

Model Performance Report - Task 1: Fine-tuning TTS for English with a Focus on Technical Vocabulary

1. Introduction

This report presents an overview of the performance of a text-to-speech model built using the SpeechT5 architecture. It includes dataset descriptions, logs from the training and inference processes, evaluation results, and a glossary of technical terms.

2. Dataset Description

2.1 Dataset Overview

- **Name:** CMU Arctic X-Vectors
- **Source:** Matthijs/cmu-arctic-xvectors on Hugging Face Datasets
- **Size:** The dataset consists of several speaker embeddings for various speakers.
- **Features:** The dataset contains x-vectors, which are embeddings capturing the voice characteristics of speakers.
- **Target Variable:** Not applicable in this context, as the focus is on generating speech rather than classification.

2.2 Data Preprocessing

- **Embedding Selection:** A specific x-vector (speaker embedding) was selected from the validation split for the inference.
- **Input Processing:** The input text was tokenized using the SpeechT5 processor.

3. Logs

3.1 Inference Logs

- **Model:** SpeechT5ForTextToSpeech
- **Processor:** SpeechT5Processor
- **Vocoder:** SpeechT5HifiGan
- **Input Text:** "OOP is centered around four main concepts: encapsulation, inheritance, polymorphism, and abstraction..."
- **Selected X-Vector Index:** 7306

3.2 Runtime Environment

- **Framework:** Hugging Face Transformers
- **Hardware:** Assumed to be running on a standard CPU or GPU (specific details not provided in the code).

4. Evaluation Results

4.1 Output Generation

- **Generated Speech File:** speech.wav
- **Sample Rate:** 16,000 Hz
- **Audio Quality:** Subjective evaluation required for quality; recommended to listen for clarity and naturalness.

4.2 Performance Metrics

- **Model Accuracy:** Not applicable as this is a generative task rather than a classification.
- **Quality of Speech:** To be evaluated through human listening tests.

5. Conclusion

- **Summary of Findings:** The model successfully generated speech from the provided text using the specified speaker embedding. The output file can be evaluated for its quality.
- **Future Work:** Consider additional evaluation metrics such as MOS (Mean Opinion Score) for assessing speech quality and exploring further fine-tuning on a larger dataset.

6. Glossary of Technical Terms

Term	Definition
Text-to-Speech (TTS)	A technology that converts written text into spoken words
Processor	A component that prepares input data for the model
Model	A mathematical representation trained to perform specific tasks
Vocoder	A device that converts audio signals into a form that can be processed
X-Vector	A fixed-dimensional representation of voice characteristics
Inference	The process of generating output from a trained model