

 **Kedreamix** Merge pull request #135 from SummerRiver3D/main · · ·

ASR

ChatTTS @ f6b530b

CosyVoice @ 6be8d0f

GPT\_SoVITS

LLM

Musetalk

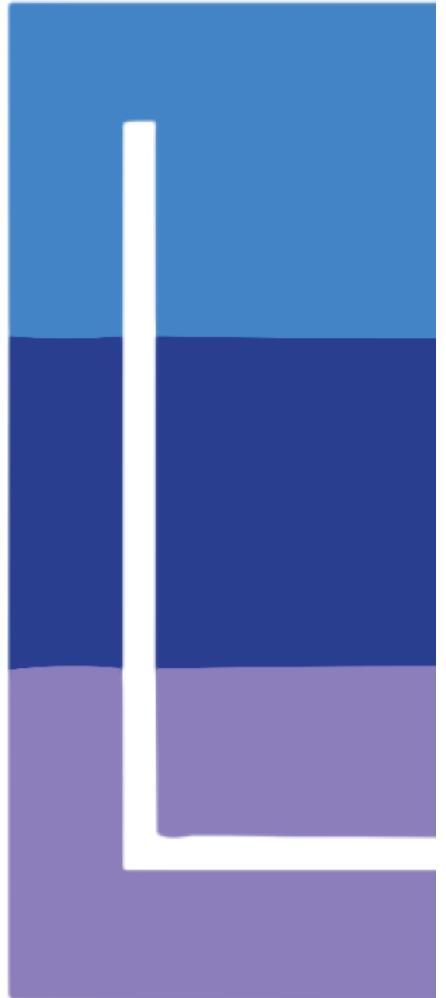
NeRF

TFG

TTS

VITS

## Digital Human Intelligent Dialogue System - Linly-Talker — 'Interac



## 2023.12 Update May 17

Users can upload any images for the conversation

## 2024.01 Update May 17

- Exciting news! I've now incorporated both the powerful GeminiPro and Qwen large models into our conversational scene. Users can now benefit from the advanced features of these models.
- The deployment invocation method for FastAPI has been updated.
- The advanced settings options for Microsoft TTS have been updated, increasing the variety of voice types. Additionally, video subtitle generation has been improved.
- Updated the GPT multi-turn conversation system to establish contextual connections in dialogue, enhancing the interactivity and fluency of conversations.

## 2024.02 Update May 17

- Updated Gradio to the latest version 4.16.0, providing the interface with additional functionalities such as capturing images from a camera and saving them directly to the application.
- ASR and THG have been updated. FunASR from Alibaba has been integrated into ASR, enhancing its speed significantly. Additionally, a new ASR model has been added.
- I have incorporated the GPT-SoVITS model, which is a voice cloning method. By fine-tuning it with just one minute of a person's speech, I can generate voices that sound like them.
- I have integrated a web user interface (WebUI) that allows for better execution of Linly-Talker.

## 2024.04 Update May 17

- Updated the offline mode for Paddle TTS, excluding Edge TTS.
- Updated ER-NeRF as one of the choices for Avatar generation.
- Updated app\_talk.py to allow for the free upload of voice and images/videos for generation without being based on a dialogue scene.

2024.05 Update 

- Updated the beginner-friendly AutoDL deployment tutorial, and also updated the codewithgpu image, allowing for one-click experience.
- Updated WebUI.py: Linly-Talker WebUI now supports multiple modules, multiple models, and multiple options

2024.06 Update 

- Integrated MuseTalk into Linly-Talker and updated the WebUI, enabling basic real-time conversation capabilities.
- The refined WebUI defaults to not loading the LLM model to reduce GPU memory usage. It directly responds with text to complete the conversation.

2024.08 Update 

- Updated CosyVoice to offer high-quality text-to-speech (TTS) functionality and voice cloning capabilities; also upgraded to Wav2Lip for better lip sync.

2024.09 Update 

- Added Linly-Talker API documentation, providing detailed interface descriptions to help users access Linly-Talker's features via the API.

2024.12 Update 

- Implemented a simple fix for the Edge-TTS bug, resolved several issues with MuseTalk, and plan to integrate fishTTS for more stable and natural-sounding responses.

