

# DeepDubber-V1: Towards High Quality and Dialogue, Narration, Monologue Adaptive Movie Dubbing Via Multi-Modal Chain-of-Thoughts Reasoning Guidance.

Anonymous ICCV submission

Paper ID 8787

## Abstract

001 Current movie dubbing technology can generate the de-  
 002 sired voice from a given speech prompt, ensuring good syn-  
 003 chronization between speech and visuals while accurately  
 004 conveying the intended emotions. However, in movie dub-  
 005 bing, key aspects such as adapting to different dubbing  
 006 styles, handling dialogue, narration, and monologue effec-  
 007 tively, and understanding subtle details like the age and  
 008 gender of speakers, have not been well studied. To ad-  
 009 dress this challenge, we propose a framework of multimodal  
 010 large language model. First, it utilizes multimodal Chain-  
 011 of-Thought (CoT) reasoning methods on visual inputs to  
 012 understand dubbing styles and fine-grained attributes. Sec-  
 013 ond, it generates high-quality dubbing through large speech  
 014 generation models, guided by multimodal conditions. Addi-  
 015 tionally, we have developed a movie dubbing dataset with  
 016 CoT annotations. The evaluation results demonstrate a per-  
 017 formance improvement over state-of-the-art methods across  
 018 multiple datasets. In particular, for the evaluation met-  
 019 rics, the SPK-SIM and EMO-SIM increases from 82.48% to  
 020 89.74%, 66.24% to 78.88% for dubbing setting 2.0 on V2C-  
 021 Animation dataset, LSE-D and MCD-SL decreases from  
 022 14.79 to 14.63, 5.24 to 4.74 for dubbing setting 2.0 on Grid  
 023 dataset, SPK-SIM increases from 64.03 to 83.42 and WER  
 024 decreases from 52.69% to 23.20% for initial reasoning set-  
 025 ting on proposed CoT-Movie-Dubbing dataset in the com-  
 026 parison with the state-of-the art models.

## 027 1. Introduction

028 Dubbing involves adding the correct human voice to a  
 029 video's dialogue, ensuring synchronization with the char-  
 030 acter's lip movements, and conveying the emotions of the  
 031 scene. It plays a vital role in film, television, animation  
 032 and gaming, enhancing immersion and effectively convey-  
 033 ing emotions and atmosphere. Existing dubbing methods  
 034 can be categorized into two groups, both of which focus on

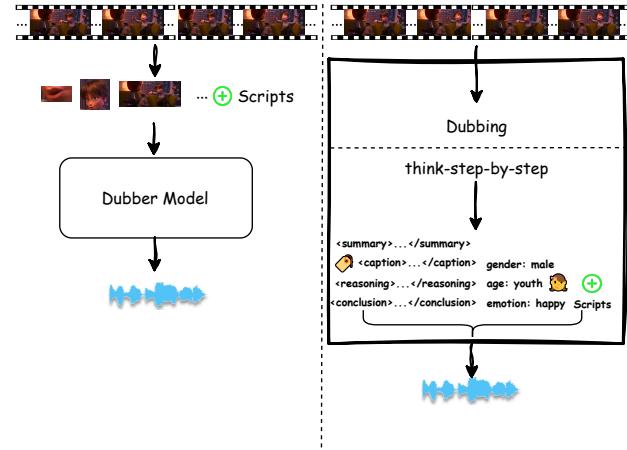


Figure 1. Current Dubbing models [14, 17, 69] (Left). Proposed Dubbing Models (Right) For dubbing types and fine-grained attributes.

035 learning different styles of key prior information to gener-  
 036 ate high-quality voices. The first group focuses on learn-  
 037 ing effective speaker style representations [7, 15, 23, 60].  
 038 The second group aims to learn appropriate prosody by  
 039 utilizing visual information from the given video input  
 040 [15, 25, 37, 70]. However, the accuracy of these priors  
 041 is insufficient and inadequate for movie dubbing in real-  
 042 world scenarios. For example, adaptive dubbing for dif-  
 043 ferent types, such as dialogue, narration, and monologue,  
 044 as well as fine-grained attributes such as expected ages and  
 045 genders, has not been thoroughly studied [17, 25].

046 With the rapid advancement of large language reasoning  
 047 models with step-by-step thinking ability [2, 19, 47–49, 52,  
 048 64] and methods that enhance reasoning capabilities to in-  
 049 terpret visual information through CoT, MLLM has increas-  
 050 ingly shown their potential in multimodal reasoning and un-  
 051 derstanding tasks [4, 11, 20, 35, 39, 43, 50, 57, 71]. These  
 052 advancements in reasoning capabilities within MLLM hold  
 053 promise for accurately providing dubbing types and fine-

054 grained attributes.

055 Therefore, we propose a multimodal large language  
056 model for high-quality movie dubbing that effectively un-  
057 derstands dubbing styles and fine-grained attributes. First,  
058 through multi-modal CoT learning, a multimodal large  
059 language model is trained to improve its reasoning abil-  
060 ity, enabling a better understanding of dubbing types (di-  
061 alogue, narration, monologue) and fine-grained attributes  
062 from video inputs. Secondly, a large multimodal speech  
063 generation model is trained with designed control mech-  
064 anisms using multiple-modal conditions. Thirdly, we create  
065 a CoT multi-modal movie dubbing dataset annotated with  
066 step-by-step reasoning instructions.

