

Text-to-Speech & Voice Cloning System

1. Introduction

1.1 Background

Speech is the most natural and efficient way of human communication. With advancements in artificial intelligence (AI), computers are increasingly capable of generating speech that is nearly indistinguishable from real human voices. **Text-to-Speech (TTS)** refers to the technology that converts written text into spoken audio. Early TTS systems relied on rule-based or concatenative approaches, which produced robotic and monotonous speech. Later, statistical models such as **Hidden Markov Models (HMMs)** improved flexibility but still lacked natural expressiveness.

The breakthrough in natural-sounding speech came with **neural network-based models** like **Tacotron**, **Deep Voice**, and **WaveNet**, which enabled end-to-end learning of text-to-speech conversion. More recent models such as **VITS** and **Glow-TTS** combine efficiency and high-quality prosody generation, allowing real-time, natural voice synthesis.

Parallel to TTS, **voice cloning** has emerged as a key area of research. Unlike traditional TTS systems that generate speech in a generic synthetic voice, voice cloning focuses on reproducing a **specific individual's voice**. Using **speaker embeddings** and **zero-shot learning**, modern models can replicate a person's voice from as little as a few seconds of recorded speech.

1.2 Importance of Speech Synthesis and Voice Cloning

The ability to synthesize speech and clone voices has significant real-world importance:

- **Accessibility** – Assisting visually impaired users by reading content aloud or helping mute individuals generate speech.
- **Personalization** – Virtual assistants (e.g., Alexa, Siri, Google Assistant) speaking in customized voices.
- **Content Creation** – Automating audiobook narration, podcast generation, or YouTube voiceovers.
- **Healthcare** – Helping patients who lose their voice due to illness restore their digital voice.
- **Entertainment** – Dubbing films, gaming characters, or creating virtual influencers with unique voices.

1.3 Applications of This Project

This project integrates **TTS** and **voice cloning** into a single platform that:

- Generates natural-sounding speech from any text.
- Clones the user's voice from short recordings to synthesize speech in that voice.
- Provides both backend APIs (FastAPI) and a user-friendly frontend (Streamlit).
- Runs locally with simple one-click startup.

2. Problem Statement

2.1 Existing Challenges in TTS

Despite rapid advancements, current Text-to-Speech (TTS) systems still face several limitations:

- **Robotic Sounding Voices:** Many TTS systems struggle with natural intonation, rhythm, and prosody, leading to synthetic or monotonous speech.
- **Limited Expressiveness:** Capturing emotions such as happiness, sadness, or emphasis in speech remains a challenge.
- **Language and Accent Support:** Most models are optimized for English, with limited performance in multilingual or accented speech scenarios.
- **High Computational Requirements:** State-of-the-art models like WaveNet or Tacotron demand large amounts of GPU memory and computational power, making real-time deployment difficult on low-resource devices.
- **Lack of Personalization:** Traditional systems generate speech in a fixed synthetic voice that does not adapt to individual users.

2.2 Need for Voice Personalization

While generic TTS can provide functional speech output, **personalized voice synthesis** is becoming increasingly important:

- **User Identity:** A personalized voice allows users to retain their unique identity in digital interactions, which is especially critical for individuals who lose their ability to speak due to medical conditions.
- **Enhanced Engagement:** Virtual assistants and chatbots become more relatable when they speak in a voice chosen or cloned by the user.
- **Accessibility and Inclusivity:** Personalized voice solutions empower differently-abled individuals by giving them a voice that feels natural and personal.
- **Content Customization:** In media and entertainment, content creators can generate unique branded voices for audiobooks, podcasts, or gaming characters.

Thus, there is a strong demand for a system that not only generates **natural-sounding speech** but also supports **voice cloning for personalization**.

3. Objectives

The primary objective of this project is to design and implement a **Text-to-Speech (TTS) and Voice Cloning system** that can generate human-like speech while supporting personalization through voice cloning.

3.1 Goals of This Project

- **Natural Text-to-Speech Conversion:** Develop a system that converts written text into fluent, natural, and expressive speech output.