2025.02 Update 

- Added OmniSenseVoice Model for Faster Speech Recognition

► Content

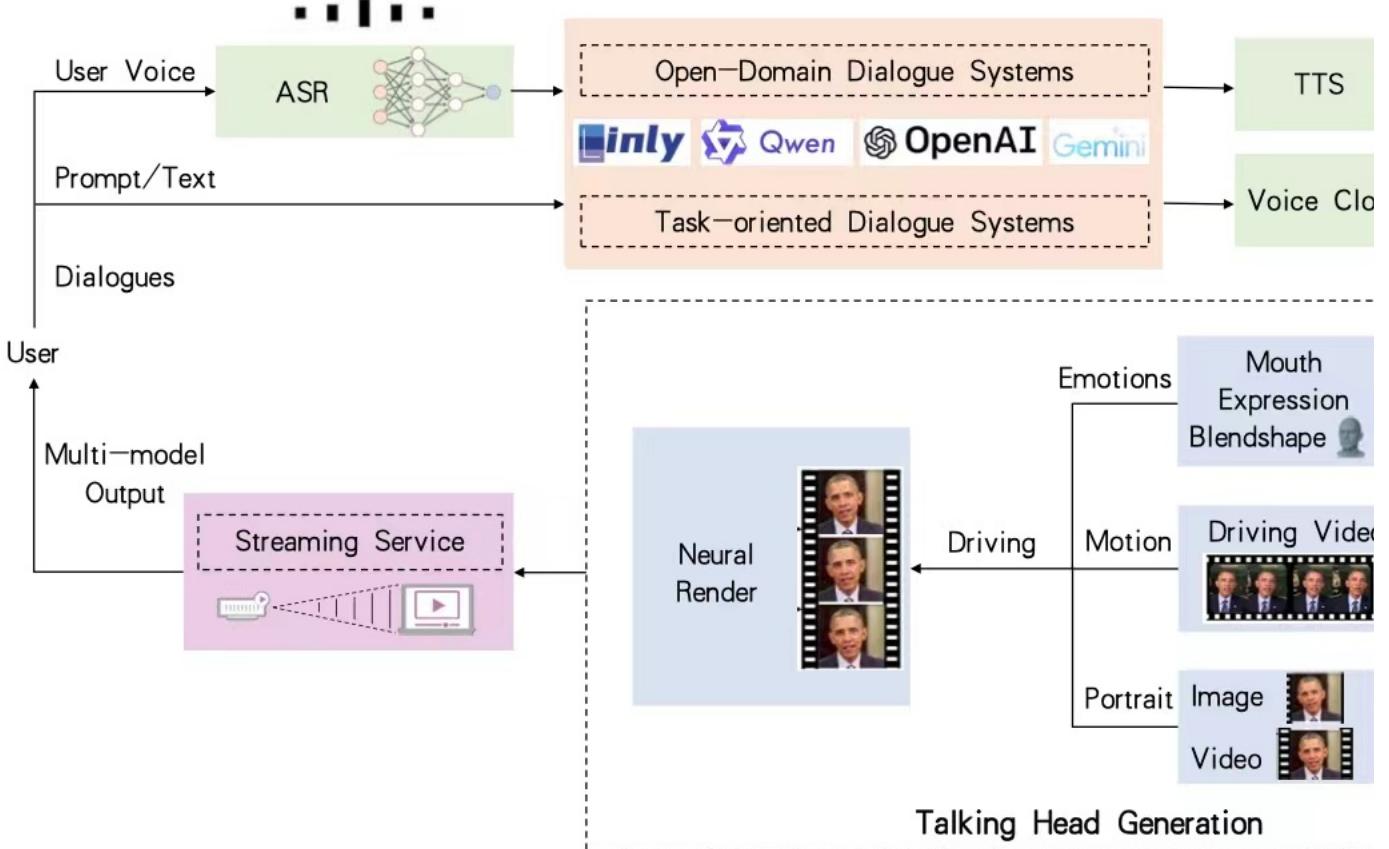
## Introduction

Linly-Talker is an innovative digital human conversation system that integrates the latest artificial intelligence technologies, including Large Language Models (LLMs) and various AI models for speech recognition, synthesis, and video generation.

The core features of the system include:

1. **Multi-Model Integration:** Linly-Talker combines major models such as Linly, GeminiPro, Qwen, as well as visual models like Whisper, SadTela, and Dalle.
2. **Multi-Turn Conversational Ability:** Through the multi-turn dialogue system powered by GPT models, Linly-Talker can understand and maintain context across multiple turns of conversation.
3. **Voice Cloning:** Utilizing technologies like GPT-SoVITS, users can upload a one-minute voice sample for fine-tuning, and the system will clone it to create a realistic digital twin.
4. **Real-Time Interaction:** The system supports real-time speech recognition and video captioning, allowing users to communicate naturally and interactively.
5. **Visual Enhancement:** With digital human generation technologies, Linly-Talker can create realistic digital human avatars, providing a more engaging and personalized interaction experience.

The design philosophy of Linly-Talker is to create a new form of human-computer interaction that goes beyond simple Q&A. By integrating advanced AI technologies, the system aims to provide a more natural, intuitive, and personalized conversational experience.



### ① Note

You can watch the demo video [here](#).

I have recorded a series of videos on Bilibili, which also represent every step of my updates and methods of use. For detailed information, pl

- 🔥🔥 Digital Human Dialogue System Linly-Talker 🔥🔥
- 🚀 The Future of Digital Humans: The Empowerment Path of Linly-Talker + GPT-SoVIT Voice Cloning Technology
- Deploying Linly-Talker on AutoDL Platform (Super Detailed Tutorial for Beginners)
- Linly-Talker Update: Offline TTS Integration and Customized Digital Human Solutions

## TO DO LIST

- Completed the basic conversation system flow, capable of voice interactions .
- Integrated the LLM large model, including the usage of Linly , Qwen , and GeminiPro .
- Enabled the ability to upload any digital person's photo for conversation.
- Integrated FastAPI invocation for Linly.
- Utilized Microsoft TTS with advanced options, allowing customization of voice and tone parameters to enhance audio diversity.
- Added subtitles to video generation for improved visualization.
- GPT Multi-turn Dialogue System (Enhance the interactivity and realism of digital entities, bolstering their intelligence)
- Optimized the Gradio interface by incorporating additional models such as Wav2Lip, FunASR, and others.
- Voice Cloning Technology (Synthesize one's own voice using voice cloning to enhance the realism and interactive experience of digital entities)
- Integrate offline TTS (Text-to-Speech) along with NeRF-based methods and models.
- Linly-Talker WebUI supports multiple modules, multiple models, and multiple options
- Added MuseTalk functionality to Linly-Talker, achieving near real-time speed with very fast communication.
- Integrated MuseTalk into the Linly-Talker WebUI.
- Added CosyVoice, which provides high-quality text-to-speech (TTS) functionality and voice cloning capabilities. Additionally, updated to WaveGlow.
- Added Linly-Talker API documentation with detailed interface descriptions.
- Real-time Speech Recognition (Enable conversation and communication between humans and digital entities using voice)

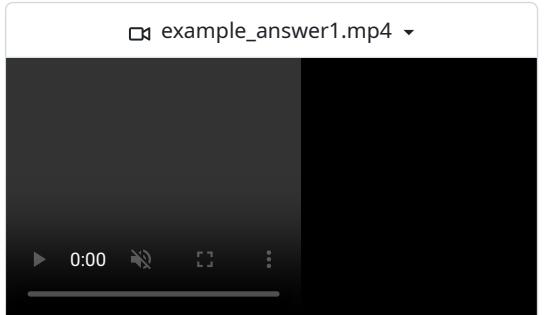
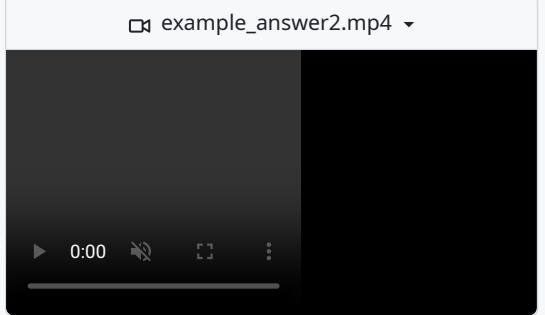
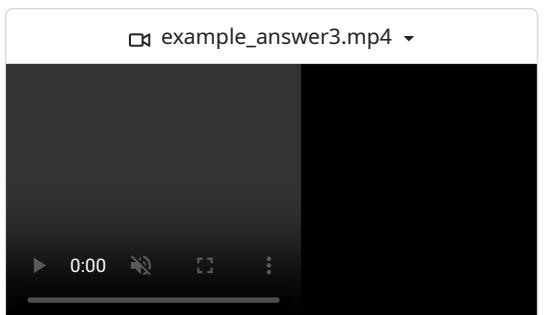
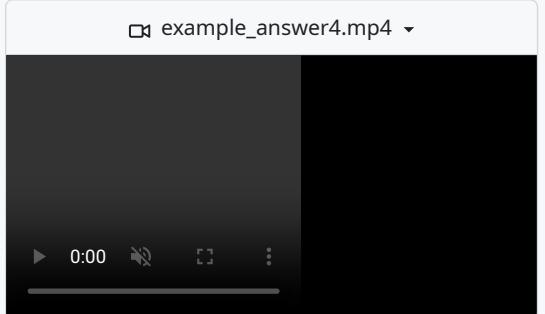
### Important

The Linly-Talker project is ongoing - pull requests are welcome! If you have any suggestions regarding new model approaches, research, or anything else, please feel free to contribute.