## 067 2. Related Work

### 068 2.1. Visual Voice Cloning

069 Current advanced dubbing technologies significantly en-  
070 hance speech-video synchronization and emotional expres-  
071 sion by integrating visual and textual information. Some  
072 works focus on improving speaker identity to handle multi-  
073 speaker scenes [14, 16, 17, 69]. For example, Speaker2Dub  
074 [69] introduces speaker embedding extracted by pre-trained  
075 GE2E to the phoneme encoder and the mel spectrogram de-  
076 coder by a learnable style affine transform, while StyleDub-  
077 ber [17] proposes a multi-scale style adapter with phoneme  
078 and utterance level to strengthen speaker characteristics.  
079 In addition, some works attempt to combine visual re-  
080 presentation to enhance prosody expressive [14, 25, 37, 70].  
081 For example, HPMDubbing [14] is a hierarchical dubbing  
082 method that bridges acoustic details with visual informa-  
083 tion: lip motion, face region, and scene. To improve con-  
084 textual prosody, MCDubber [70] enlarges the modeling ob-  
085 ject from a single sentence to the previous and follow-  
086 ing sentences, incorporating more contextual video scenes.  
087 Although speaker identity and prosody modeling have re-  
088 ceived attention, existing works still suffer from poor lip-  
089 sync and lifeless emotional expression, which is unaccept-  
090 able in dubbing.

### 091 2.2. Flow-Matching Speech Generation

092 Flow Matching [41] is a simulation-free method to train  
093 Continuous Normalizing Flows (CNFs) [8] models, which  
094 model arbitrary probability path and capture the proba-  
095 bility trajectories represented by diffusion processes [59].  
096 Due to its advantages of high sampling speed and genera-  
097 tion quality, flow matching has attracted significant atten-  
098 tion in speech generation [21, 36, 45]. Recently, Matcha-  
099 TTS [45] and DiTTo-TTS [38] have introduced optimal-  
100 transport conditional flow matching (OT-CFM) for train-  
101 ing, which yields an ODE-based decoder to improve the  
102 fidelity of the mel spectrograms. Then, F5-TTS [9] lever-  
103 ages the Diffusion Transformer with ConvNeXt V2 [63] to

104 better tackle text-speech alignment during in-context learn-  
105 ing. However, these works are limited in the field of TTS  
106 and cannot be applied to the V2C task. Therefore, we study  
107 the integration with MLLM reasoning models and TTS for  
108 the V2C task.

### 109 2.3. Chain-of-thought Reasoning

110 Visual reasoning demands the model’s visual perception  
111 capability and high-level cognition ability [34, 44]. Sev-  
112 eral tasks have been applied to evaluate the visual reason-  
113 ing ability of Visual-Language Models (VLMs), including  
114 VQA [32, 40] requiring models to answer visual content  
115 and textual questions, and Visual Entailment [1, 12, 58]  
116 requiring models to determine the consistency of text de-  
117 scriptions and visual content, etc. With the development of  
118 LLMs, vision-language models leverage the advanced rea-  
119 soning abilities of LLMs to interpret visual tasks [42, 67].  
120 Some vision-language models enhance visual reasoning by  
121 optimizing the visual encoding strategy [22, 26, 31, 42, 68,  
122 72] to produce cognition-focused visual tokens. Then, with  
123 the rapid advancement of large language reasoning models  
124 with the step-by-step thinking ability [2, 19, 47–49, 52, 64],  
125 vision-language task is studied through step-step reasoning  
126 for a variety of multimodal large language models. How-  
127 ever, step-step reasoning mechanism is not well-studied in  
128 the movie dubbing, therefore, we propose DeepDubbber for  
129 movie dubbing with internal multimodal chain-of-thoughts  
130 reasoning guidance.

## 131 3. Method

### 132 3.1. Overview

133 Given a silent video clip  $V_l$ , a corresponding subtitle  $T_v$ ,  
134 and the goal of generating a fully dubbed video, the pro-  
135 posed model (DeepDubbber) aims to produce speech  $\hat{S}$  that  
136 matches the video, ensures contextual and prosodic rele-  
137 vance, and maintains speech-video synchronization with the  
138 help of MLLM. The model can be formalized as follows:

$$\hat{S} = F_{dubbber}(V_l, T_v) \quad (1)$$

139 DeepDubbber consists of two modeling stages: i) Multi-  
140 modal reasoning and understanding through in-context  
141 learning and mixed preference optimization. ii) Speech  
142 Generation Stage: This stage incorporates a conditional DiT  
143 based speech generator.

### 144 3.2. Multi-modal Chain-of-Thought Learning via 145 MLLM

#### 146 3.2.1. Stage 1.1: Training Multi-modal Chain-of- 147 Thought via Supervised Learning.

148 The core functionality of DeepDubbber is to extract key se-  
149 mantic features from the visual stream that are crucial for  
150 the dubbing process. These features include scene type,

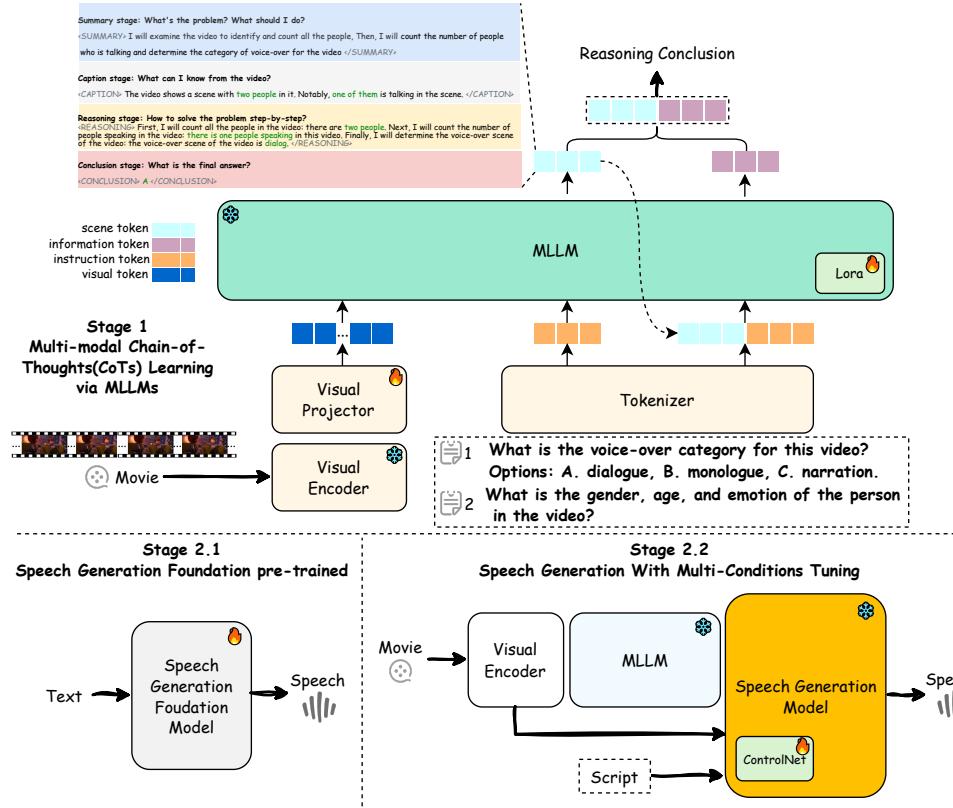


Figure 2. DeepDubber pipeline with multi-stage, multi-modal training.