- **Voice Cloning Capability:** Enable cloning of a user's voice using short audio recordings, allowing the system to reproduce any given text in that voice.
- **Multilingual and Multi-Speaker Support:** Integrate a TTS model that can handle multiple languages and support different speaker styles.
- **Backend Implementation (FastAPI):** Provide RESTful APIs for text-to-speech and voice cloning, ensuring modularity and scalability.
- **Frontend Interface (Streamlit):** Design an easy-to-use UI so non-technical users can access the system without interacting directly with code.
- **Audio Preprocessing:** Include preprocessing pipelines (conversion, resampling, normalization) to ensure consistent audio quality across different input formats.
- **Automation and Ease of Use:** Create a one-click startup mechanism (`start_project.bat`) that automatically launches both backend and frontend services.
- **Practical Usability:** Ensure the system runs on local hardware without requiring cloud infrastructure, making it accessible for experimentation and personal use.

4. Theoretical Background

4.1 Text-to-Speech (TTS) Fundamentals

Text-to-Speech (TTS) is the process of converting written text into spoken audio. The fundamental pipeline of TTS consists of:

1. **Text Analysis (Front-End):**
 - Normalization (expanding abbreviations, handling numbers, punctuation).
 - Grapheme-to-phoneme conversion (mapping letters to sounds).
2. **Acoustic Modeling (Middle Layer):**
 - Predicting speech features such as mel-spectrograms, which represent the frequency distribution of speech over time.
3. **Vocoder (Back-End):**
 - Converting mel-spectrograms into raw audio waveforms.
 - Early systems used signal-processing-based vocoders (e.g., Griffin-Lim), but modern TTS uses neural vocoders like **WaveNet**, **HiFi-GAN**, and **WaveGlow** for higher realism.

The primary goals of TTS are **naturalness**, **intelligibility**, and **speaker diversity**.

4.2 Neural TTS Models

Traditional TTS systems were rule-based or concatenative, often producing robotic or unnatural speech. Modern **neural TTS models** have revolutionized this space:

- **Tacotron & Tacotron 2**
 - End-to-end sequence-to-sequence models using attention.
 - Input: characters/phonemes → Output: mel-spectrograms.
 - Vocoder: WaveNet/HiFi-GAN used to generate speech.

- Known for natural intonation, but slow inference.
- **Glow-TTS**
 - Based on **normalizing flows**, enabling fast, parallel speech generation.
 - Removes autoregressive dependencies (faster than Tacotron).
 - Provides controllability over speech duration and prosody.
- **VITS (Variational Inference TTS)**
 - Combines **variational autoencoders (VAEs)**, **normalizing flows**, and **GAN-based vocoders** into a single end-to-end model.
 - Produces highly natural speech with efficient training and inference.
 - Capable of zero-shot speaker adaptation when combined with speaker embeddings.

These models represent the **evolution from text → spectrogram → waveform** pipelines toward **integrated architectures** like VITS.

4.3 Voice Cloning Fundamentals

Voice cloning is the task of synthesizing speech in the **unique voice of a target speaker** using limited data.

- **Types of Voice Cloning:**
 - **Speaker Adaptation:** Requires several minutes/hours of training data from a speaker.
 - **Zero-Shot Voice Cloning:** Uses only a short sample (5–10 seconds) to extract speaker embeddings and generate speech in that voice.
- **Key Components:**
 1. **Speaker Encoder:** Learns to map audio samples into fixed-length **speaker embeddings** (numerical representation of voice identity).
 2. **TTS Model:** Takes text + speaker embedding as input → generates mel-spectrogram.
 3. **Vocoder:** Converts mel-spectrogram to waveform.

Voice cloning enables **personalized assistants, audiobooks, dubbing, and assistive tools**, but also raises ethical concerns like misuse in deepfakes.