### Tip

If you encounter any issues during deployment, please consult the [Common Issues Summary](#) section, where I have compiled a list of all potential problems and their solutions.

## Example

文字/语音对话	数字人回答
应对压力最有效的方法是什么？	
如何进行时间管理？	
撰写一篇交响乐音乐会评论，讨论乐团的表演和观众的整体体验。	
翻译成中文：Luck is a dividend of sweat. The more you sweat, the luckier you get.	

## Setup Environment

### Note

AutoDL has released an image, which can be used directly at <https://www.codewithgpu.com/i/Kedreamix/Linly-Talker/Kedreamix-Linly-Talker>

```
docker pull registry.cn-beijing.aliyuncs.com/codewithgpu2/kedreamix-linly-talker:afGA8RPDLf
```

For Windows, I've included an all-in-one Python package. You can run the steps in sequence to install the necessary dependencies and down

## [Windows All-in-One Package](#)

Download the code:

```
git clone https://github.com/Kedreamix/Linly-Talker.git --depth 1  
cd Linly-Talker  
git submodule update --init --recursive
```

---

If you are using Linly-Talker, you can set up the environment directly with Anaconda, which covers almost all the dependencies required by the

```
conda create -n linly python=3.10  
conda activate linly  
  
# PyTorch Installation Option 1: Using conda  
# CUDA 11.8  
# conda install pytorch==2.4.1 torchvision==0.19.1 torchaudio==2.4.1 pytorch-cuda=11.8 -c pytorch -c nvidia  
# CUDA 12.1  
# conda install pytorch==2.4.1 torchvision==0.19.1 torchaudio==2.4.1 pytorch-cuda=12.1 -c pytorch -c nvidia  
# CUDA 12.4  
# conda install pytorch==2.4.1 torchvision==0.19.1 torchaudio==2.4.1 pytorch-cuda=12.4 -c pytorch -c nvidia  
  
# PyTorch Installation Option 2: Using pip  
# CUDA 11.8  
# pip install torch==2.4.1 torchvision==0.19.1 torchaudio==2.4.1 --index-url https://download.pytorch.org/whl/cu118  
# CUDA 12.1  
# pip install torch==2.4.1 torchvision==0.19.1 torchaudio==2.4.1 --index-url https://download.pytorch.org/whl/cu121  
# CUDA 12.4  
# pip install torch==2.4.1 torchvision==0.19.1 torchaudio==2.4.1 --index-url https://download.pytorch.org/whl/cu124  
pip install torch==2.0.1 torchvision==0.15.2 torchaudio==2.0.2 --index-url https://download.pytorch.org/whl/cu118  
  
conda install -q ffmpeg==4.2.2 # ffmpeg==4.2.2  
  
# Upgrade pip  
python -m pip install --upgrade pip  
# Change the PyPI source to speed up the installation of packages  
pip config set global.index-url https://pypi.tuna.tsinghua.edu.cn/simple  
  
pip install tb-nightly -i https://mirrors.aliyun.com/pypi/simple  
pip install -r requirements_webui.txt  
  
# Install dependencies related to musetalk  
pip install --no-cache-dir -U openmim  
mim install mmengine  
mim install "mmcv==2.1.0"  
mim install "mmdet>=3.1.0"  
mim install "mmpose>=1.1.0"  
  
#💡 The ttsfrd from CosyVoice can be replaced with WeTextProcessing, so a few steps can be omitted, while ensuring compatibility  
  
#⚠ Note: You must first download CosyVoice-ttsfrd. Complete the model download before proceeding with these steps.  
# mkdir -p CosyVoice/pretrained_models # Create directory CosyVoice/pretrained_models  
# mv checkpoints/CosyVoice_ckpt/CosyVoice-ttsfrd CosyVoice/pretrained_models # Move directory  
# unzip CosyVoice/pretrained_models/CosyVoice-ttsfrd/resource.zip # Unzip  
# This .whl library is only compatible with Python 3.8  
# pip install CosyVoice/pretrained_models/CosyVoice-ttsfrd/ttsfrd-0.3.6-cp38-cp38-linux_x86_64.whl  
  
# Install NeRF-based dependencies, which might have several issues and can be skipped initially  
pip install "git+https://github.com/facebookresearch/pytorch3d.git"  
# If you encounter problems installing PyTorch3D, you can use the following command to install it:  
# python scripts/install_pytorch3d.py  
pip install -r TFG/requirements_nerf.txt  
  
# If you encounter issues with pyaudio  
sudo apt-get update  
sudo apt-get install libasound-dev portaudio19-dev libportaudio2 libportaudiocpp0  
  
# Note the following modules. If installation fails, you can enter the directory and use pip install . or python setup.py install  
# NeRF/freqencoder  
# NeRF/gridencoder
```

```
# NeRF/raymarching
# NeRF/shencoder

# If you encounter sox compatibility issues
# ubuntu
sudo apt-get install sox libsox-dev
# centos
sudo yum install sox sox-devel
```

### ① Note

The installation process is very slow.

Below are some older installation methods, which might cause dependency conflicts, but they generally don't produce many bugs. For an easier setup, consider using the newer methods above.

To install the environment using Anaconda and PyTorch, follow the steps below:

```
conda create -n linly python=3.10
conda activate linly

# PyTorch Installation Method 1: Conda Installation (Recommended)
conda install pytorch==1.12.1 torchvision==0.13.1 torchaudio==0.12.1 cudatoolkit=11.3 -c pytorch

# PyTorch Installation Method 2: Pip Installation
pip install torch==1.12.1+cu113 torchvision==0.13.1+cu113 torchaudio==0.12.1 --extra-index-url https://download.pytorch.org/whl/cu113

conda install -q ffmpeg # ffmpeg==4.2.2
pip install -r requirements_app.txt
```

If you want to use models like voice cloning, you may need a higher version of PyTorch. However, the functionality will be more diverse. You can use the following command:

```
conda create -n linly python=3.10
conda activate linly

pip install torch==2.0.1 torchvision==0.15.2 torchaudio==2.0.2 --index-url https://download.pytorch.org/whl/cu118

conda install -q ffmpeg # ffmpeg==4.2.2
pip install -r requirements_app.txt

# Install dependencies for voice cloning
pip install -r VITS/requirements_gptsovits.txt
```

If you wish to use NeRF-based models, you may need to set up the corresponding environment:

```
# Install dependencies for NeRF
pip install "git+https://github.com/facebookresearch/pytorch3d.git"
pip install -r TFG/requirements_nerf.txt

# If there are issues with PyAudio, you can install the corresponding dependencies
# sudo apt-get install libasound-dev portaudio19-dev libportaudio2 libportaudiocpp0

# Note the following modules. If installation is unsuccessful, you can navigate to the path and use pip install . or python setup.py install
# NeRF/freqencoder
# NeRF/gridencoder
# NeRF/raymarching
# NeRF/shencoder
```

If you are using PaddleTTS, you can set up the corresponding environment with:

```
pip install -r TTS/requirements_paddle.txt
```

If you are using the FunASR speech recognition model, you can install the environment with:

```
pip install -r ASR/requirements_funasr.txt
```

If using the MuesTalk model, you can set up the environment with the following commands:

```
pip install --no-cache-dir -U openmim
mim install mmengine
mim install "mmcv>=2.0.1"
mim install "mmdet>=3.1.0"
mim install "mmpose>=1.1.0"
pip install -r TFG/requirements_musetalk.txt
```

### ① Note

Next, you need to install the corresponding models. You can download them using the following methods. Once downloaded, place the files in the specified directory.

