# 数据集

1. SE-cap https://arxiv.org/abs/2312.10381

数据集：内部的EMOSpeech中文数据集，https://github.com/thuhcsi/SECap

仓库里有600条 wav 和对应的描述文本 ✔️

2. AlignCap: Aligning Speech Emotion Captioning to Human Preferences

https://aclanthology.org/2024.emnlp-main.224.pdf

数据集：

a. MER2023 (A large-scale video emotion reason dataset) → MER23SEC

b. NNIME (A Chinese interactive multimodal emotion corpus)：用来评估模型在其它数据集上的可迁移性的。

c. EMOSEC (他们自己提出的41 hours of Chinese-English Speech Emotion Captioning datasets) https://zenodo.org/records/10948423 ✔️

(在zenodo上只放了情绪的描述文本merged_file.json, 但它是基于 ESD 和 IEMOCAP构建的，ESD和IEMPCAP都能单独找到)



```
{
    "wav": "0003_000898.wav",
    "caption_ch1": "音频中的人情绪是生气的，说话的语调是平常的，语气是愤怒的，语速是中等的，声音变化是逐渐降低的。",
    "caption_ch2": "这是一个女人在讲话，听起来他的语气是积极的。",
    "caption_en1": "The audio is of a woman speaking, in a normal tone, with a fast speaking speed, saying, \"当然知道我们各有千秋\".",
    "caption_en2": "Based on the voice, it sounds like this person is angry in the audio, and their emotions are also reflected in their tone of
},
```

d. ESD: **Publicly Available Emotional Speech Dataset (ESD) for Speech Synthesis and Voice   Conversion** ✔️

https://github.com/HLTSingapore/Emotional-Speech-Data

e. IEMOCAP: https://sail.usc.edu/iemocap/index.html

IEMOCAP还在申请中......

3. EmotionCaps: Enhancing Audio Captioning Through Emotion-Augmented
   Data Generation

https://zenodo.org/records/13755932 （情感增强的音频字幕）

但这个数据集面向环境声为主，并非专门的"语音"集合，但是可以作为让 caption 更贴合情绪语气的过渡数据 ✅

| segment_id | caption |
| --- | --- |
| 0-00uubiDNU_210000 | Loud, chaotic music fills the air, creating a disorienting and overwhelming atmosphere. |
| 006SLC2yxCk_2000 | Amidst the chaotic symphony of a baby crying, gurgling, and a woman speaking, a tick and generic impact sounds add to the frenzy. |
| 007Cl0YLNRk_30000 | Amidst the unpleasant wind, a chaotic mix of voices, water, shouts, and laughter fills the air, accompanied by the intrusive wind noise captured by the microphone. |

4. MER2024

https://huggingface.co/datasets/MERChallenge/MER2024

5,030 条带标注的视频切片（含音频），都有对应字幕和离散情感标签。但只有其中332 条有与情感相关的描述/理由 ✅