speaker gender, speaker age, and speaker emotion, which are inferred step by step, as illustrated in Figure 4. Inspired by the success of MLLM, such as VLMs [42, 62, 65], we leverage multimodal instruction tuning to enhance CoT reasoning, thereby improving the quality of movie dubbing. The CoT reasoning process is formulated as follows:

$$C_1^{CoT} = F_{mcot}^1(V, \text{Instruct}_1), \quad (2)$$

$$C_2^{QA} = F_{qa}^2(V, \text{Instruct}_2) \quad (3)$$

where  $V$  represents the input video clip,  $\text{Instruct}_i$  denotes the  $i$ -th instruction provided to guide the reasoning process,  $F_{mcot}^i$  is the function representing the  $i$ -th step of multimodal CoT reasoning using an MLLM,  $C_1^{CoT}$  and  $C_2^{QA}$  are intermediate outputs from the reasoning process, where  $C_1^{CoT}$  is the result of the first CoT reasoning step, and  $C_2^{QA}$  is obtained through a question-answering (QA) step.

To optimize the response generation of the model, we define the multimodal reasoning process as follows:

$$M_{mllm} : (\text{Video}_{clip}, \text{CoT}_{instruction}, \text{QA}_{instruction}) \mapsto \text{Response}, \quad (4)$$

$$\min_{\theta_{\text{response}}} \mathbb{E}_{(\text{Video}_{clip}, \text{CoT}_{instruction}, \text{QA}_{instruction}) \sim \mathcal{D}} \left[ \mathcal{L}_{res} \left( M_{mllm}(\text{Video}_{clip}, \text{CoT}_{instruction}, \text{QA}_{instruction}), \text{Response}_{gt} \right) \right], \quad (5)$$

where  $M_{mllm}$  represents the multimodal large language model performing CoT reasoning and QA,  $\text{Video}_{clip}$  is the input video segment,  $\text{CoT}_{instruction}$  and  $\text{QA}_{instruction}$  are instructions guiding the CoT reasoning and QA process, respectively,  $\text{Response}$  is the generated output of the model,  $\mathcal{D}$  denotes the distribution of the training dataset,  $\theta_{\text{response}}$  represents the parameters of the model to optimize,  $\mathcal{L}_{res}$  is the loss function measuring the difference between the predicted response of the model and the ground truth response  $\text{Response}_{gt}$ . This approach ensures that the MLLM effectively reasons over multimodal inputs, facilitating high-quality dubbing through structured stepwise reasoning.

### 3.2.2. Stage 1.2: Training Multi-modal Chain-of-Thought via Reinforcement Learning.

The reward is the source of the training signal that decides the direction of RL optimization. To train CoT-MLLM, we adopt a rule-based reward system like Deepseek-R1 [19] that consists mainly of two types of rewards: accuracy rewards and format rewards. We employ a format reward model that enforces the model to put its reasoning process between `<SUMMARY> </SUMMARY>`, `<CAPTION> </CAPTION>`, `<REASONING> </REASONING>`, and `<CONCLUSION> </CONCLUSION>` tags. We use the Mixed Preference Op-

172  
173  
174  
175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195

timization (MPO) [62] method to learn the relative preferences between pairs of responses and enhance the reasoning capability of MLLM across different instructions. The mixed preference optimization is applied to further enhance the ability of the multimodal CoT reasoning. The training objective is represented as the following.

**Training Objective.** The MPO objective combines three loss components and F&O rewards: preference loss ( $L_p$ ), quality loss ( $L_q$ ), generation loss ( $L_g$ ), format loss ( $L_f$ ), and accuracy loss ( $L_o$ ). The total loss is formulated as:

$$L = w_p L_p + w_q L_q + w_g L_g + w_f L_f + w_c L_c, \quad (6)$$

where  $w_*$  represents the weight for each loss. We use DPO [55] for preference loss and BCO [10] for quality loss. The details of three terms of the loss are then represented as the following:

**Preference Loss.** The DPO [55] loss models the relative preference between the chosen and rejected responses without requiring a reward model. The loss function is:

$$L_p = -\log \sigma \left( \beta \log \frac{\pi_\theta(y_c | x)}{\pi_0(y_c | x)} - \beta \log \frac{\pi_\theta(y_r | x)}{\pi_0(y_r | x)} \right), \quad (7)$$

where  $\beta$  is the KL penalty coefficient,  $x$  is the user query,  $y_c$  is the chosen response,  $y_r$  is the rejected response, and  $\pi_\theta$  is the policy model initialized from  $\pi_0$ .

**Quality Loss.** The BCO [10] loss measures the absolute quality of individual responses using a binary classifier. The total loss is:

$$L_q = L_q^+ + L_q^-, \quad (8)$$

where the chosen and rejected loss terms are:

$$L_q^+ = -\log \sigma \left( \beta \log \frac{\pi_\theta(y_c | x)}{\pi_0(y_c | x)} - \delta \right), \quad (9)$$

$$L_q^- = -\log \sigma \left( -\left( \beta \log \frac{\pi_\theta(y_r | x)}{\pi_0(y_r | x)} - \delta \right) \right), \quad (10)$$

and  $\delta$  is the reward shift for stabilizing training.