- [Baidu \(百度云盘\)](#) (Password: linl)
- [huggingface](#)
- [modelscope](#)
- [Quark\(夸克网盘\)](#)

I made a script that can download all the models mentioned below without requiring much input from the user. This method is suitable for stable and reliable model download.

1. **Choose Download Method:** Users can choose to download models from three different sources: ModelScope, Huggingface, or Baidu Netdisk.
2. **Download Models:** Based on the user's selection, the script executes the corresponding download command.
3. **Move Model Files:** After downloading, the script moves the model files to the specified directory.
4. **Error Handling:** Error checks are included in each step of the operation. If any step fails, the script will output an error message and stop execution.

```
sh scripts/download_models.sh
```

### HuggingFace Download

If the download speed is too slow, consider using a mirror site. For more information, refer to [Efficiently Obtain Hugging Face Models Using Mirrors](#).

```
# Download pre-trained models from HuggingFace
git lfs install
git clone https://huggingface.co/Kedreamix/Linly-Talker --depth 1
# git lfs clone https://huggingface.co/Kedreamix/Linly-Talker --depth 1

# pip install -U huggingface_hub
# export HF_ENDPOINT=https://hf-mirror.com # Use a mirror site
huggingface-cli download --resume-download --local-dir-use-symlinks False Kedreamix/Linly-Talker --local-dir Linly-Talker
```

### ModelScope Download

```
# Download pre-trained models from ModelScope
# 1. Using git
git lfs install
git clone https://www.modelscope.cn/Kedreamix/Linly-Talker.git --depth 1
# git lfs clone https://www.modelscope.cn/Kedreamix/Linly-Talker.git

# 2. Download using Python code
pip install modelscope
from modelscope import snapshot_download
model_dir = snapshot_download('Kedreamix/Linly-Talker')
```

### Move All Models to the Current Directory

If you downloaded from Baidu Netdisk, you can refer to the directory structure at the end of the document to move the models.

```
# Move all models to the current directory
# Checkpoints contain SadTalker and Wav2Lip
mv Linly-Talker/checkpoints/* ./checkpoints

# Enhanced GFPGAN for SadTalker
# pip install gfpgan
# mv Linly-Talker/gfpgan ./gfpgan

# Voice cloning models
```

```

mv Linly-Talker/GPT_SoVITS/pretrained_models/* ./GPT_SoVITS/pretrained_models/
# Qwen large language model
mv Linly-Talker/Qwen ./

# MuseTalk model
mkdir -p ./Musetalk/models
mv Linly-Talker/MuseTalk/* ./Musetalk/models

```

For the convenience of deployment and usage, an `configs.py` file has been updated. You can modify some hyperparameters in this file for cu:

```

# Device Running Port
port = 7870

# API Running Port and IP
# Localhost port is 127.0.0.1; for global port forwarding, use "0.0.0.0"
ip = '127.0.0.1'
api_port = 7871

# Linly Model Path
mode = 'api' # For 'api', Linly-api-fast.py must be run first
mode = 'offline'
model_path = 'Linly-AI/Chinese-LLaMA-2-7B-hf'

# SSL Certificate (required for microphone interaction)
# Preferably an absolute path
ssl_certfile = "./https_cert/cert.pem"
ssl_keyfile = "./https_cert/key.pem"

```

This file allows you to adjust parameters such as the device running port, API running port, Linly model path, and SSL certificate paths for ease

## API Documentation

In the [api/README.md](#) file, we provide detailed information about the usage and configuration of the Linly-Talker API. This documentation incli

recognition, and generating speech.

For detailed API interface descriptions, please refer to the `api/README.md` file.

## ASR - Speech Recognition

For detailed information about the usage and code implementation of Automatic Speech Recognition (ASR), please refer to [ASR - Bridging the C](#)

### Whisper

To implement ASR (Automatic Speech Recognition) using OpenAI's Whisper, you can refer to the specific usage methods provided in the GitHub

### FunASR

The speech recognition performance of Alibaba's FunASR is quite impressive and it is actually better than Whisper in terms of Chinese languag

### Coming Soon

Welcome everyone to provide suggestions, motivating me to continuously update the models and enrich the functionality of Linly-Talker.

## TTS - Text To Speech

For detailed information about the usage and code implementation of Text-to-Speech (TTS), please refer to [TTS - Empowering Digital Humans v](#)

### Edge TTS

To use Microsoft Edge's online text-to-speech service from Python without needing Microsoft Edge or Windows or an API key, you can refer to t

#### ⚠ Warning

Due to some issues with the Edge TTS repository, it seems that Microsoft has restricted certain IPs. For more details, refer to [403 error is bac](#)

### PaddleTTS

In practical use, there may be scenarios that require offline operation. Since Edge TTS requires an online environment to generate speech, we have to find alternative solutions.

## Coming Soon

Welcome everyone to provide suggestions, motivating me to continuously update the models and enrich the functionality of Linly-Talker.

## Voice Clone

For detailed information about the usage and code implementation of Voice Clone, please refer to [Voice Clone - Stealing Your Voice Quietly Due to GPT-SoVITS](#).

### GPT-SoVITS (Recommend)

Thank you for your open source contribution. I have also found the [GPT-SoVITS](#) voice cloning model to be quite impressive. You can find the project on GitHub.

## XTTS

Coqui XTTS is a leading deep learning toolkit for Text-to-Speech (TTS) tasks, allowing for voice cloning and voice transfer to different languages.

🐸 TTS is a library for advanced text-to-speech generation.

🚀 Over 1100 pre-trained models for various languages.

🛠 Tools for training new models and fine-tuning existing models in any language.

📚 Utility programs for dataset analysis and management.

- Experience XTTS online <https://huggingface.co/spaces/coqui/xtts>
- Official GitHub repository: <https://github.com/coqui-ai/TTS>

## CosyVoice

CosyVoice is an open-source multilingual speech understanding model developed by Alibaba's Tongyi Lab, focusing on high-quality speech synthesis.

CosyVoice supports one-shot voice cloning technology, enabling the generation of realistic and natural-sounding voices with details such as pronunciation and intonation.