4.4 Coqui YourTTS Model Overview

For this project, we use **YourTTS** from **Coqui-AI's TTS library**:

- **Model Family:** Built upon **VITS**, extended for multilingual and multi-speaker TTS.
- **Features:**
 - **Multilingual support:** English (en), French (fr-fr), Portuguese (pt-br).
 - **Multi-speaker synthesis:** Can generate speech in multiple speaker styles.
 - **Zero-shot voice cloning:** With just a few seconds of sample audio, the model can clone voices.
 - **Speaker embeddings:** Uses a pretrained speaker encoder to capture speaker identity.
- **Why YourTTS?**
 - End-to-end architecture → easier to integrate.
 - Open-source and optimized for real-time inference.

- Suitable for running locally without requiring massive GPUs.

Thus, **YourTTS** provides a balance between **accuracy**, **efficiency**, and **practical usability**, making it ideal for this project.

5. Methodology

The development of this project followed a modular approach where each component was designed, tested, and integrated to form a complete end-to-end system.

5.1 Model Selection

We selected **YourTTS** from the **Coqui TTS library** as the core model for both **text-to-speech** and **voice cloning**.

- **Reason for Selection:**
 - Supports **multilingual synthesis** (en, fr-fr, pt-br).
 - Enables **multi-speaker TTS** with pre-trained voices.
 - Provides **zero-shot voice cloning** using short audio samples.
 - Lightweight compared to larger models like Tacotron2 + WaveNet, making it feasible to run locally.

Thus, the model serves as a unified solution for both **TTS** and **cloning tasks**.

5.2 Preprocessing Pipeline

To ensure reliable model performance, we built a **preprocessing module** that handles raw audio before feeding it into the TTS engine.

- **Audio Format Conversion:** Converts `.m4a` or `.mp3` recordings into `.wav` format using **FFmpeg/Pydub**.
- **Resampling:** Ensures all input audio is resampled to **16 kHz mono**, which is the expected format for YourTTS.
- **Normalization:** Removes unwanted noise, trims silence, and standardizes loudness to produce cleaner embeddings.
- **Output Handling:** Preprocessed files are stored in the `data/samples/` folder for consistent access by the voice cloning module.

This pipeline ensures that audio inputs from different devices/platforms are standardized for model compatibility.

5.3 Backend (FastAPI for API Services)

We implemented the backend using **FastAPI**, which provides a lightweight and high-performance framework for serving ML models.

- **Endpoints:**
 - `/tts/` → Accepts text and generates speech as `.wav`.
 - `/clone/` → Accepts text + audio sample, performs preprocessing, and generates cloned speech.
- **Features:**
 - Model is **loaded once** on server startup, reducing inference time.
 - Uses `FileResponse` to return generated audio directly to clients.
 - Error handling ensures that missing parameters or invalid files are managed gracefully.

This makes the system accessible through simple **REST API calls**, allowing both frontend and external applications to use it.

5.4 Frontend (Streamlit UI)

To provide a user-friendly interface, we developed a **Streamlit-based frontend**.

- **Modules in the UI:**
 - **Text-to-Speech Panel:** User enters text → audio is generated and played back.
 - **Voice Cloning Panel:** User uploads a short audio sample and enters text → cloned speech is produced.
 - **Audio Player:** Embedded playback widgets allow users to directly listen to results without opening external tools.
- **Why Streamlit?**
 - Rapid prototyping with minimal code.
 - Easy integration with Python backend via `requests`.
 - Intuitive for non-technical users.

Thus, even users without coding knowledge can use the system seamlessly.

5.5 Automation (start_project.bat)

One of the major goals was **ease of use**. Manually running multiple terminals for backend and frontend was inefficient.

- To solve this, we created a **Windows batch script** `start_project.bat` that:
 1. Activates the Python virtual environment.
 2. Starts the FastAPI backend (`uvicorn app:app --reload`).
 3. Launches the Streamlit frontend (`streamlit run ui/app_ui.py`).

