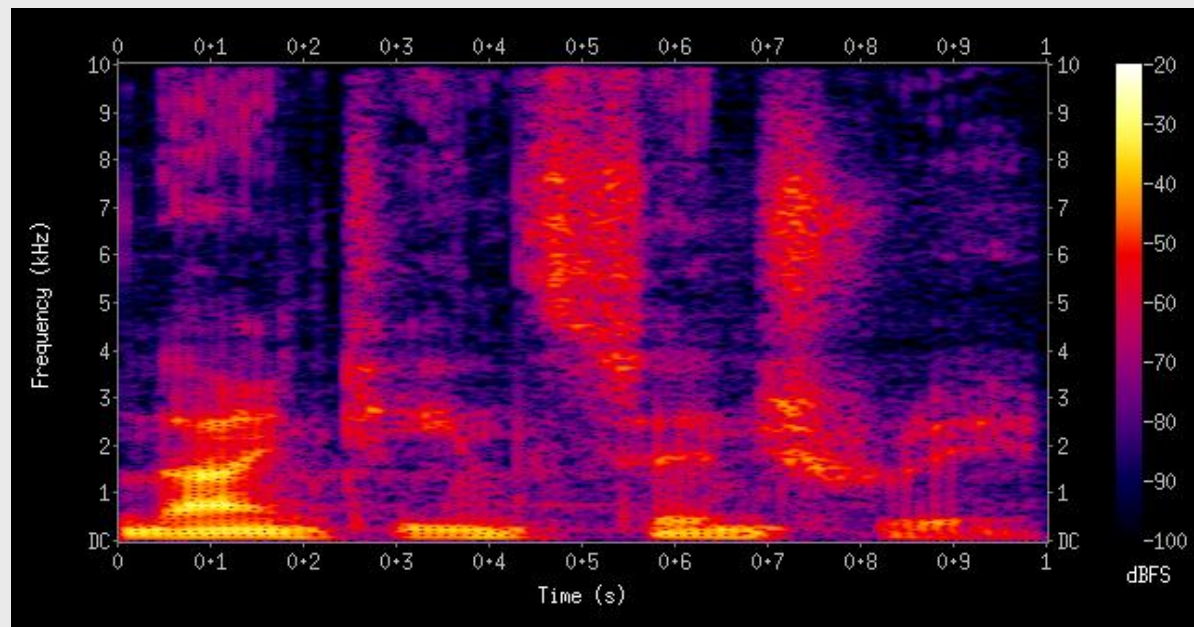# VOICE CLONING

# AI Voice Cloning

- **AI can easily mimic someone's speaking style in text with transformers**

- **This is possible because the AI is able to represent text with a good embedding vector**

- **If the AI wanted to represent the sound of someone's voice, what would the embedding vector represent?**
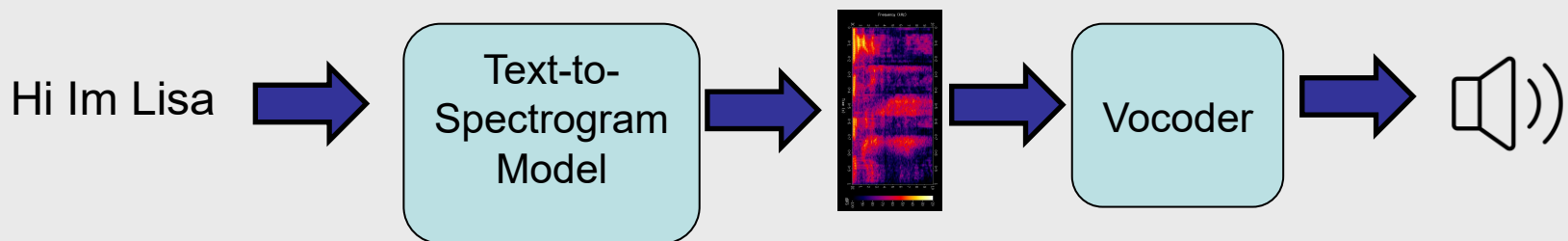
# Spectrograms: How Computers Represent Voices

- **Spectrogram – represents the frequencies in an audio signal as it varies in time**
  - **Like an audio fingerprint**

# Text-to-Speech (TTS) Models

- **Text to speech (TTS) models are neural networks that convert text into spectrograms**

- **The spectrograms are then turned into voice audio using neural networks called vocoders**