GitHub project link: [CosyVoice GitHub](#)

CosyVoice includes several pre-trained speech synthesis models, mainly:

1. **CosyVoice-300M**: Supports zero-shot and cross-lingual speech synthesis in Chinese, English, Japanese, Cantonese, Korean, and other languages.
2. **CosyVoice-300M-SFT**: A model focused on supervised fine-tuning (SFT) inference.
3. **CosyVoice-300M-Instruct**: A model that supports command-based inference, capable of generating speech with specific tones, emotions, and styles.

## Key Features

1. **Multilingual Support**: Capable of handling various languages including Chinese, English, Japanese, Cantonese, and Korean.
2. **Multi-style Speech Synthesis**: Allows control over the tone and emotion of the generated speech through commands.
3. **Streaming Inference Support**: Future updates will include streaming inference modes, such as KV caching and SDPA, for real-time optimization.

Currently, Linly-Talker integrates three features from CosyVoice: pre-trained voice cloning, 3s rapid cloning, and cross-lingual cloning. Stay tuned!

	PROMPT TEXT
Pre-trained Voice	中文女 音色 ('中文女', '中文男', '日语男', '粤语女', '英文女', '英文男', '韩语女')
3s Language Cloning	希望你以后能够做的比我还好呦。

## Cross-lingual Cloning

在那之后，完全收购那家公司，因此保持管理层的一致性，利益与即将加入家族的资产保持一致。这就是我们有时不买下全

## Coming Soon

Welcome everyone to provide suggestions, motivating me to continuously update the models and enrich the functionality of Linly-Talker.

## THG - Avatar

Detailed information about the usage and code implementation of digital human generation can be found in [THG - Building Intelligent Digital](#)

## SadTalker

Digital persona generation can utilize SadTalker (CVPR 2023). For detailed information, please visit <https://sadtalker.github.io>.

Before usage, download the SadTalker model:

```
bash scripts/sadtalker_download_models.sh
```

[Baidu \(百度云盘\)](#) (Password: linly)

[Quark\(夸克网盘\)](#)

If downloading from Baidu Cloud, remember to place it in the `checkpoints` folder. The model downloaded from Baidu Cloud is named `sadt`

## Wav2Lip

Digital persona generation can also utilize Wav2Lip (ACM 2020). For detailed information, refer to <https://github.com/Rudrabha/Wav2Lip>.

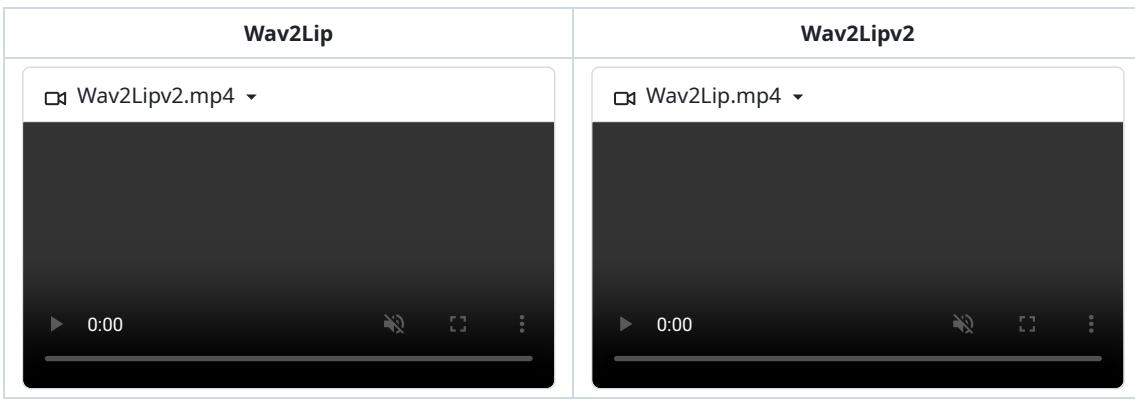
Before usage, download the Wav2Lip model:

Model	Description	Link to the model
Wav2Lip	Highly accurate lip-sync	<a href="#">Link</a>
Wav2Lip + GAN	Slightly inferior lip-sync, but better visual quality	<a href="#">Link</a>
Expert Discriminator	Weights of the expert discriminator	<a href="#">Link</a>
Visual Quality Discriminator	Weights of the visual disc trained in a GAN setup	<a href="#">Link</a>

## Wav2Lipv2

Inspired by the repository [https://github.com/primepeake/wav2lip\\_288x288](https://github.com/primepeake/wav2lip_288x288), Wav2Lipv2 uses a newly trained 288 model to achieve higher qualit

Additionally, by employing YOLO for facial detection, the overall effect is improved. You can compare and test the results in Linly-Talker. The mo



## ER-NeRF

ER-NeRF (ICCV 2023) is a digital human built using the latest NeRF technology. It allows for the customization of digital characters and can reco Updated: Taking inspiration from the likeness of Obama, for better results, consider cloning and customizing the voice of digital personas for ir

## MuseTalk

MuseTalk is a real-time, high-quality audio-driven lip synchronization model capable of running at over 30 frames per second on an NVIDIA Tes

MuseTalk is trained to operate within the latent space of ft-mse-vae and offers the following features:

- **Unseen Face Synchronization:** It can modify unseen faces based on input audio, with a face region size of 256 x 256.
- **Multi-language Support:** Supports audio inputs in various languages, including Chinese, English, and Japanese.
- **High-performance Real-time Inference:** Achieves real-time inference at over 30 frames per second on an NVIDIA Tesla V100.
- **Facial Center Point Adjustment:** Allows the adjustment of the facial region's center point, significantly impacting the generated results.
- **HDTF Dataset Training:** Provides model checkpoints trained on the HDTF dataset.
- **Upcoming Training Code Release:** Training code will be released soon, facilitating further development and research.

MuseTalk offers an efficient and versatile tool for precise audio synchronization with facial expressions in virtual humans, marking a significant In Linly-Talker, MuseTalk has been integrated to perform inference on videos based on MuseV, achieving an ideal speed for conversations with

## Coming Soon

Welcome everyone to provide suggestions, motivating me to continuously update the models and enrich the functionality of Linly-Talker.

## LLM - Conversation

For detailed information about the usage and code implementation of Large Language Models (LLM), please refer to [LLM - Empowering Digital](#)

## Linly-AI

Linly-AI is a Large Language model developed by CVI at Shenzhen University. You can find more information about Linly-AI on their GitHub rep

Download Linly models: <https://huggingface.co/Linly-AI/Chinese-LLaMA-2-7B-hf>

You can use `git` to download:

```
git lfs install
git clone https://huggingface.co/Linly-AI/Chinese-LLaMA-2-7B-hf
```

Alternatively, you can use the `huggingface` download tool `huggingface-cli`:

```
pip install -U huggingface_hub

# Set up mirror acceleration
# Linux
export HF_ENDPOINT="https://hf-mirror.com"
# Windows PowerShell
$env:HF_ENDPOINT="https://hf-mirror.com"
```

```
huggingface-cli download --resume-download Linly-AI/Chinese-LLaMA-2-7B-hf --local-dir Linly-AI/Chinese-LLaMA-2-7B-hf
```

## Qwen

Qwen is an AI model developed by Alibaba Cloud. You can check out the GitHub repository for Qwen here: <https://github.com/QwenLM/Qwen>

If you want to quickly use Qwen, you can choose the 1.8B model, which has fewer parameters and can run smoothly even with limited GPU memory.

