

**Model Type and Task:** Kokoro-82M is a **Text-to-Speech (TTS)** model designed to convert written text into spoken words. This task involves generating human-like speech audio from textual input, making it a **sequence-to-sequence generation** problem where a sequence of text tokens is transformed into a sequence of audio frames.

**Use Cases:** Potential applications for Kokoro-82M include:

- **Virtual Assistants:** Providing natural and expressive voices for AI assistants.
- **Audiobook Narration:** Converting written books into spoken form for auditory consumption.
- **Accessibility Tools:** Assisting visually impaired users by reading out on-screen text.
- **Language Learning:** Offering pronunciation examples for language learners.
- **Content Creation:** Generating voiceovers for videos and interactive media.

**Evaluation Metrics:** Assessing the performance of TTS models like Kokoro-82M involves both objective and subjective measures:

- **Mean Opinion Score (MOS):** A subjective metric where human listeners rate the naturalness and quality of the synthesized speech on a scale, typically from 1 to 5.
- **Word Error Rate (WER):** An objective measure that calculates the accuracy of the generated speech by comparing the synthesized audio to the original text, focusing on the correctness of word pronunciation.
- **Inference Speed:** Evaluates how quickly the model can generate speech, which is crucial for real-time applications.

### **Kokoro-82M Notes:**

- **Model Overview:**
  - Kokoro-82M is an open-weight Text-to-Speech (TTS) model developed by hexgrad.
  - It contains 82 million parameters, offering a balance between model size and performance.
  - Despite its compact size, it delivers quality comparable to larger TTS models.
  - The model is licensed under Apache 2.0, allowing for flexible deployment in various environments.
- **Performance and Efficiency:**
  - Kokoro-82M is noted for being significantly faster and more cost-efficient than larger models.
  - Its lightweight architecture enables deployment in production settings and personal projects without heavy computational requirements.
- **Releases and Updates:**
  - As of January 2, 2025, version 0.19 weights were released in full FP32 precision under the Apache 2.0 license.
  - Ten unique voice packs have been made available to enhance versatility.
  - An ONNX version of v0.19 is also accessible, facilitating integration into various platforms.

- **Language Support:**
  - The model supports multiple languages, including:
    - English
    - Chinese (partial support)
    - Japanese (partial support)
    - German (partial support)
- **Community and Development:**
  - A Rust-based implementation, "Kokoros," offers fast inference capabilities, enabling integration into various applications.
  - The Rust version supports streaming mode and style mixing, allowing for diverse voice outputs.
  - A Discord community has been established for collaboration and support among users and developers.
- **Deployment and Accessibility:**
  - A no-code demo is available on Hugging Face Spaces, allowing users to experience the model's capabilities without setup.
  - The model can be installed via pip, simplifying integration into Python projects.
  - DeepInfra provides deployment options with pricing set at \$0.80 per million characters, making it accessible for various applications.
- **Additional Resources:**
  - Sample outputs, evaluation metrics, and voice details are documented in the model's repository.
  - The community actively contributes to the model's development, ensuring continuous improvements and feature additions.

The model saves its audio files in .wav –

## Why .wav?

1. **Uncompressed Quality:** WAV files store raw, uncompressed audio, which preserves the highest quality without data loss. This is especially useful in TTS models where clarity and fidelity are important.
2. **Standard Format for TTS:** Many Text-to-Speech models generate .wav because it's widely supported across platforms and easy to process further (e.g., for editing or converting to other formats).
3. **Consistent Sample Rate & Bit Depth:** Unlike MP3, which is a compressed format with varying quality, WAV files ensure that the generated speech maintains a consistent sample rate and bit depth, which is crucial for further processing, analysis, or training other models.

To assess the model, I experimented with:

- **Different languages:** English, Japanese, and Chinese
- **Varying sentence structures:** Short and long sentences, humorous and serious tones
- **Pronunciation and clarity:** Checking phonetic accuracy across different inputs

I ran multiple inputs through the model and recorded its outputs to analyze consistency and quality.

## Findings

- The model performs very well in **English**, and **Chinese** speech felt notably more fluid and natural compared to **Japanese**.
- In **Japanese**, the output was intelligible but had a **robotic, overly enunciated quality**, almost like listening to a historical play rather than natural speech.
- **Longer sentences** sometimes led to minor pronunciation inconsistencies, but overall, the speech quality remained stable.
- The **model is deterministic**, meaning the same input always produces the same output.

## Example Results

The results(audio files) are available in the gen\_audio folder in GitHub.

## Model Provenance

- **Trained by:** The Kokoro-82M model was developed and trained by hexgrad.
- **Purpose:** Designed as an open-weight, efficient Text-to-Speech (TTS) model, Kokoro-82M aims to deliver high-quality speech synthesis while maintaining a lightweight architecture.
- **Data Used:** The model was trained on less than 100 hours of audio data over fewer than 20 epochs. Specific datasets have not been disclosed, but the concise training regimen highlights the model's efficiency.
- **Licensing:** Released under the Apache 2.0 license, Kokoro-82M is accessible for both research and commercial applications, ensuring flexibility in deployment.

## Conclusion

Kokoro-82M is a surprisingly efficient TTS model, delivering high-quality speech generation while remaining lightweight. However, improvements could be made in handling tonal languages like Chinese and processing longer inputs more smoothly. Overall, it's a great option for fast and cost-effective TTS applications.