Hi Im Lisa → **Text-to-Spectrogram Model** → → **Vocoder** → 🔊

# First Good TTS Model: Tacotron

## TACOTRON: TOWARDS END-TO-END SPEECH SYNTHESIS

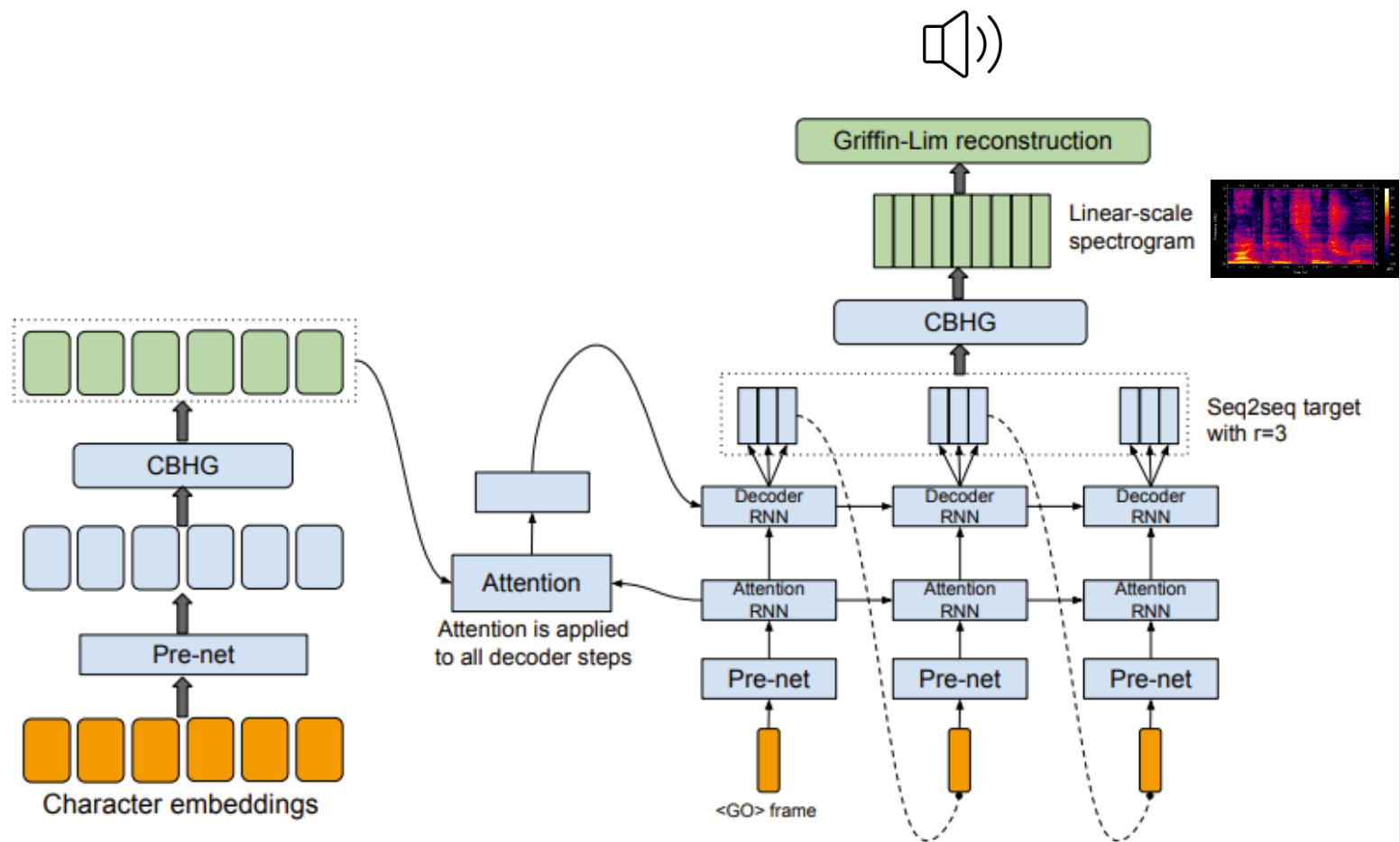Yuxuan Wang*, RJ Skerry-Ryan*, Daisy Stanton, Yonghui Wu, Ron J. Weiss[†], Navdeep Jaitly,

Zongheng Yang, Ying Xiao*, Zhifeng Chen, Samy Bengio[†], Quoc Le, Yannis Agiomyrgiannakis,