You can download the Qwen 1.8B model from this link: [https://huggingface.co/Qwen/Qwen-1\\_8B-Chat](https://huggingface.co/Qwen/Qwen-1_8B-Chat)

You can use `git` to download:

```
git lfs install  
git clone https://huggingface.co/Qwen/Qwen-1_8B-Chat
```

Alternatively, you can use the `huggingface` download tool `huggingface-cli`:

```
pip install -U huggingface_hub  
  
# Set up mirror acceleration  
# Linux  
export HF_ENDPOINT="https://hf-mirror.com"  
# Windows PowerShell  
$env:HF_ENDPOINT="https://hf-mirror.com"  
  
huggingface-cli download --resume-download Qwen/Qwen-1_8B-Chat --local-dir Qwen/Qwen-1_8B-Chat
```

## Gemini-Pro

Gemini-Pro is an AI model developed by Google. To learn more about Gemini-Pro, you can visit their website: <https://deepmind.google/technologies/gemini-pro/>

If you want to request an API key for Gemini-Pro, you can visit this link: <https://makersuite.google.com/>

## ChatGPT

From OpenAI, requires API application. For more information, please visit <https://platform.openai.com/docs/introduction>.

## ChatGLM

From Tsinghua University, for more information please visit <https://github.com/THUDM/ChatGLM3>.

## GPT4Free

For free access to GPT-4 and other models, you can refer to <https://github.com/xtekky/gpt4free>. This resource provides methods to utilize these models.

## LLM Multiple Model Selection

In the `webui.py` file, easily select the model you need.  For the first run, make sure to download the model first. Refer to Qwen1.8B.

## Coming Soon

Welcome everyone to provide suggestions, motivating me to continuously update the models and enrich the functionality of Linly-Talker.

## Optimizations

Some optimizations:

- Use fixed input face images, extract features beforehand to avoid reading each time
- Remove unnecessary libraries to reduce total time
- Only save final video output, don't save intermediate results to improve performance
- Use OpenCV to generate final video instead of mimwrite for faster runtime

## Gradio

Gradio is a Python library that provides an easy way to deploy machine learning models as interactive web apps.

For Linly-Talker, Gradio serves two main purposes:

1. **Visualization & Demo:** Gradio provides a simple web GUI for the model, allowing users to see the results intuitively by uploading an image
2. **User Interaction:** The Gradio GUI can serve as a frontend to allow end users to interact with Linly-Talker. Users can upload their own image

Specifically, we create a Gradio Interface in app.py that takes image and text inputs, calls our function to generate the response video, and disp

In summary, Gradio provides visualization and user interaction interfaces for Linly-Talker, serving as effective means for showcasing system cap

If considering real-time conversation, it may be necessary to switch to a different framework or customize Gradio. Looking forward to workin

## Start WebUI

Previously, I had separated many versions, but it became cumbersome to run multiple versions. Therefore, I have added a WebUI feature to pr

### WebUI

The current features available in the WebUI are as follows:

- Text/Voice-based dialogue with virtual characters (fixed characters with male and female roles)
- Dialogue with virtual characters using any image (you can upload any character image)
- Multi-turn GPT dialogue (incorporating historical dialogue data to maintain context)
- Voice cloning dialogue (based on GPT-SoVITS settings for voice cloning, including a built-in smoky voice that can be cloned based on the vo
- Digital Persona Text/Voice Playback (based on input text/voice)
- Multiple modules+Multiple models+Multiple choices
  - Multiple role selections: Female/Male/Custom (each part can automatically upload images) Coming Soon
  - Multiple TTS model selections: EdgeTTS / PaddleTTS / GPT-SoVITS / CosyVoice / Coming Soon
  - Multiple LLM model selections: Linly / Qwen / ChatGLM / GeminiPro / ChatGPT / Coming Soon
  - Multiple Talker model selections: Wav2Lip / Wav2Lipv2 / SadTalker / ERNeRF / MuseTalk/ Coming Soon
  - Multiple ASR model selections: Whisper / FunASR / Coming Soon

## 角色选择

女性角色

男性角色

自定义角色

Text To Speech Method (Edge-TTS利用微软的TTS, PaddleSpeech是离线的TTS, 不过第一次运行会自动下载模型)

Edge-TTS

PaddleTTS

GPT-SoVITS克隆声音

## 语音识别模型选择

Whisper-tiny

Whisper-base

FunASR

Comming Soon!!!

## 数字人模型选择

SadTalker

Wav2Lip

ER-NeRF

MuseTalk

Comming Soon!!!

## LLM 模型选择

Qwen



You can directly run the web UI to obtain results. The page you will see is as follows:

```
# WebUI  
python webui.py
```

## Linly-Talker WebUI

个性化角色互动

数字人多轮智能对话

MuseTalk数字人实时对话

Linly 智能对话系统 (Linly-Talker) MuseTalker 数字人实时

[知乎] [bilibili] [GitHub] [个人主页]

Linly-Talker是一款创新的数字人对话系统，它融合了最新的人工智能技术，包括大型语言模型（LLM）、自动语音识别（ASR）、文本转

MuseV Video

MuseV: need help? please visit MuseVDemo to generate Video

<https://huggingface.co/spaces/AnchorFake/MuseVDemo>

#### □ Reference Video



对话

语音对话

录制

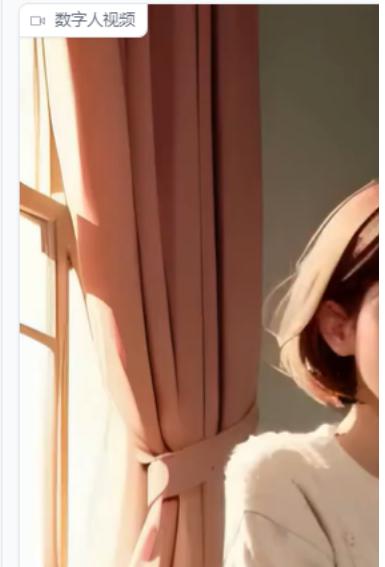
输入文字/问题

你好，我是Linly 智能对话系统 (Linly Talk  
来实时对话，感谢你的关注哦

语音

MuseTalk Video

□ 数字人视频



This time, we've updated the interface. We can freely select the fine-tuned model of GPT-SoVITS to achieve voice cloning. Simply upload a referen-