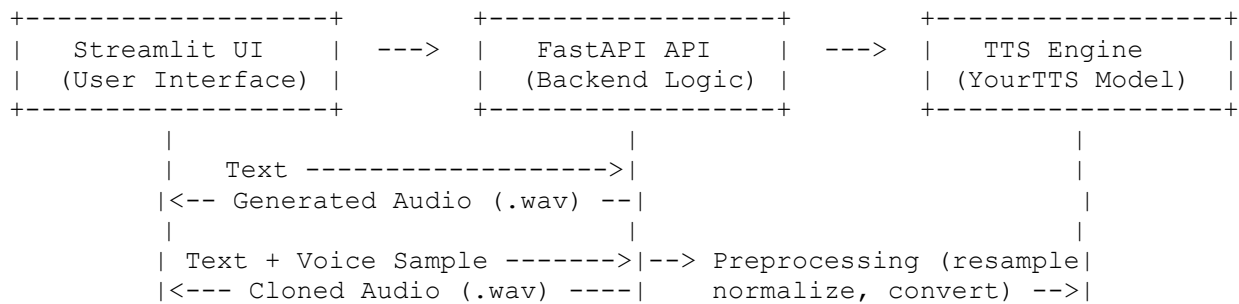
This makes the project **one-click runnable**. Users just double-click the batch file to start both backend and frontend services simultaneously.

Great — this section will tie everything together and show the **big picture** of how your project runs. Here’s the full write-up:

6. System Architecture

The system is designed in a **modular architecture** where each component is responsible for a specific task (TTS, voice cloning, preprocessing, frontend, backend). This modularity ensures that parts can be updated or replaced without affecting the whole system.

6.1 Workflow Diagram



6.2 Component Interaction

- **Frontend (Streamlit UI):**
 - User inputs text for TTS.
 - User uploads audio + text for voice cloning.
 - Sends requests to backend and plays audio response.
- **Backend (FastAPI):**
 - Exposes `/tts/` and `/clone/` endpoints.
 - Handles incoming requests, validates inputs.
 - Calls respective modules (`tts_engine.py`, `voice_cloning.py`).
- **TTS Engine (`tts_engine.py`):**
 - Loads YourTTS model once at startup.
 - Generates speech from text with selected speaker/language.
- **Voice Cloning (`voice_cloning.py`):**
 - Uses a reference audio sample.
 - Extracts **speaker embeddings**.
 - Generates cloned speech in the same voice.
- **Preprocessing (`preprocessing.py`):**
 - Converts input audio (`.m4a`, `.mp3`) into `.wav`.
 - Ensures correct format: 16kHz, mono, normalized.
- **Automation (`start_project.bat`):**

- Starts both backend and frontend in one click.

6.3 Implementation

- **Text-to-Speech Flow:**
 1. User enters text in UI.
 2. API `/tts/` is called.
 3. `tts_engine.py` generates speech.
 4. Resulting `.wav` is returned to UI for playback.
- **Voice Cloning Flow:**
 1. User uploads `.wav` sample + enters text.
 2. API `/clone/` is called.
 3. `preprocessing.py` prepares audio.
 4. `voice_cloning.py` generates cloned voice.
 5. Resulting `.wav` is returned to UI for playback.

6.4 Folder Structure

```
text_speech_project/
├── app.py                # FastAPI backend entry point
├── start_project.bat     # One-click starter script
├── requirements.txt      # Dependencies
├── config/
│   └── settings.py      # Config paths, model names
├── src/
│   ├── tts_engine.py    # Text-to-Speech module
│   ├── voice_cloning.py # Voice cloning module
│   ├── preprocessing.py # Audio preprocessing (resampling, conversion)
│   └── utils.py         # Helper utilities
├── ui/
│   └── app_ui.py        # Streamlit frontend
├── data/
│   ├── output/          # Generated speech files
│   └── samples/         # User-provided recordings
└── convert_m4a_to_wav.py # Standalone converter utility
```

6.5 Important Scripts

- **tts_engine.py**
 - Loads the YourTTS model.
 - Generates natural speech from input text.