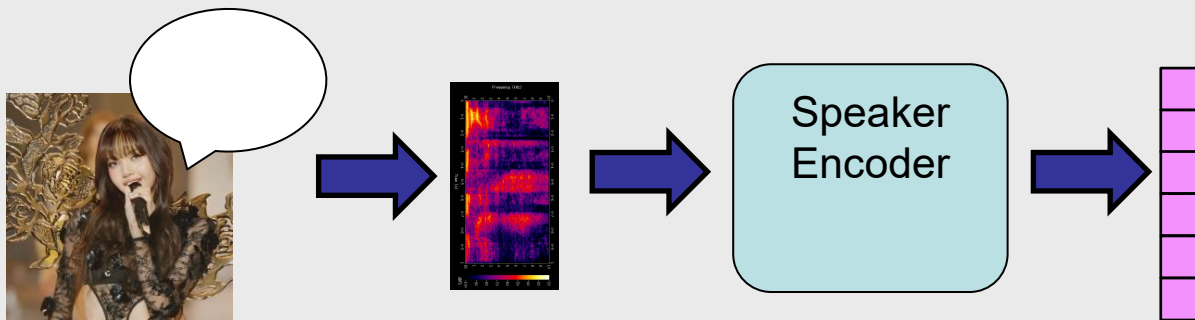Rob Clark, Rif A. Saurous*

Google, Inc.
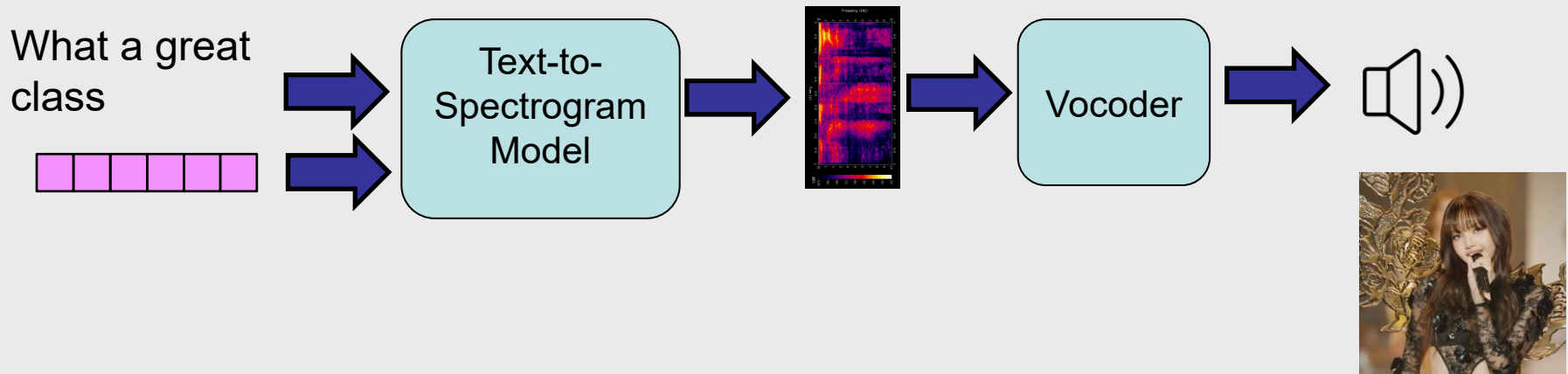{yxwang,rjryan,rif}@google.com

# Tacotron Architecture



Hi Im Lisa

# Representing a Speaker

- **Spectrograms are like voice fingerprints**
- **Voice cloning neural networks called "speaker encoders" take spectrograms and turn them into vectors**
- **Trained using contrastive learning**