## TTS Method语音方法调节

Edge-TTS

PaddleTTS

GPT-SoVITS

### GPT模型路径

- Gnews
- chinese-hubert-base
- chinese-roberta-wwm-ext-large
- ➜ Gnews-e15.ckpt
- ➜ s1bert25hz-2kh-longer-epoch=68e-step=50232.ckpt

### SoVITS模型路径

- Gnews
- chinese-hubert-base
- chinese-roberta-wwm-ext-large
- ➜ Gnews\_e8\_s96.pth
- ➜ s2D488k.pth
- ➜ s2G488k.pth

## 加载模型

请上传3~10秒内参考音频，超过会报错！



0:00

0:04



1x



使用语音问答的麦  
克风

参考音频的文本

才之后很快就能够  
立起来这个呢是要  
国内

参考音频的语种

中文



## 语音识别 - 克隆参考音频

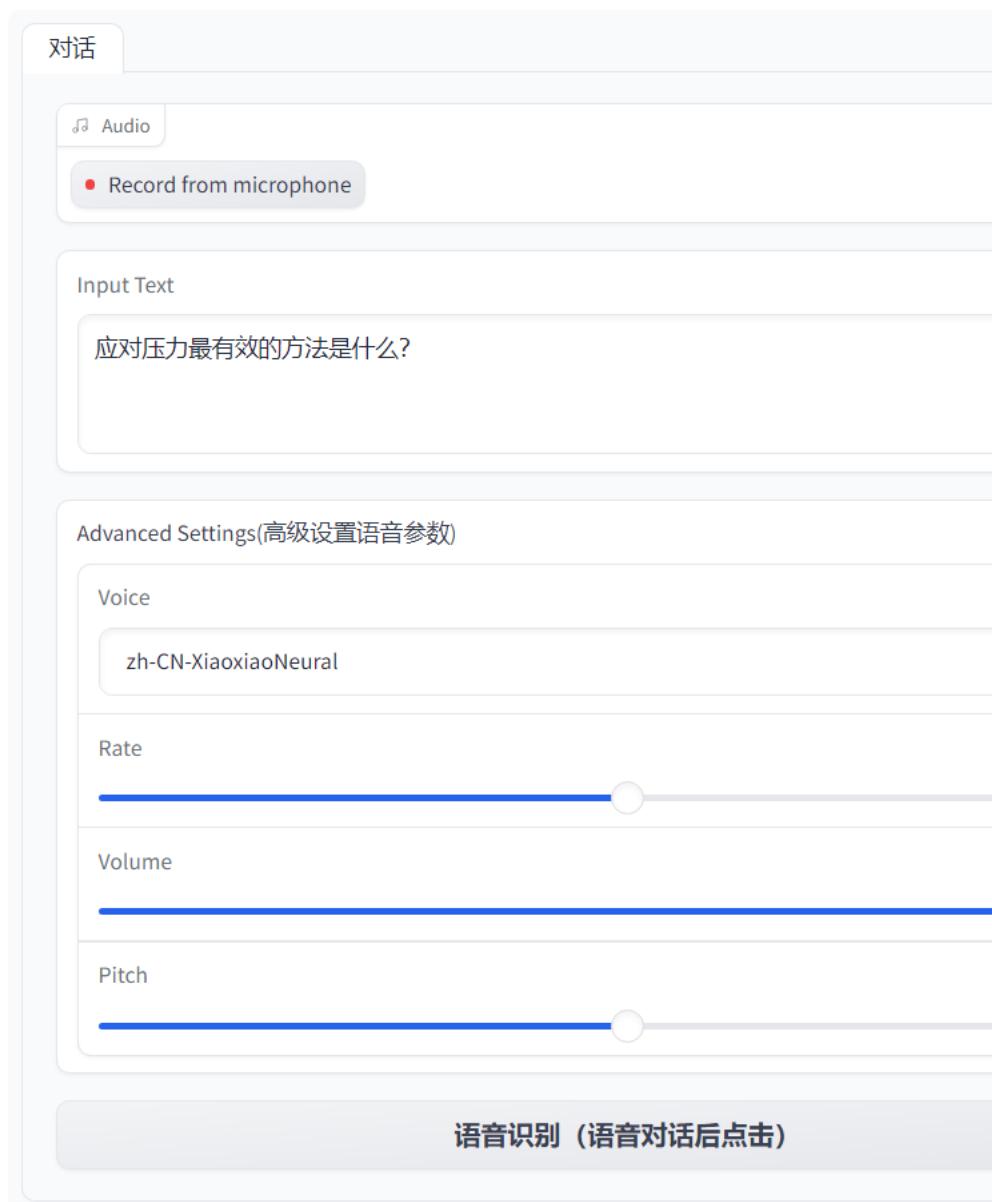
### Old Verison

There are three modes for the current startup, and you can choose a specific setting based on the scenario.

The first mode involves fixed Q&A with a predefined character, eliminating preprocessing time.

```
python app.py
```

Linly-Talker 是一款智能 AI 对话



The first mode has recently been updated to include the Wav2Lip model for dialogue.

```
python appv2.py
```

The second mode allows for conversing with any uploaded image.

```
python app_img.py
```

Linly-Talker 是一款智能 AI 对话

Source image

Source image



对话

Audio

Record from microphone

Input Text

文字对话

如何进行时间管理?

Advanced Settings

语音识别 (语音对话后点击)

## Text Examples

Examples

应对压力最有效的方法是什么? 如何进行时间管理?

The third mode builds upon the first one by incorporating a large language model for multi-turn GPT conversations.

```
python app_multi.py
```

Advanced Settings(高级设置)

数字人问答



 生成数字人视频 (对话后)

Now, the part of voice cloning has been added, allowing for freely switching between cloned voice models and corresponding person images. If

```
python app_vits.py
```

A fourth method has been added, which does not fixate on a specific scenario for conversation. Instead, it allows for direct input of voice or the

ER-NeRF is trained on videos of a single individual, so a specific model needs to be replaced to render and obtain the correct results. It come

```
python app_talk.py
```

# Linly 智能对话系统 (Linly-Talker)

[知乎] [bilibili] [GitHub] [个人主页]

Linly-Talker 是一款智能 AI 对话系统，结合了大型语言模型 (LLMs) 与视觉模型，是一种新颖的人工智能交互方式

The screenshot shows the Linly-Talker application interface. At the top left, there are tabs for "图片人物" (Image Person) and "视频人物" (Video Person). Below these are two sections: one for "Source image" (with a placeholder for dragging or clicking to upload) and another for "Generated video" (with a placeholder for "Generated video"). In the center, there's a "语音" (Voice) section showing a waveform from 0:00 to 0:03, with controls for play/pause, volume, and speed (1x). Below this is an "Input Text" field containing the text "应对压力最有效的方法是什么？". At the bottom left, there's a note about the Text To Speech Method (Edge-TTS vs PaddleTTS), with Edge-TTS selected. On the right side of the interface, a large video frame shows a close-up of a person's face.