**Generation Loss.** The SFT loss helps the model learn to generate preferred responses. The loss is defined as:

$$L_g = -\frac{\log \pi_\theta(y_c | x)}{|y_c|}. \quad (11)$$

**Format Reward.** The Format loss helps the model to learn to generate in preferred format. The loss is defined as:

$$L_f = -\sum [f_{\text{true}} \log(p_f) + (1 - f_{\text{true}}) \log(1 - p_f)] \quad (12)$$

$f_{\text{true}} \in \{0, 1\}$ : format correct or not,  $p_f$ : the probability of the format being correct predicted by the model.

**Outcome Reward.** The Accuracy loss helps the model learn to generate preferred answer. The loss is defined as:

$$L_o = -\sum [o_{\text{true}} \log(p_o) + (1 - o_{\text{true}}) \log(1 - p_o)] \quad (13)$$

$o_{\text{true}} \in \{0, 1\}$ : the answer is correct or not,  $p_o$ : the probability of the format being correct predicted by the model.

### 3.3. Multi-Conditioned Speech Generation

#### 3.3.1. Speech Generation Foundation Pre-Training

In the second stage of DeepDubber, we first train the foundational speech generation model. To optimize this process, we aim to learn the parameters  $\theta_{\text{generation}}$  by minimizing the composite Conditional Flow Matching (CFM) loss  $\mathcal{L}_{\text{cfm}}$ . The speech generation process is formulated as:

$$\begin{aligned} M_{\text{speech}} : & (\text{Video}_\text{clip}, \text{Speech}_\text{prompt}, \text{Caption}_\text{condition}, \text{Transcript}_\text{text}) \\ & \mapsto \text{Speech}_\text{target}, \end{aligned} \quad (14) \quad 248$$

$$\begin{aligned} \min_{\theta_{\text{generation}}} \mathbb{E}_{(\text{Video}_\text{clip}, \text{Speech}_\text{prompt}, \text{Caption}_\text{condition}, \text{Transcript}_\text{text}) \sim \mathcal{D}} \\ & [\mathcal{L}_{\text{cfm}}(M_{\text{speech}}(\text{Video}_\text{clip}, \text{Speech}_\text{prompt}, \\ & \text{Caption}_\text{condition}, \text{Transcript}_\text{text}))]. \end{aligned} \quad (15) \quad 249$$

We adopt the same architecture as F5-TTS [9], which employs a diffusion transformer (DiT) as the backbone. The model is trained to output a vector field  $v_\tau$  using the CFM objective  $\mathcal{L}_{\text{cfm}}$  [41], defined as:

$$\mathcal{L}_{\text{cfm}}(\theta) = \mathbb{E}_{\tau, q(x_1), p(x|x_1)} \|u_\tau(x|x_1) - v_\tau(x; \theta)\|^2 \quad (16) \quad 254$$

where  $p_\tau$  represents the probability path at time  $\tau$ ,  $u_\tau$  is the designated vector field for  $p_\tau$ ,  $x_1$  is a random variable corresponding to the training data,  $q(x_1)$  denotes the distribution of the training data. By optimizing  $\mathcal{L}_{\text{cfm}}$ , the model learns to generate high-quality speech synchronized with the visual and textual cues, ensuring natural and contextually appropriate dubbing.

#### 3.3.2. Speech Generation With Multi-Conditions Tuning

Next, in the ControlNet-transformer tuning stage, the video frames, along with an instruction are fed as inputs to the MLLM model. The sequence of video features and the video understanding conclusion are then combined and passed into the speech generation model. In this context, the provided conclusion helps guide the V2S generation process, as shown in stage 2.2 of Figure 2. The proposed speech generation model takes as input the script, silent video, video understanding conclusion, and an optional reference speech, and generates video-aligned speech context sequences, which can be described as:

$$\hat{S} = F_{\text{generator}}(V_l, C_v \{C_1, \dots, C_n\}_{n=4}, T_v) \quad (17) \quad 274$$

where  $C_v$  is the combination of  $\{C_1, \dots, C_n\}$ ,  $n = 4$ , which represents the scene type condition  $C_s$ , speaker gender condition  $C_g$ , speaker age condition  $C_a$ , and speaker emotion condition  $C_e$ . These conditions are combined and encoded by an encoder [56]. The  $V_l$  represents the visual features derived from the input video frames, which are encoded by CLIP [53]. We implement a cross-attention mechanism that facilitates the integration of understanding conclusion features  $C_v$  and visual features  $V_l$ . Furthermore,

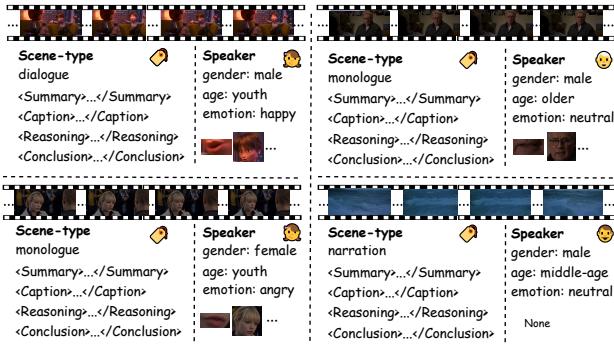


Figure 3. Proposed dataset with multi-type annotations, including annotation for lips, faces, scene-type, speaker gender, speaker age, voice emotion.

