# Voice Cloning with Tacotron2

Creating a voice cloning system using deep learning involves several steps, including data preparation, model training, and synthesis. Here, I'll provide a simplified Python code outline using the Tacotron2 and WaveGlow models from NVIDIA, which are commonly used for text-to-speech (TTS) applications.

## Prerequisites

1. **Install Dependencies**: Make sure you have Python installed, along with PyTorch, and NVIDIA's Tacotron2 and WaveGlow repositories.

```bash
pip install torch
pip install numpy scipy
git clone https://github.com/NVIDIA/tacotron2.git
cd tacotron2
pip install -r requirements.txt
git clone https://github.com/NVIDIA/waveglow.git
cd waveglow
pip install -r requirements.txt
```

## Step 1: Data Preparation

You need a dataset consisting of audio files and their corresponding transcripts. For simplicity, we'll assume you have a dataset in the format required by Tacotron2.

## Step 2: Preprocess the Data

Preprocessing involves converting audio files into mel-spectrograms and normalizing the transcripts.

```python
# Preprocessing script (simplified version)
import os
from tacotron2.data_function import TextMelLoader, TextMelCollate
from torch.utils.data import DataLoader

# Set the paths
dataset_path = "path_to_your_dataset"
filelist = "filelist.txt"  # This file should contain lines of "<audio_path>|<transcript>"

# Create data loader
trainset = TextMelLoader(filelist, dataset_path)
collate_fn = TextMelCollate(n_frames_per_step=1)
train_loader = DataLoader(trainset, num_workers=1, shuffle=True,
                          sampler=None, batch_size=32, pin_memory=False,
                          drop_last=True, collate_fn=collate_fn)
```

## Step 3: Train Tacotron2 Model

Train the Tacotron2 model on the preprocessed data.

```python
import torch
from tacotron2.model import Tacotron2
from tacotron2.train import train

# Initialize the model
model = Tacotron2()

# Load the dataset
train_loader = ...

# Train the model
train(model, train_loader)
```

## Step 4: Synthesize Speech with Tacotron2 and WaveGlow

Once the Tacotron2 model is trained, use it to generate mel-spectrograms from text, and then use WaveGlow to convert mel-spectrograms into audio.

```python
import torch
from tacotron2.model import Tacotron2
from tacotron2.text import text_to_sequence
from waveglow.denoiser import Denoiser

# Load Tacotron2 and WaveGlow models
tacotron2 = Tacotron2()
waveglow = torch.load('path_to_waveglow_model')
denoiser = Denoiser(waveglow)

# Function to generate speech
def generate_speech(text, tacotron2, waveglow, denoiser):
    sequence = np.array(text_to_sequence(text, ['english_cleaners']))[None, :]
    sequence = torch.from_numpy(sequence).to(device='cuda', dtype=torch.long)

    with torch.no_grad():
        mel_outputs, mel_outputs_postnet, _, alignments = tacotron2.inference(sequence)
        audio = waveglow.infer(mel_outputs_postnet, sigma=0.666)
        audio = denoiser(audio, strength=0.01)[:, 0]

    return audio.cpu().numpy()

# Generate speech
text = "Hello, how are you?"
audio = generate_speech(text, tacotron2, waveglow, denoiser)

# Save to file
import soundfile as sf
sf.write('output.wav', audio, 22050)
```

## Notes

1. **Data Quality**: The quality of the cloned voice depends heavily on the quality and quantity of your training data.

2. **Training Time**: Training a Tacotron2 model from scratch can take a significant amount of time and computational resources. Consider using pre-trained models if available.
3. **Fine-Tuning**: If you have a pre-trained model, you can fine-tune it on your specific dataset to adapt it to your voice cloning task.

This outline provides a basic framework. For a production-level system, you'd need to handle more aspects such as model checkpoints, data augmentation, error handling, and more robust preprocessing.