

Something Cool with TTS



Machine Learning Singapore #MSLG

- Martin Andrews = GeminiTTS @ mda.net

13-May-2026

About Me

- Machine Intelligence / Startups / Finance
 - Moved from NYC to Singapore in Sep-2013
- 2014 = 'fun' :
 - Machine Learning, Deep Learning, NLP
 - Robots, drones
- Since 2015 = 'serious' :: NLP + deep learning
 - Including Papers...
 - & GDE ML; ML-Singapore co-organiser...
 - & Red Dragon AI...



Outline

- Gemini-TTS
 - Announcement
 - AI Studio playtime
 - Coding in Colab
 - Resources
- How are TTS models trained?
 - Example Open Source model
 - Voice Cloning Colab
- Wrap-up & QR-code

Gemini TTS

Google Official Announcement

(15-April-2026)

- "Gemini 3.1 Flash TTS: the next generation of expressive AI speech"
 - Improved speech quality and controllability
 - New audio tags for more expressive speech generation
 - Built for global scale
 - Watermarked with SynthID
- Gemini API Speech Generation Documentation



Gemini 3.1 Flash TTS



Audio Profile

Describe the voice persona (e.g., warm, energetic, professional)

Director's note

Style

Pace

Accent

↑↓ Style

🕒 Pace

🌐 Accent

Voice

🔍 Search voices

Achernar

Soft Higher pitch

Achird

Friendly Lower middle pitch

Algenib

Gravelly Lower pitch

Alaieba

Gemini-TTS in AI Studio

- AI Studio Link
 - sadly, don't seem to be able to link to pre-populated example
- Check out:
 - Voice samples
 - Direction via RHS panel
 - REST samples

Gemini-TTS Colab DEMO

- Gemini-TTS Colab Link
- DEMO
 - Generate Speech:
 - DJ in London
 - Whisper
 - Analyse voice:
 - Rewrite prompt for Singapore DJ + Play
 - Play existing audio + Analyse (3 times)
 - Synthesise by approximation...
 - Scottish version
 - Iterate to match it more closely x4
 - finally listen to genuine article



A screenshot of a web-based interface for Gemini-TTS. The top bar shows the Gemini logo, the title "Gemini-TTS_playground", and navigation icons. Below the title is a menu with "File", "Edit", "View", "Insert", "Runtime", "Tools", and "Help". A search bar contains "Commands" and buttons for "+ Code", "+ Text", and "Run all". The main area is split into two panels. The top panel shows a status message "Generating audio... (Attempt 1)" and "Audio generated successfully and saved to Iapetu". Below this is a media player with a play button, a progress bar at "0:00 / 0:13", and volume controls. The bottom panel shows a code editor with Python code. The code includes a comment "# Adjust the Performance Notes, etc, to match th", a function definition "def full_cycle_of_matching(previous_context, pre", and a call to "matching_analysis = run_prompts_and_files([". The code also contains a docstring with instructions: "The following audio file is a recording of Max, Please listen to Max, to understand the qualitie", "audio_reference,", and "The following is the speaker delivery by 'Sam W'". The docstring continues with "The instructions they were given were:", "--- Previous instructions start ---", "{previous_context.strip()}", "--- Previous instructions end ---", "The transcript will remain unchanged.", and "match_audio_fname(previous_iter),". The code ends with "Please rewrite the previous instructions for 'Sa", "* the Audio Profile by-line;", "* the Scene;", and "* the Performance notes for style. pace and acce".

The Nine Rules

1. Always prepend a synthesize-speech preamble

This single paragraph is what reliably triggers the speech-synthesis classifier path instead of the "read it all" path.

```
1 Synthesize speech for the performance defined below. The profile, scene,  
2 performance notes, and context are direction only. Do NOT speak them.  
3 Speak ONLY the lines under ##### TRANSCRIPT.
```

Skip this and you'll hit intermittent failures where Gemini reads your whole prompt out loud.

2. Use ##### TRANSCRIPT as the delimiter, exact header

We tested alternatives. Lower-level headers (##### TRANSCRIPT) work, but ##### TRANSCRIPT is what the official docs example uses, and it's what we saw the most reliable classifier behavior with.

3. Use short section labels, not the docs' verbose ones

The current iteration of the official Gemini docs use ### DIRECTOR'S NOTES and ### SAMPLE CONTEXT . **Don't**. In our testing, the literal string DIRECTOR'S got spoken aloud.

Use these instead:

- ### PERFORMANCE (for Style / Pace / Accent)
- ### CONTEXT (for sample context)

Apostrophes and multi-word section headers are classifier hazards.

4. Classify emotional register BEFORE picking audio tags

Don't use a universal tag template. Pick one register first.