# Voice Cloning

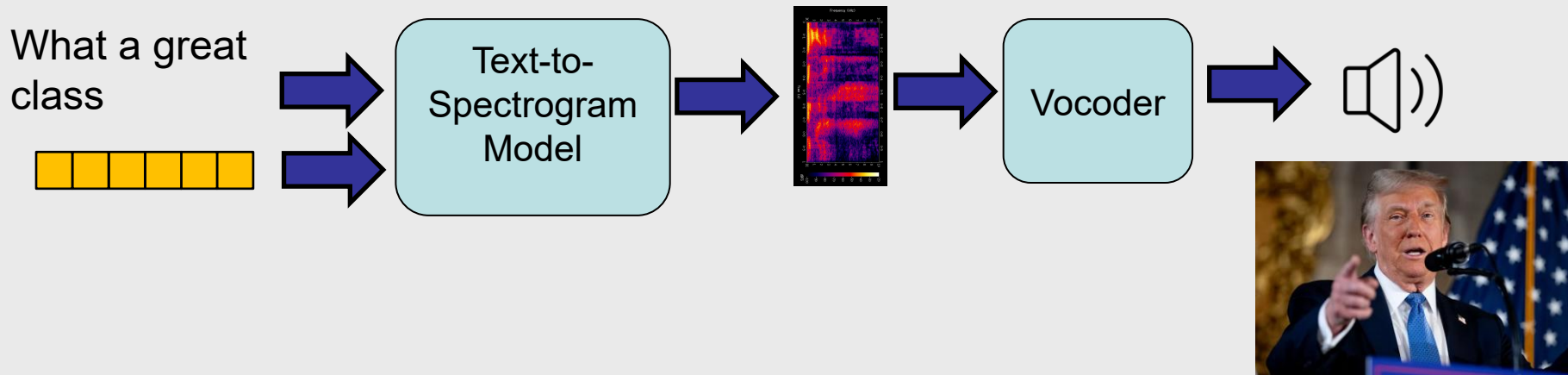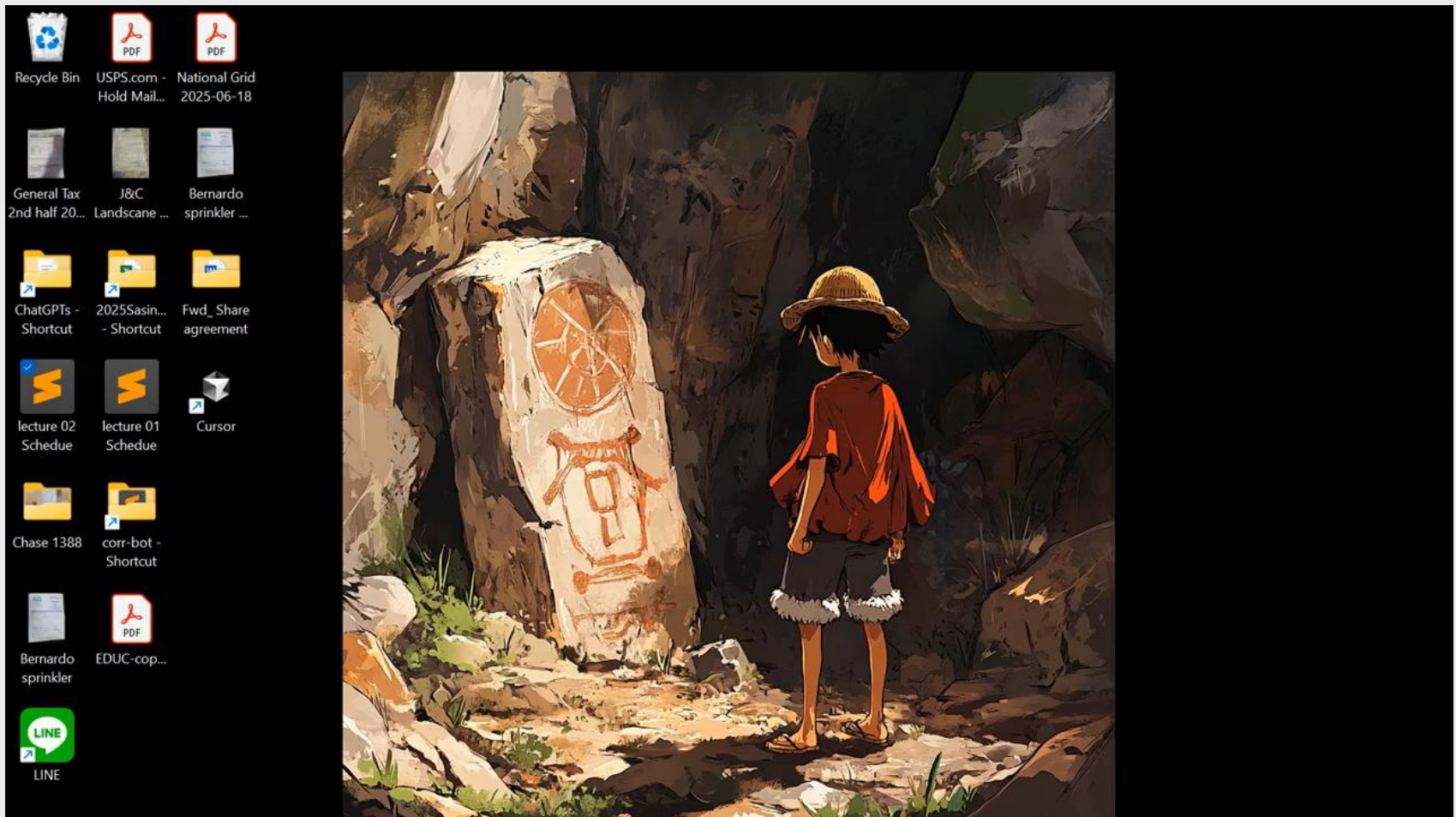- **Once we have the speaker vector, we can just add it to the TTS model**