MuseTalk has been integrated into Linly-Talker, enabling efficient preprocessing of MuseV-generated videos. Once preprocessed, these videos can be used directly in the system.

To run the application, use the following command:

```
python app_musetalk.py
```

# Linly 智能对

Linly-Talker是一款创新的数字人对话系统，它融合了最新

MuseV Video

MuseV: need help? please visit MuseVDemo to generate Video  
<https://huggingface.co/spaces/AnchorFake/MuseVDemo>

Reference Video



将视频拖放到此处

- 或 -

点击上传

BBox\_shift 推荐值下限，在生成初始结果后生成相应的 bbox，可以根据该参考值进行调整。一般来说，在我们的实验观察（面部部分移动）通常会增加嘴巴的张开度，而负值（向上半部分移动）通常会减小嘴巴的张开度。然而，需要注意的是，这并不是绝对的规则，用户可以根据期望效果来调整该参数。

BBox\_shift value, px

0



bbox\_shift\_

TTS Method语音方法调节

## Folder structure

### ① Note

The folder structure of the weight files is as follows:

## 角色选择

显示菜单

女性角色

男性角色

自定义角色

- Baidu (百度云盘) : You can download the weights from [here](#) (Password: linly).
- huggingface : You can access the weights at [this link](#).
- modelscope : The weights will be available soon at [this link](#).
- Qurak(夸克网盘) : You can download the weights from [here](#)

```
Linly-Talker/
├── checkpoints
│   ├── audio_visual_encoder.pth
│   ├── hub
│   │   └── checkpoints
│   │       └── s3fd-619a316812.pth
│   ├── lipsync_expert.pth
│   ├── mapping_00109-model.pth.tar
│   ├── mapping_00229-model.pth.tar
│   ├── May.json
│   ├── May.pth
│   ├── Obama_ave.pth
│   ├── Obama.json
│   ├── Obama.pth
│   ├── ref_eo.npy
│   ├── ref.npy
│   ├── ref.wav
│   ├── SadTalker_V0.0.2_256.safetensors
│   ├── visual_quality_disc.pth
│   ├── wav2lip_gan.pth
│   └── wav2lip.pth
├── gfp gan
│   └── weights
│       ├── alignment_WFLW_4HG.pth
│       └── detection_Resnet50_Final.pth
└── GPT_SoVITS
    └── pretrained_models
        ├── chinese-hubert-base
        │   ├── config.json
        │   ├── preprocessor_config.json
        │   └── pytorch_model.bin
        ├── chinese-roberta-wwm-ext-large
        │   ├── config.json
        │   ├── pytorch_model.bin
        │   └── tokenizer.json
        ├── README.md
        ├── s1bert25hz-2kh-longer-epoch=68e-step=50232.ckpt
        ├── s2D488k.pth
        ├── s2G488k.pth
        └── speech_paraformer-large_asr_nat-zh-cn-16k-common-vocab8404.pytorch
├── MuseTalk
    ├── models
    │   ├── dwpose
    │   │   └── dw_ll_ucoco_384.pth
    │   ├── face-parse-bisent
    │   │   └── 79999_iter.pth
    │   └── musetalk
    │       ├── musetalk.json
    │       └── pytorch_model.bin
    ├── README.md
    ├── sd-vae-ft-mse
    │   ├── config.json
    │   └── diffusion_pytorch_model.bin
    └── whisper
        └── tiny.pt
└── Qwen
    └── Qwen-1_8B-Chat
        ├── assets
        │   └── logo.jpg
```

```

    |   |   ├── qwen_tokenizer.png
    |   |   ├── react_showcase_001.png
    |   |   ├── react_showcase_002.png
    |   |   └── wechat.png
    |   ├── cache_autogptq_cuda_256.cpp
    |   ├── cache_autogptq_cuda_kernel_256.cu
    |   ├── config.json
    |   ├── configuration_qwen.py
    |   ├── cpp_kernels.py
    |   ├── examples
    |   |   └── react_prompt.md
    |   ├── generation_config.json
    |   ├── LICENSE
    |   ├── model-00001-of-00002.safetensors
    |   ├── model-00002-of-00002.safetensors
    |   ├── modeling_qwen.py
    |   ├── model.safetensors.index.json
    |   ├── NOTICE
    |   ├── qwen_generation_utils.py
    |   ├── qwen.tiktoken
    |   ├── README.md
    |   ├── tokenization_qwen.py
    |   └── tokenizer_config.json
    └── Whisper
        ├── base.pt
        └── tiny.pt
    └── FunASR
        ├── punc_ct-transformer_zh-cn-common-vocab272727-pytorch
        |   ├── configuration.json
        |   ├── config.yaml
        |   ├── example
        |   |   └── punc_example.txt
        |   ├── fig
        |   |   └── struct.png
        |   ├── model.pt
        |   ├── README.md
        |   └── tokens.json
        ├── speech_fsmn_vad_zh-cn-16k-common-pytorch
        |   ├── am.mvn
        |   ├── configuration.json
        |   ├── config.yaml
        |   ├── example
        |   |   └── vad_example.wav
        |   ├── fig
        |   |   └── struct.png
        |   ├── model.pt
        |   └── README.md
        └── speech_seaco_paraformer_large_asr_nat-zh-cn-16k-common-vocab8404-pytorch
            ├── am.mvn
            ├── asr_example_hotword.wav
            ├── configuration.json
            ├── config.yaml
            ├── example
            |   ├── asr_example.wav
            |   └── hotword.txt
            ├── fig
            |   ├── res.png
            |   └── seaco.png
            ├── model.pt
            ├── README.md
            ├── seg_dict
            └── tokens.json
    └── README.md

```

## Reference

### ASR

- <https://github.com/openai/whisper>
- <https://github.com/alibaba-damo-academy/FunASR>

### TTS

- <https://github.com/rany2/edge-tts>

- <https://github.com/PaddlePaddle/PaddleSpeech>

## LLM

- <https://github.com/CVI-SZU/Linly>
- <https://github.com/QwenLM/Qwen>
- <https://deepmind.google/technologies/gemini/>
- <https://github.com/THUDM/ChatGLM3>
- <https://openai.com>

## THG

- <https://github.com/OpenTalker/SadTalker>
- <https://github.com/Rudrabha/Wav2Lip>
- <https://github.com/Fictionarry/ER-NeRF>

## Voice Clone

- <https://github.com/RVC-Boss/GPT-SoVITS>
- <https://github.com/coqui-ai/TTS>

## License

### ⓘ Caution

When using this tool, please comply with all applicable laws, including copyright, data protection, and privacy laws. Do not use, modify, distr

`Linly-Talker` follows the MIT License. In addition to adhering to the MIT License, ensure that you comply with all license agreements for any i

## Star History