284  $T_v$  represents the embedded script. Furthermore, we added  
 285 a duration loss  $\mathcal{L}_{dur}$  to constrain the duration consistency,  
 286 which can be described as:

$$287 \quad \mathcal{L}_{dur} = \ell(f(V_l, C_l), dur), \quad (18)$$

288 The final loss function is constructed as follows:

$$289 \quad \mathcal{L}_g = \mathbb{E}_{\tau, q(x_1), p(x|x_1)} \|u_\tau(x|x_1, v_l, c_v, t_v) - v_\tau(x; \theta)\|^2 + \mathcal{L}_{dur} \quad (19)$$

290 In the training stage, the visual condition  $V_l$ , video under-  
 291 standing conclusion condition  $C_v$ , and video script condition  
 292  $T_v$  are each set to  $\phi$  with a 5 % probability. Extending  
 293 classifier-free guidance from the script condition to visual  
 294 input and visual understanding enhances both conditional  
 295 control precision and speech quality. The guidance scales  
 296  $\lambda_V$ ,  $\lambda_C$ , and  $\lambda_T$ , correspond to the video clip, video con-  
 297 clusion, and video-related script, respectively, and measure  
 298 the alignment between the sampling results and conditions.  
 299 Inspired by [33], during inference, the modified velocity es-  
 300 timate is as follows:

$$301 \quad g'_\theta = g_\theta(x_\tau, \phi, \phi, \phi) \\ + \lambda_V \cdot (v_0(x_\tau, c_v, c_c, c_t) - v_0(x_\tau, \phi, c_c, c_t)) \\ + \lambda_C \cdot (v_0(x_\tau, \phi, c_c, c_t) - v_0(x_\tau, \phi, \phi, c_t)) \\ + \lambda_T \cdot (v_0(x_\tau, \phi, \phi, c_t) - v_0(x_\tau, \phi, \phi, \phi)) \quad (20)$$

## 303 4. Experiments

### 304 4.1. Datasets

305 **Emilia** [24] is a comprehensive multilingual speech gener-  
 306 ation dataset containing a total of 101,654 hours of speech  
 307 data across six languages. The English portion of this  
 308 dataset, comprising approximately 46,800 hours, is utilized  
 309 to train our foundational TTS model.

310 **V2C-Animation** [6] is a specialized dataset designed for  
 311 animated movie dubbing, consisting of 10,217 clips from 26

films with synchronized text, speech, and video. The dataset  
 312 is partitioned into 60% for training, 10% for validation, and  
 313 30% for testing.

314 **GRID** is a dubbing benchmark for multi-speaker dubbing  
 315 [18]. The whole dataset has 33 speakers, each with 1000  
 316 short English samples. All participants are recorded in stu-  
 317 dio with unified background. The number of train and test  
 318 data are 32,670 and 3280, respectively.

319 **Chian-of-Thought Movie Dubbing Dataset.** We build a  
 320 7.2 hour multimodal CoT movie dubbing dataset for gen-  
 321 erating high-quality and accurate movie dubbing. Based  
 322 on CoT reasoning and CoT-like guidance [65], we uti-  
 323 lize a professional annotation team to label the follow-  
 324 ing dataset. We develop a CoT reasoning framework  
 325 to guide subsequent movie dubbing tasks, as illustrated  
 326 in Figures 3 and 4. Specifically, a step-by-step in-  
 327 struction process with video input is designed to en-  
 328 able efficient and accurate movie scene type classifica-  
 329 tion. As shown in Figure 4, <SUMMARY></SUMMARY>  
 330 provides a high-level overview of the entire scene, while  
 331 <CAPTION></CAPTION> describes the characters in the  
 332 video. During the <REASONING></REASONING> stage,  
 333 the reasoning process is divided into four steps:  
 334

**Step 1.** Count the numbers of people in the video.

**Step 2.** Distinguish whether the people in the video are  
 335 talking or not.

**Step 3.** Distinguish whether the movie contains dia-  
 336logue, narration, or monologue.

**Step 4.** Conclusion and give the answer.

337 And then <CONCLUSION></CONCLUSION> stage  
 338 give the final answer. Each stage is initiated at the model's  
 339 discretion, without external prompt engineering frame-  
 340 works or additional prompting. Specifically, we provide the  
 341 model with four pairs of special tags, these tags correspond  
 342 to summarizing the response approach, describing relevant  
 343 image content, conducting reasoning, and preparing a final  
 344 answer, respectively. The proposed dataset consists of two  
 345 difficulty levels: (1) Level-1, where people are talking in  
 346 the videos, includes 7,276 video clips for training and 1,100  
 347 video clips for testing. (2) Level-2, where animals are talk-  
 348 ing in the videos, includes 3,486 video clips for training and  
 349 388 video clips for testing. Notably, due to OpenAI's fine-  
 350 tuning policy, all Level-1 video clips have been filtered out,  
 351 and 328 video clips remain for Level-2 training.

### 352 4.2. Evaluation Metrics

353 We evaluate using both objective and subjective metrics.  
 354 To assess pronunciation accuracy, we use Word Error Rate  
 355 (WER) with Whisper-V3 [54] as the ASR model. Tim-  
 356 bre consistency is evaluated with speaker encoder cosine  
 357 similarity (SPK-SIM) [17]. We also calculate mel cepstral  
 358 distortion dynamic time warping (MCD) and speech length  
 359 variance (MCD-SL) [3] for spectral and length differences.



Figure 4. The reasoning stages of movie scene type CoT annotations.

**Table 1. Objective evaluation of the initial reasoning setting.** For speech generation setting, we use the target speaker’s speech as voice prompt if the predict scene type is correct and use random speaker’s speech as voice prompt if the predict scene type is not correct.

Models Name	Scores on dialogue(A), monologue(B) and narration(c)					Speech Generation				
	Ave.Acc(%) ↑	Ave.Recall(%) ↑	A.Recall(%) ↑	B.Recall(%) ↑	C.Recall(%) ↑	SPK-SIM(%) ↑	WER(%) ↓	MCD ↓	MCD-SL ↓	
MLLMs based										
Qwen [51]	MMLM-1B [27]	84.09	82.97	86.50	68.40	94.00	83.17	23.60	8.59	8.60
	MMLM-4B [29]	81.73	80.98	83.33	<b>75.20</b>	84.40	83.34	23.41	8.53	8.53
InternLM [5]	MMLM-2B [28]	84.18	81.23	<b>90.50</b>	59.20	94.00	82.97	23.20	8.58	8.60
	MMLM-8B [30]	<b>86.00</b>	<b>85.84</b>	86.33	73.20	<b>98.00</b>	<b>83.42 (+30.28%)</b>	<b>23.20 (+55.70%)</b>	<b>8.54 (+0.93%)</b>	<b>8.54 (+3.94%)</b>
Dubbing Models										
HPMDubbing [14]	-	-	-	-	-	61.06	199.40	8.82	11.88	
Speaker2Dub [69]	-	-	-	-	-	61.73	84.42	8.75	10.78	
StyleDubber [17]	-	-	-	-	-	64.03	52.69	8.62	8.89	

