

Real-Time Audio Transcription with NVIDIA Parakeet RNNT 1.1B and Gradio

Introduction

In today's digital landscape, real-time speech recognition has become increasingly important across various applications - from virtual assistants and content creation tools to accessibility features and meeting transcription services. This tutorial demonstrates how to build a real-time audio transcription system using NVIDIA's Parakeet RNNT 1.1B model and create an interactive web interface with Gradio.

The Parakeet RNNT (Recurrent Neural Network Transducer) is a state-of-the-art speech recognition model that offers excellent accuracy while being efficient enough for real-time applications. Combined with Gradio's user-friendly interface development capabilities, we can create a powerful yet accessible transcription tool.

What We'll Build

By the end of this tutorial, you'll have:

- A real-time audio transcription system powered by NVIDIA's Parakeet RNNT 1.1B model
- A web interface built with Gradio that allows users to:
 - Transcribe audio from microphone input in real-time
 - Upload audio files for transcription
 - Download transcription results
 - Customize transcription settings

Prerequisites

- Python 3.8+
- GPU with CUDA support (recommended but not required)
- Basic understanding of Python and deep learning concepts

Setting Up Your Environment

First, let's set up a virtual environment and install the necessary packages:

```
bash
```

```
# Create and activate a virtual environment
```

```
python -m venv parakeet_env
```

```
source parakeet_env/bin/activate # On Windows: parakeet_env\Scripts\activate
```

```
# Install dependencies
```

```
pip install torch torchaudio gradio numpy soundfile
```

```
pip install git+https://github.com/NVIDIA/NeMo.git@main#egg=nemo_toolkit[asr]
```

Downloading the Parakeet RNNT 1.1B Model

NVIDIA's Parakeet RNNT 1.1B model is available through the NeMo framework. Let's create a script to download and prepare the model:

```
python
```

```
import os
```

```
import torch
```

```
from nemo.collections.asr.models import EncDecRNNTBPEModel
```

```
def download_model():
```

```
    # Create model directory if it doesn't exist
```

```
    os.makedirs("models", exist_ok=True)
```

```
    # Download and return the Parakeet RNNT 1.1B model
```

```
    model = EncDecRNNTBPEModel.from_pretrained(model_name="nvidia/parakeet-rnnt-1.1b")
```

```
    model.save_to("models/parakeet_rnnt_1.1b.nemo")
```

```
    return model
```

```
if __name__ == "__main__":
```

```
    model = download_model()
```

```
    print("Model downloaded successfully!")
```

Save this script as `download_model.py` and run it to download the model:

```
bash
```

```
python download_model.py
```

Creating the Transcription Engine

Now, let's create the core transcription functionality:

python

```
import torch
import numpy as np
from nemo.collections.asr.models import EncDecRNNTBPEModel

class ParakeetTranscriber:
    def __init__(self, model_path="models/parakeet_rnnt_1.1b.nemo"):
        # Load the model
        self.device = 'cuda' if torch.cuda.is_available() else 'cpu'
        print(f"Using device: {self.device}")

        # Load the model
        self.model = EncDecRNNTBPEModel.restore_from(model_path).to(self.device)
        self.model.eval()

        # Set default parameters
        self.sample_rate = 16000 # The model expects 16kHz audio

    def transcribe_file(self, audio_path):
        """Transcribe an audio file"""
        with torch.no_grad():
            transcription = self.model.transcribe([audio_path])
        return transcription[0]

    def transcribe_audio(self, audio_array, sample_rate):
        """Transcribe audio from a numpy array"""
        # Resample if necessary
        if sample_rate != self.sample_rate:
            # Implement resampling (using torchaudio or another library)
            # For simplicity, we'll assume the correct sample rate is provided
            pass

        with torch.no_grad():
            # Process the audio array
            audio_tensor = torch.tensor(audio_array).unsqueeze(0).to(self.device)

            # Transcribe
            transcription = self.model.transcribe(audio_tensor,
                                                  sampling_rate=sample_rate)

        return transcription[0]
```

Save this as `transcriber.py`.

Building the Gradio Interface

Now, let's create a user-friendly interface using Gradio:


```

import gradio as gr
import numpy as np
import torch
from transcriber import ParakeetTranscriber

# Initialize the transcriber
transcriber = ParakeetTranscriber()

def transcribe_file(audio_file):
    """Transcribe an uploaded audio file"""
    if audio_file is None:
        return "Please upload an audio file."

    transcription = transcriber.transcribe_file(audio_file)
    return transcription

def transcribe_microphone(audio, state=""):
    """Transcribe microphone input with streaming capability"""
    if audio is None:
        return state

    # Get the audio data and sample rate
    audio_data, sample_rate = audio

    # Transcribe the audio segment
    transcription = transcriber.transcribe_audio(audio_data, sample_rate)

    # Update and return the accumulated transcription
    state += " " + transcription
    return state.strip()

def reset_transcription():
    """Reset the transcription state"""
    return ""

# Create the Gradio interface
with gr.Blocks(title="Real-Time Audio Transcription") as demo:
    gr.Markdown("# 🎧 Real-Time Audio Transcription with NVIDIA Parakeet RNNT 1.1B")
    gr.Markdown("Transcribe audio in real-time using NVIDIA's Parakeet RNNT 1.1B model")

    with gr.Tabs():
        with gr.TabItem("Microphone Transcription"):
            with gr.Row():
                with gr.Column():
                    audio_input = gr.Audio(
                        sources=["microphone"],

