TL;DR MultiTalk extends Wan 2.1 with label-aware audio injection, partial fine-tuning, and multi-task training so you can drive two or more performers from separate speech tracks, keep prompts responsive, and render 15-second clips at 480p/720p.

#### What is MeiGen MultiTalk?

MeiGen's MultiTalk is a research and open-source framework for generating multi-person conversational video that stays lip-synced to multi-stream audio. The team builds on Wan 2.1 I2V-14B, injects audio labels through a novel Label Rotary Position Embedding (L-RoPE), and keeps the base model's instruction-following intact via selective, partial parameter fine-tuning. MultiTalk can animate humans, stylised avatars, and cartoon characters and supports both short clips and 15-second streaming segments.

#### Links:

- Paper (arXiv): https://arxiv.org/abs/2505.22647
- GitHub: https://github.com/MeiGen-Al/MultiTalk
- · Project page: https://meigen-ai.github.io/multi-talk/
- Launch post: https://www.linkedin.com/posts/arminas-valunas-b4477255\_meigen-just-introduced-multitalk-a-new-ugcPost-7342832157330976769-zBX1

### Why it matters for production teams

- Orchestrate conversational explainers or interviews without filming on set, binding the right voice track to the right digital actor.
- Rapidly localise dialogue: swap out speech tracks (real or TTS) while preserving interaction prompts and body language.
- Mix humans and stylised avatars in the same shot—handy for brand mascots, product walk-throughs, and hybrid live-action/animated content.
- Scale creative testing with minimal compute: TeaCache, LoRA accelerators, and INT8 options bring render times and VRAM down to RTX 4090 class hardware.

#### Core innovations

- Label Rotary Position Embedding (L-RoPE): tags each audio stream so the diffusion backbone knows which character to drive, stopping cross-talk and mismatched lip-sync.
- Partial parameter training: fine-tunes a subset of layers to retain Wan's prompt-following while specialising for audio-person binding.
- **Multi-task curriculum:** stages training over talking-head, talking-body, and multi-person data to balance fidelity, instruction following, and motion diversity.
- Time-aligned conditioning stacks: aligns audio features, reference frames, and text prompts per denoising step for tighter conversational timing.

#### Generation modes and controls

- --mode streaming stitches multiple 81-frame chunks for up to ~15 seconds; --mode clip keeps to a single chunk for punchy responses.
- Guidance scales: keep audio CFG around 3-5 for lip-sync, and adjust text guidance to bias toward prompt style or spontaneity.
- TeaCache caching and SageAttention 2.2 support yield ~2-3x faster sampling; combine with LoRA distillations (FusioniX, lightx2v) for 4-8 step renders.
- JSON scene files pair prompts, image references, and multi-channel audio; add { "audio\_mode": "tts" } to integrate Kokoro-82M or other TTS tracks.

## Setup cheatsheet

```
Environment (CUDA 12.4 wheel shown):
conda create -n multitalk python=3.10
conda activate multitalk
pip install torch==2.5.1 torchvision==0.20.1 torchaudio==2.5.1 \
  --index-url https://download.pytorch.org/whl/cu124
pip install flash_attn==2.7.4.post1 psutil packaging
pip install -r requirements.txt
conda install -c conda-forge ffmpeg librosa
Model prep:
huggingface-cli download Wan-AI/Wan2.1-I2V-14B-480P --local-dir ./weights/Wan2.1-I2V-14B-480P
huggingface-cli download TencentGameMate/chinese-wav2vec2-base --local-dir ./weights/chinese-wav2vec2-base
huggingface-cli download hexgrad/Kokoro-82M --local-dir ./weights/Kokoro-82M
huggingface-cli download MeiGen-AI/MeiGen-MultiTalk --local-dir ./weights/MeiGen-MultiTalk
ln -sf $(pwd)/weights/MeiGen-MultiTalk/multitalk.safetensors \
 weights/Wan2.1-I2V-14B-480P/
Single-GPU multi-person inference:
python generate_multitalk.py \
  --ckpt_dir weights/Wan2.1-I2V-14B-480P \setminus
  --wav2vec_dir weights/chinese-wav2vec2-base \
  --input_json examples/multitalk_example_2.json \
 --sample_steps 40 \
  --mode streaming \
  --use_teacache \
  --save_file multi_long_exp
```

#### **Benchmarks and caveats**

- Trained on 81-frame sequences at 25 FPS; longer clips may drift unless you restitch with streaming mode or future checkpoints.
- 480p runs on a single 24-32 GB GPU; 720p currently expects multi-GPU with FSDP + sequence parallelism.
- TeaCache accelerates inference but may introduce slight colour shifts beyond threshold 0.5—keep between 0.2 and 0.5 for broadcast work.
- Instruction following benefits from partial tuning, yet extreme prompts can still override audio semantics; audit outputs when scripting complex choreography.

# **Ecosystem roadmap**

- July 2025 updates add INT8 quantisation (Optimum-Quanto) and SageAttention 2.2, shrinking VRAM and enabling faster CFG sweeps.
- TeaCache + APG support unlocks low-latency preview loops for creative reviews; expect official 720p scripts soon.
- August 2025's InfiniteTalk sibling targets infinite-length dubbing—keep an eye on cross-compatibility for longer scripted series.
- Community examples already splice singers, news anchors, and anime avatars; contribute scenes via the GitHub discussions to speed fine-tuning recipes.

Need help standing up multi-speaker virtual talent or integrating MultiTalk into your workflow?

### References

- MultiTalk (GitHub)
- MultiTalk (arXiv)
- InfiniteTalk (GitHub)
- InfiniteTalk (arXiv)