364 Emotion similarity (EMO-SIM) is assessed using a speech  
365 emotion recognition model [66]. For alignment with video,  
366 we use Lip Sync Error Distance (LSE-D) and Lip Sync Er-  
367 rror Confidence (LSE-C) metrics on the Grid benchmark,  
368 based on the pre-trained SyncNet model [13]. For sub-  
369 jective evaluation, we conduct human evaluations of the Mean  
370 Opinion Score (MOS) for naturalness (NMOS) and similar-  
371 ity (SMOS), rated on a 1-to-5 scale with 95% confidence  
372 intervals. Following [69], participants evaluate the dubbing  
373 quality of 30 randomly selected speech samples from each  
374 test set.

### 375 4.3. Benchmark Results

376 We compare our approach with a TTS model [9] and  
377 three recent video dubbing models. HPMDubbing [14] in-  
378 troduces an emotional prosody adaptor that enables fine-  
379 grained alignment of the speaker’s emotions. StyleDubber  
380 [17], on the other hand, designs a multimodal phoneme-  
381 level style adaptor that generates stylized voice tones based  
382 on facial expressions. Speaker2Dubber [69] combines char-

acter emotions, phoneme prosody, and lip movements to en-  
383 sure consistency in both prosody and duration throughout  
384 the dubbing process.

**385 Results on movie scene type reasoning and speech gener-  
386 ation.** As shown in Table 1, the MMLM-8B achieves su-  
387 perior performance across all benchmarks in the classification  
388 of movie scene types. Our method outperforms the SOTA  
389 dubbing methods (StyleDubber and Speaker2Dub) on SPK-  
390 SIM, WER and MCD/MCD-SL. In detail, SPK-SIM im-  
391 proved from 64. 03% to 70. 83%, WER decreased from 52.  
392 69% to 27. 68%. And, as shown in 3, the MMLM-8B main-  
393 tains the competitive performance which slightly lower than  
394 GPT-4o [46]. These results demonstrate the effectiveness of  
395 our multimodal reasoning stages in enhancing multimodal  
396 movie dubbing performance.

**397 Results on V2C-Animation benchmark.** As shown in Ta-  
398 ble 3, compared to the state-of-the-art models [14, 17, 69],  
399 our model achieves improvements across evaluation metrics  
400 in the same setting [69]. Our model achieves the best per-  
401 formance across all metrics. In detail, SPK-SIM increased

**Table 2. Ablation of objective evaluation under the initial reasoning setting on the proposed dataset.** For the speech generation setting, we use the target speaker’s speech as a voice prompt if the predicted scene type is correct, and a random speaker’s speech if the predicted scene type is incorrect. F&O Reward refers to format reward and outcome reward.

Methods	Setting	Scores on dialogue(A), monologue(B) and narration(c)					Speech Generation		
		Ave.Acc(%) ↑	Ave.Recall(%) ↑	A.Recall(%) ↑	B.Recall(%) ↑	C.Recall(%) ↑	SPK-SIM(%) ↑	WER(%) ↓	MCD ↓
<b>Base Model</b>									
InternVL2.5-8B [30]	QA	39.45%	44.24%	28.83%	5.20%	99.20%	81.73%	24.82%	8.91
<b>Our Models</b>									
SFT[30]	QA	82.27%	79.13%	89.00%	50.00%	<b>98.40%</b>	82.89%	23.52%	8.66
	Reasoning	85.09%	82.10%	<b>91.50%</b>	56.80%	98.00%	83.34%	23.41%	8.55
MPO [61]	RL	84.00%	81.36%	89.67%	56.40%	98.00%	83.02%	23.49%	8.59
MPO + F&O Reward [19]	RL	<b>86.00% (+4.53%)</b>	<b>85.84% (+8.50%)</b>	86.33%	<b>73.20%</b>	98.00%	<b>83.42% (+0.64%)</b>	<b>23.20% (+1.36%)</b>	<b>8.54 (+1.39%)</b>
<b>MPO + F&amp;O Reward [19]</b>									

**Table 3. Objective evaluation of the Reasoning Stage** based on scores in dialogue (A), monologue (B), and narration (C) under two different levels, as explained in 4.1.

Model Name	Setting1					Setting2				
	Ave.Acc(%) ↑	Ave.Recall(%) ↑	A.Recall(%) ↑	B.Recall(%) ↑	C.Recall(%) ↑	Ave.Acc(%) ↑	Ave.Recall(%) ↑	A.Recall(%) ↑	B.Recall(%) ↑	C.Recall(%) ↑
Closed Source										
GPT-4o[46]	-	-	-	-	-	<b>73.97</b>	<b>64.58</b>	<b>91.11</b>	44.68	57.97
Open Source										
Qwen [51]	MMLM-1B [27]	84.09	82.97	86.50	68.40	94.00	61.86	53.78	80.44	12.77
	MMLM-4B [29]	81.73	80.98	83.33	<b>75.20</b>	84.40	62.63	55.51	79.56	15.96
InternLM [5]	MMLM-2B [28]	84.18	81.23	<b>90.50</b>	59.20	94.00	53.09	51.37	61.33	15.96
	MMLM-8B [30]	<b>86.00</b>	<b>85.84</b>	86.33	73.20	<b>98.00</b>	67.53	59.29	79.56	<b>60.64</b>
<b>InternLM [5]</b>										