What a great class → Text-to-Spectrogram Model → Vocoder → 🔊

# Voice Cloning

- **Once we have the speaker vector, we can just add it to the TTS model**

What a great class

# ElevenLabs

- **ElevenLabs is one of the most advanced voice cloning tools: https://elevenlabs.io/**


- **Many voice cloning features**
  - **Voice library (use voices others cloned)**
  - **Instant voice cloning (10 seconds of audio)**
  - **Professional voice cloning (30 minutes of audio)**
  - **Voice design (create voice by text prompt)**
  - **Voice changer (change the voice in an audio file)**

# Voice Changer

# Permission Granted to Voice Clone

# Permission Granted to Voice Clone



James Earl Jones Signed Over Rights For AI To Recreate Darth Vader's Voice

By Tim Lammers, Contributor. ⓘ I cover Hollywood and entertainment.

Follow Author

Published Sep 09, 2024, 06:08pm EDT, Updated Sep 09, 2024, 11:55pm EDT

# Permission Not Granted to Voice Clone



Robert Downey Jr. vows there will never be a digital AI replica of him on-screen

News    By Eric Hal Schwartz published October 31, 2024

Ultron isn't Iron Man's only AI foe

# Voice Cloning Ethics

- **If using a voice clone, be transparent**

- **Don't use anyone's voice without their permission**

- **Don't make the voice clone say offensive things**

- **Don't create deepfake audio that will cause panic…**
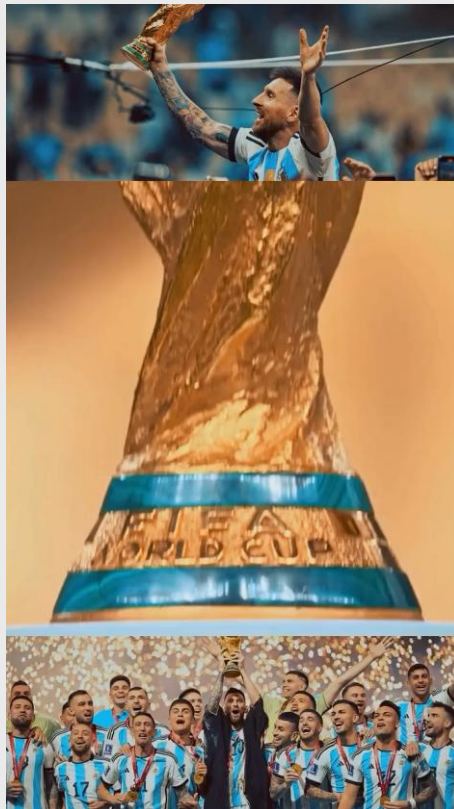
# How to Clone a Voice

- **ElevenLabs allows us to clone voices with only a few seconds of audio**

- **Can obtain audio files from YouTube**
  - **Choose a YouTube video of the subject speaking alone or mostly alone**
  - **Download audio of YouTube video (https://tuberipper.net/73/)**
  - **Clip the audio to 10 to 30 second clips (https://audiotrimmer.com/)**

# Video Narrations

- **If we give the frames of a video to AI plus some instructions, it can give us a narration for the video**

- **We already have methods to gives multiple images and context text to AI and get a text description**

- **We just need to extract the images from the video**
  - **We will usually sample around 10-20 images per video**
  - **This is near the limit of what OpenAI allows**

# Old HW Problem

- **Make a narration of this video**
- **Bonus – clone a voice and have it narrate the video**



Richard Ogu
Yale MAM
Class of 2024

# Coding Session

- **We will clone a voice**

- **We will clone a GitHub repository with some useful code to get us started**

- **We will build an app that narrates a video using the cloned voice**

- **AI tools used**
  - **OpenAI**
  - **ElevenLabs**