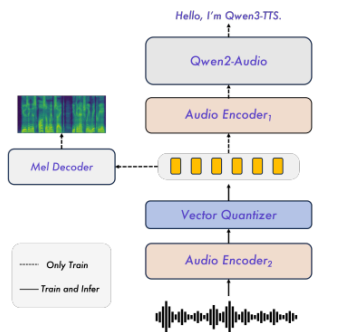
Other Resources

- [Simon Willison Blog](#)
 - includes a vibe-coded App
- [LiveKit Gemini TTS prompting guide](#)
 - actual real-world advice
 - ... incorporated into my Colab DEMO

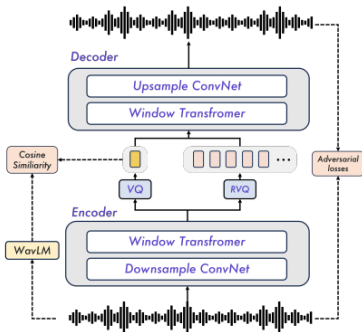
TTS Model Training

Architectural Highlights

- Qwen-TTS-Tokenzier-12Hz
 - multi-codebook tokenizer
 - jointly optimized semantic and acoustic streams
 - high-level semantic content
 - acoustic detail, prosody, and others

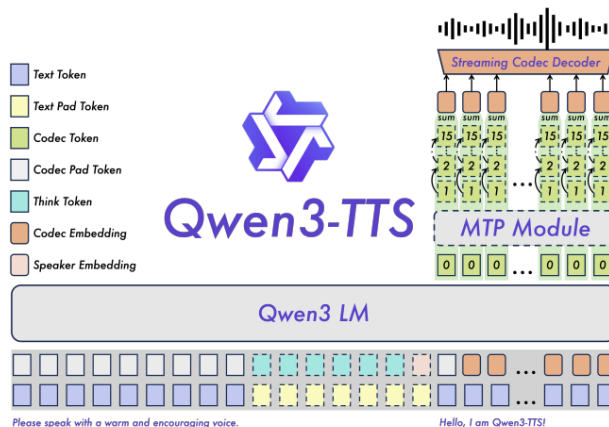


(a) Qwen-TTS-Tokenzier-25Hz



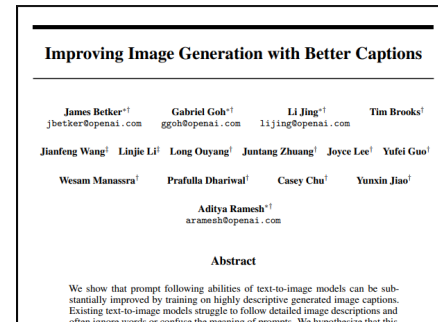
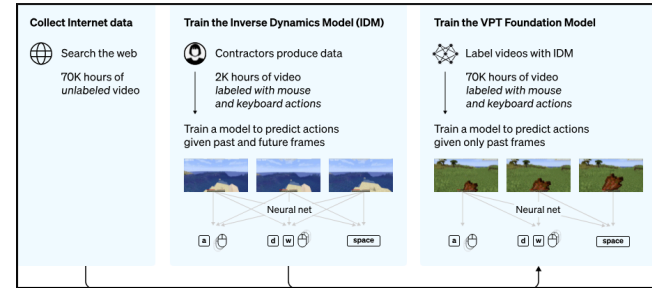
(b) Qwen-TTS-Tokenzier-12Hz

- LLM with MTP for audio tokens
 - LLM carries text
 - + (optional) speaker embeddings
 - MTP generates 16 audio tokens per step
 - audio frame rate much higher than LLM



How to train 'performance notes' ?

- Train voice for self-reconstruction
 - can transfer voice characteristics
 - separate from Transcript
 - measure voice similarity between I/O
- Train a model to *generate* 'performance notes'
 - ~same idea as :
 - Learning to play Minecraft with Video PreTraining (2022)
 - Improving Image Generation with Better Captions (2023)
- Can then create synthetic data at scale
 - still need examples of:
 - same-speaker, different-setting
 - same-speaker, different-emotion



Qwen Colab DEMO

- Actual Voice-Cloning Colab
- DEMO
 - Overall set-up
 - GUI section skipped
 - In PyTorch:
 - ingestion of voice for cloning
 - generation of cloned voice

A screenshot of a Google Colab notebook interface. The top bar shows the notebook title "Qwen3-TTS_voice-cloning" with a star and refresh icon, and a menu with "File", "Edit", "View", "Insert", "Runtime", "Tools", and "Help". Below the menu is a search bar for "Commands" and buttons for "+ Code", "+ Text", and "Run all". The notebook content is displayed in a dark theme. A cell titled "Actually, let's just try a voice clone..." contains a list item: "See : [https://github.com/QwenLM/Qwen3-TTS#voice-clone](\"https://github.com/QwenLM/Qwen3-TTS#voice-clone\")". Below this, two code cells are visible. The first code cell contains Python code for importing torch and soundfile, and loading the Qwen3-TTS model. The second code cell contains code for setting reference audio and text, and creating a voice clone prompt.

```
import torch
import soundfile as sf
from qwen_tts import Qwen3TTSModel

model = Qwen3TTSModel.from_pretrained(
    # "Qwen/Qwen3-TTS-12Hz-1.7B-Base",
    # ". /models/Qwen3-TTS-12Hz-1.7B-Base",
    device_map="cuda:0",
    dtype=torch.bfloat16,
    #attn_implementation="flash_attention_2", # pe
)
```

```
ref_audio = "https://redcatlabs.com/audio/2026-05-13"
ref_text = ""
Okay, so Google has just dropped Gemma 4.
And this is four new models with multimodality, thin
And honestly, that alone would get me covering this.
"".replace("\n", " ")

prompt_items = model.create_voice_clone_prompt(
    ref_audio=ref_audio,
```

Wrap-Up

- TTS models are getting better very quickly
 - realism is now table-stakes
 - controllability is improving
- Other factors : Latency, streaming, ...
 - the right combo still difficult to find
- Lots of potential applications just opened up!