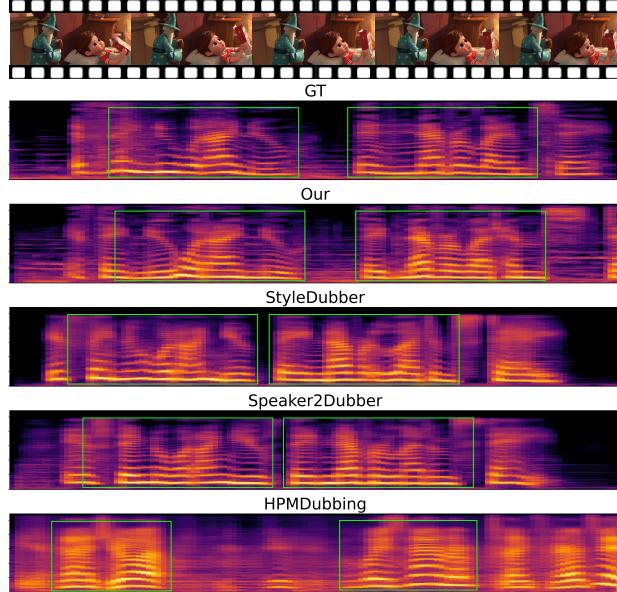


Figure 5. Visualization of speech samples generated by state-of-the-art models and our. The green rectangles highlight key regions that have significant differences in overall expressiveness.

from 79.81% to 83.30%, EMO-SIM improved from 59.71% to 64.93%, MCD decreased from 9.11 to 8.80, and WER decreased from 26. 48% to 24. 71%. It shows that our framework achieves performance improvement in pronunciation accuracy and consistency of speech duration.

**Results on GRID benchmark.** As shown in Table 4, our model achieves the best lip-sync performance on the GRID

benchmark with the same evaluation of the state-of-the-art models [69], which decreased from 14.79 to 14.63. And MCD decreased from 4.95 to 4.73. Unlike V2C-Animation, samples in GRID are recorded in a studio environment, which does not involve exaggerated prosody variations or background noise. As a result, the WER of all comparison methods is generally better on the GRID compared to V2C-Animation. As shown in Table 4, our model achieves the best lip-sync performance on the GRID benchmark. In addition, our method also achieves competitive results in WER, slightly lower than the best fine-tuned F5-TTS model, Speaker2Dub and StyleDubber. However, these models have WER results (11.94%, 12.11% and 11.97%) that exceed the ground-truth WER result (13. 67%), suggesting that the intelligibility has reached an acceptable range for humans.

**Results on Speaker Zero-shot test.** As shown in Table 6, this setting uses the speech of unseen speakers as reference speech to measure the generalizability of the dubbing model [69]. Here, we use the speech from GRID as reference speech to measure V2C. We compare LSE-C/D, SPK-SIM, and WER, along with subjective evaluations in the same evaluation setting with the state-of-the-art models [69]. As shown in Table 6, our method outperforms Style-Dubber and Speaker2Dub in both SPK-SIM and WER. In detail, the SIP-SIM improved from 78.30% to 83.55%, the WER decreased from 16.57% to 15.49%. Furthermore, the proposed method still maintains competitive performance in speech-visual synchronization (see LSE-C and LSE-D), slightly lower than HPMDubbing.

**Table 4. Objective results on V2C-Animation and Grid benchmark.** For the Dub 1.0 setting, we use the ground truth speech as reference speech, for the Dub 2.0 setting, we use the non-ground truth speech from the same speaker within the dataset as the reference speech which is more aligned with practical usage in dubbing.

benchmark	Setting	Dub 1.0								Dub2.0						
		Methods	Visual	SPK-SIM(%) ↑	WER(%) ↓	EMO-SIM(%) ↑		MCD ↓	MCD-SL ↓	SPK-SIM(%) ↑	WER(%) ↓	EMO-SIM(%) ↑		MCD ↓	MCD-SL ↓	
			GT	-	100.00	17.38	100.00	0.00	0.00			100.00	17.38	100.00	0.00	
V2C	F5-TTS [9]	x		89.3	24.41	76.78	8.32	8.32	83.11	24.83	64.91	10.86	10.87			
	HPMDubbing [14]	✓		73.64	151.02	39.85	8.59	8.32	73.01	150.83	34.69	9.11	12.15			
	Speaker2Dub [69]	✓		82.15	31.23	65.92	10.68	11.21	79.53	31.28	59.71	11.16	11.70			
	StyleDubber [17]	✓		82.48	27.36	66.24	10.06	10.52	79.81	26.48	59.08	10.56	11.05			
	Ours	✓		<b>89.74</b>	<b>22.51</b>	<b>78.88</b>	<b>6.98</b>	<b>6.99</b>	<b>83.30</b>	<b>24.71</b>	<b>64.93</b>	<b>8.80</b>	<b>8.80</b>			
	Methods	Visual	SPK-SIM(%) ↑	WER(%) ↑	LSE-C ↓	LSE-D ↓	MCD ↓	MCD-SL ↓	SPK-SIM(%) ↑	WER(%) ↑	LSE-C ↓	LSE-D ↓	MCD ↓	MCD-SL ↓		
Grid	GT	-	100.00	13.67	7.18	13.36	0.00	0.00	100.00	13.67	7.18	13.36	0.00	0.00		
	F5-TTS [9]	x	<b>96.51</b>	<b>11.94</b>	5.51	14.70	<b>4.23</b>	<b>4.24</b>	94.45	16.75	5.10	14.71	4.89	4.90		
	HPMDubbing [14]	✓	93.64	16.78	<b>6.35</b>	14.78	4.57	4.85	92.84	17.40	<b>6.34</b>	14.79	4.95	5.24		
	Speaker2Dub [69]	✓	96.11	12.11	5.64	14.82	7.85	8.01	94.91	12.89	5.56	14.84	7.57	7.73		
	StyleDubber [17]	✓	96.40	11.97	6.19	14.81	7.71	7.81	<b>95.25</b>	<b>11.97</b>	6.16	14.83	7.34	7.43		
	Ours	✓	95.73	14.71	4.87	<b>14.63</b>	4.48	4.49	94.71	16.08	4.46	<b>14.63</b>	<b>4.73</b>	<b>4.74</b>		