- Supports multiple speakers and languages.
- **voice_cloning.py**
 - Uses reference audio to extract speaker embeddings.
 - Clones the voice to speak arbitrary text.
- **app.py**
 - Backend service using FastAPI.
 - Defines API endpoints `/tts/` and `/clone/`.
 - Calls TTS and cloning functions.
- **preprocessing.py**
 - Ensures all input audio is in required `.wav` format.
 - Resamples and normalizes recordings.
- **app_ui.py (Streamlit UI)**
 - Provides user-friendly web interface.
 - Text input box for TTS.
 - File upload + text for voice cloning.
 - Embedded audio players for playback.
- **start_project.bat**
 - Activates virtual environment.
 - Starts backend (FastAPI).
 - Starts frontend (Streamlit).
 - Simplifies usage into a one-click run.

7. API Endpoints

Our backend, built with **FastAPI**, exposes two main endpoints:

◆ 1. `/tts/`

Method: POST

Inputs:

- `text` → String (required)
- `speaker` → String (optional, e.g., "male-en-2")
- `language` → String (optional, e.g., "en")

Output:

- Returns generated audio as `.wav` file.

Usage Example (cURL):

```
curl -X POST "http://127.0.0.1:8000/tts/" -F "text=Hello, world!"
```

🔗 2. /clone/

Method: POST

Inputs:

- text → String (required)
- file → Audio file (.wav preferred, but .m4a supported after conversion)

Output:

- Returns cloned audio as .wav file.

Usage Example (cURL):

```
curl -X POST "http://127.0.0.1:8000/clone/" \  
-F "text=This is my cloned voice" \  
-F "file=@sample.wav"
```

8. Challenges & Solutions

Building this project required solving multiple technical challenges.

8.1 Python Version Issues

Problem:

- Coqui-TTS requires **Python ≥3.7, <3.11**.
- Initial environment was Python 3.12 → incompatible.

Solution:

- Installed **Python 3.10** specifically.
- Created virtual environment:
- `python3.10 -m venv venv310`
- `venv310\Scripts\activate`

8.2 Torch Installation Issues

Problem:

- Attempted installation of `torch==2.1.2`, but wheel not found for CUDA version.

Solution:

- Installed the **latest supported PyTorch** directly from the official index:
- `pip install torch torchvision torchaudio --index-url https://download.pytorch.org/whl/cu121`
- This ensured GPU acceleration compatibility.

8.3 Voice Sample Conversion (m4a → wav)

Problem:

- Many users record audio in .m4a, but Coqui-TTS only accepts **16kHz .wav mono**.
- Pydub raised errors because **FFmpeg** was missing.

Solution:

- Installed **FFmpeg** and integrated a preprocessing step.
- Provided a script:
- `from pydub import AudioSegment`
- `sound = AudioSegment.from_file("sample.m4a", format="m4a")`
- `sound.export("sample.wav", format="wav")`
- This conversion step was also included in `preprocessing.py`.

8.4 Backend–Frontend Connection

Problem:

- Running Streamlit frontend resulted in:
`requests.exceptions.ConnectionError: Failed to establish a new connection`

Cause:

- FastAPI backend wasn't running before the frontend started.

Solution:

- Always start backend first:
- `uvicorn app:app --reload`
- Later automated using **start_project.bat**, which runs both backend and frontend together.

9. Requirements

◆ Software & Libraries

- **Python 3.10**

- **PyTorch** (with CUDA if GPU available)
- **Coqui-TTS (0.22.0)**
- **FastAPI + Uvicorn**
- **Streamlit**
- **Librosa, Soundfile, Pydub** (for audio handling)
- **FFmpeg** (for .m4a → .wav conversion)
- **python-multipart** (for file uploads in FastAPI)

◆ Hardware

- **CPU:** Works fine but slower.
- **GPU (CUDA 12.1):** Significantly faster inference.
- **RAM:** Minimum 8GB recommended.