```

```

        type="numpy",
        streaming=True,
        show_label=False
    )
    reset_btn = gr.Button("Reset Transcription")

    with gr.Column():
        text_output = gr.Textbox(
            label="Transcription",
            placeholder="Speak into your microphone...",
            lines=10
        )

    # Set up event handlers
    audio_input.stream(
        fn=transcribe_microphone,
        inputs=[audio_input, text_output],
        outputs=text_output
    )
    reset_btn.click(fn=reset_transcription, outputs=text_output)

    with gr.TabItem("File Transcription"):
        with gr.Row():
            with gr.Column():
                file_input = gr.Audio(type="filepath", label="Upload Audio File")
                transcribe_btn = gr.Button("Transcribe")

            with gr.Column():
                file_output = gr.Textbox(
                    label="Transcription",
                    placeholder="Upload an audio file and click 'Transcribe'",
                    lines=10
                )

        # Set up event handler
        transcribe_btn.click(
            fn=transcribe_file,
            inputs=file_input,
            outputs=file_output
        )

    gr.Markdown("#### Notes")
    gr.Markdown("- For best results, use clear audio with minimal background noise")
    gr.Markdown("- Supported audio formats: WAV, MP3, OGG, FLAC")
    gr.Markdown("- The microphone transcription works in real-time with streaming enabled")

# Launch the app

```

```
if __name__ == "__main__":  
    demo.launch()
```

Save this as `app.py`.


Advanced Features

Adding Punctuation and Capitalization

To improve the readability of our transcriptions, we can add punctuation and proper capitalization:

python

```
from transformers import pipeline  
  
class TextEnhancer:  
    def __init__(self):  
        # Load the punctuation and capitalization model  
        self.punctuator = pipeline("text2text-generation",  
                                   model="oliverguhr/fullstop-punctuation-multilang-large")  
  
    def enhance(self, text):  
        """Add punctuation and proper capitalization to text"""  
        if not text.strip():  
            return text  
  
        enhanced = self.punctuator(text)[0]['generated_text']  
        return enhanced  
  
# Then in the transcribe functions, add:  
enhancer = TextEnhancer()  
transcription = enhancer.enhance(transcription)
```



Supporting Multiple Languages

We can extend our application to support multiple languages:

python

```
# Update the Gradio interface to include language selection
with gr.Blocks(title="Multilingual Audio Transcription") as demo:
    # ... previous code ...

    language_selector = gr.Dropdown(
        choices=["English", "Spanish", "French", "German", "Chinese"],
        value="English",
        label="Language"
    )

    # Update the transcription functions to use the selected language
    # ...
```

Performance Optimization

For better real-time performance, consider the following optimizations:

1. **Batch Processing:** Process audio in batches for better GPU utilization
2. **Model Quantization:** Reduce model size using quantization techniques
3. **Streaming Buffer Management:** Optimize how audio buffers are processed

Example code for model quantization:

python

```
def load_quantized_model(model_path):
    model = EncDecRNNTBPEModel.restore_from(model_path)
    quantized_model = torch.quantization.quantize_dynamic(
        model, {torch.nn.Linear}, dtype=torch.qint8
    )
    return quantized_model
```

Deployment Considerations

When deploying your application:

1. **Hardware Requirements:** Ensure sufficient GPU memory (4GB+ recommended)
2. **Server Setup:** Consider using NGINX as a reverse proxy
3. **Scaling:** Implement a queue system for handling multiple requests
4. **Security:** Add authentication for production deployments

Example Docker setup:

dockerfile

```
FROM python:3.9-slim
```

```
WORKDIR /app
```

```
# Install system dependencies
```

```
RUN apt-get update && apt-get install -y \  
    build-essential \  
    git \  
    && rm -rf /var/lib/apt/lists/*
```

```
# Copy requirements and install
```

```
COPY requirements.txt .
```

```
RUN pip install --no-cache-dir -r requirements.txt
```

```
# Copy application files
```

```
COPY . .
```

```
# Expose port
```

```
EXPOSE 7860
```

```
# Command to run the application
```

```
CMD ["python", "app.py"]
```

Conclusion

In this tutorial, we've built a powerful real-time audio transcription system using NVIDIA's Parakeet RNNT 1.1B model and created an intuitive interface with Gradio. This combination provides both the accuracy of a state-of-the-art speech recognition model and the accessibility of a user-friendly web interface.

The applications for this technology are vast - from creating meeting transcription tools and closed captioning systems to building voice-controlled applications and accessibility aids. As speech recognition technology continues to advance, we can expect even more accurate and efficient models to emerge.

Further Reading and Resources

- [NVIDIA NeMo Framework Documentation](#)
- [Gradio Documentation](#)
- [Speech Recognition Fundamentals](#)
- [RNNT Architecture Explained](#)

About the Author

[Your Name] is a [Your Profession] specializing in [Your Specialization]. With a background in [Your Background], [Your Name] is passionate about making advanced AI technologies accessible to developers of all skill levels.

#ASR #SpeechRecognition #NVIDIA #Gradio #DeepLearning #AudioProcessing #MachineLearning
#AITools