9.3 Libraries

The following Python libraries were used in this project:

- **Coqui-TTS (0.22.0)** → Core Text-to-Speech and Voice Cloning framework.
- **PyTorch (2.2.x)** → Deep learning backend for running TTS models.
- **Torchaudio** → Audio processing utilities integrated with PyTorch.
- **Librosa** → Audio analysis, spectrogram handling, and feature extraction.
- **Soundfile** → For reading and writing .wav files.
- **Pydub** → Audio conversion and manipulation (used for .m4a → .wav).
- **FFmpeg** → Required backend tool for handling audio formats.
- **FastAPI** → Backend framework for serving REST API endpoints.
- **Uvicorn** → ASGI server for running FastAPI.
- **Streamlit** → Frontend web interface for user interaction.
- **Requests** → Communication between frontend (Streamlit) and backend (FastAPI).
- **Python-Multipart** → Required for handling file uploads in FastAPI.

10. How to Run the Project

Step-by-Step Guide

1. **Clone or Download the Project**
2. `git clone <repo-url>`
3. `cd text_speech_project`
4. **Create and Activate Virtual Environment**
5. `python -m venv venv310`
6. `venv310\Scripts\activate` # Windows
7. **Install Dependencies**
8. `pip install -r requirements.txt`
9. **Ensure FFmpeg is Installed**

- Download from ffmpeg.org
- Add bin/ folder to **PATH** environment variable.

10. Start the Project (One Click)

- Run:
- start_project.bat
- This will:
 - ✓ Activate virtual environment
 - ✓ Start FastAPI backend on port 8000
 - ✓ Start Streamlit frontend on port 8501

11. Access the Application

- API Docs: <http://127.0.0.1:8000/docs>
- UI: <http://localhost:8501>

11. Results & Applications

11.1 TTS Results

- Input: "Hello, welcome to the text-to-speech system."
- Output: A .wav audio file generated in **natural human-like speech**.
- Supports multiple speakers (male, female) and languages (en, fr-fr, pt-br).

11.2 Voice Cloning Results

- Given a **5–10 sec voice sample** (in .wav), the system can:
 - Clone the voice style and tone.
 - Generate speech in the **user's voice** for any new text.

11.3 Real-World Applications

- **Assistive Technology** → Helping visually impaired individuals access written content.
- **Personalized Voice Assistants** → Custom AI assistants that sound like the user.
- **Audiobook Generation** → Convert e-books into speech with different narrators.
- **Content Creation** → YouTubers, podcasters, or educators can generate speech without recording.
- **Customer Service Bots** → More natural and human-like voice interactions.

12. Future Improvements

- **Custom Training for More Voices**
 - Allow users to fine-tune and train the model with longer datasets for higher accuracy.
- **Deploying on Cloud/Docker**
 - Package the system into a **Docker container** for deployment on AWS, Azure, or GCP.
 - Provide API access to multiple users simultaneously.
- **Packaging as Desktop App**

- Convert into an .exe (Windows) or .dmg (Mac) desktop application.
 - Enable offline usage with a simple GUI.
- **Support for More Languages**
 - Extend the model with multilingual datasets for global applications.
- **History and Playback in UI**
 - Save generated files in UI for easy re-access.

13. Conclusion

This project successfully integrates **Text-to-Speech (TTS)** and **Voice Cloning** into a unified framework. By combining **Coqui-TTS**, **FastAPI**, and **Streamlit**, it provides both:

- A **robust backend API** for text-to-speech and voice cloning.
- A **user-friendly frontend** for seamless interaction.

The system highlights key aspects of modern speech synthesis:

- ✓ Natural-sounding multilingual speech.
- ✓ Personalization through voice cloning.
- ✓ Easy local deployment with minimal setup.

By addressing challenges such as **Python version compatibility**, **Torch installation**, **audio format conversion**, and **backend–frontend integration**, the project demonstrates a complete end-to-end pipeline.

With future extensions (Docker deployment, desktop app packaging, and custom training), this system has the potential to be used in **assistive tech, education, media, and AI assistants**.