**Table 5. Subjective evaluation on V2C-Animation and GRID benchmarks.**

Dataset	V2C-Animation		GRID	
	NMOS ↑	SMOS ↑	NMOS ↑	SMOS ↑
Methods				
GT	4.98±0.01	-	4.99±0.01	-
F5-TTS [9]	4.20±0.68	3.83±0.63	<b>4.43±0.03</b>	<b>3.32±0.05</b>
HPMDubbing [14]	1.04±0.01	1.02±0.01	3.50±0.10	2.77±0.12
Speaker2Dub [69]	2.93±0.21	2.58±0.19	4.04±0.07	3.00±0.10
StyleDubber [17]	2.68±0.21	2.39±0.21	4.01±0.03	3.06±0.07
Ours	<b>4.37±0.35</b>	<b>3.91±0.45</b>	4.33±0.07	3.14±0.08

**Table 6. Results on zero-shot test**, which use unseen speaker as reference speech.

Setting	Dubbing Setting 3.0				
	LSE-C ↑	LSE-D ↓	SPK-SIM (%) ↑	EMO-SIM (%) ↑	MCD ↓
Methods					
HPMDubbing [14]	1.72	<b>11.74</b>	68.14	126.85	1.29±0.60
Speaker2Dub [69]	2.21	12.67	76.10	16.57	3.38±0.14
StyleDubber [17]	2.15	12.76	78.30	19.07	3.30±0.15
Ours	<b>2.21</b>	12.59	<b>83.55</b>	<b>15.49</b>	<b>4.12±0.16</b>

**Table 7. Results of ablation study on the proposed dataset with 2.0 setting.**

	LSE-C ↑	LSE-D ↓	SPK-SIM (%) ↑	EMO-SIM (%) ↑	MCD ↓	MCD-SL ↓
w/o Clip	1.99	12.73	82.63	63.24	8.85	8.86
w/o Dur	<b>2.06</b>	12.81	82.71	64.32	8.82	8.83
w/o Conclusion	1.89	12.66	82.59	63.18	8.81	8.82
Proposed	2.01	<b>12.61</b>	<b>82.99</b>	<b>64.74</b>	<b>8.76</b>	<b>8.77</b>

#### 4.4. Ablation Studies

**Ablation Studies on Reasoning Stages.** To compare the impact of SFT, MPO and MPO with F&O rewards on improving multimodal reasoning ability, we used constructed CoT and QA pairs as training data to fine-tune InternVL-8B. As shown in Table 2, the results indicate that the model trained with MPO with F&O rewards consistently outper-

forms that trained with Zeroshot, SFT and MPO. For example, the MPO (with F&O rewards) trained model achieves an acc of 86.00% on the movie scene reasoning benchmark, surpassing its SFT (QA) counterpart by 4.53%. Furthermore, the MPO (with F&O rewards) trained model also performs better on the recall rate of each category.

**Ablation studies on Speech Generation.** The ablation results in Table 7 indicate that each condition contributes to overall performance. Removing the video clip control causes all metrics to drop significantly, highlighting its importance for speech-video alignment. Adding the video understanding conclusion control improves SPK-SIM and EMO-SIM. Furthermore, removing the duration predictor results in the largest drop in LSE-D performance, emphasizing that learning duration-level consistency is crucial for synchronizing speech and video. Additionally, with the MMLM conclusion conditions, SPK-SIM and EMO-SIM increase by 0.4% and 1.56%, respectively, demonstrating the effectiveness of multimodal reasoning stages in enhancing multimodal movie dubbing performance.

## 5. Conclusion

In this paper, we propose a multi-stage, multimodal large language framework consisting of two-stage models and an accompanying multi-stage training strategy to improve the initial reasoning capabilities in movie dubbing. Additionally, we have created a corresponding movie dubbing dataset with CoT annotations. In the evaluation, the results show an improvement in performance compared to state-of-the-art methods across a variety of datasets.

## References

- [1] Saeed Amizadeh, Hamid Palangi, Oleksandr Polozov, Yichen Huang, and Kazuhito Koishida. Neuro-symbolic visual reasoning: disentangling “visual” from “reasoning”. In

- 480        *Proceedings of the 37th International Conference on Ma- 537  
481 chine Learning*. JMLR.org, 2020. 2 538  
482 [2] Anthropic. Claude 3.7, 2025. 1, 2 539  
483 [3] Eric Battenberg, RJ Skerry-Ryan, Soroosh Mariooryad, 540  
484 Daisy Stanton, David Kao, Matt Shannon, and Tom Bagby. 541  
485 Location-relative attention mechanisms for robust long-form 542  
486 speech synthesis, 2020. 5 543  
487 [4] Mirco Bonomo and Simone Bianco. Visual rag: Expanding 544  
488 multimodal visual knowledge without fine-tuning, 2025. 1 545  
489 [5] Zheng Cai and Maosong Cao et.al. Internlm2 technical 546  
490 report, 2024. 6, 7 547  
491 [6] Qi Chen, Yuanqing Li, Yuankai Qi, Jiaqiu Zhou, Mingkui 548  
492 Tan, and Qi Wu. V2c: Visual voice cloning, 2021. 5 549  
493 [7] Qi Chen, Mingkui Tan, Yuankai Qi, Jiaqiu Zhou, Yuanqing 550  
494 Li, and Qi Wu. V2c: Visual voice cloning. In *Proceedings of 551  
495 the IEEE/CVF Conference on Computer Vision and Pattern 552  
496 Recognition*, pages 21242–21251, 2022. 1 553  
497 [8] Tian Qi Chen, Yulia Rubanova, Jesse Bettencourt, and David 554  
498 Duvenaud. Neural ordinary differential equations. In 555  
499 *NeurIPS*, pages 6572–6583, 2018. 2 556  
500 [9] Yushen Chen, Zhihang Niu, Ziyang Ma, Keqi Deng, Chun- 557  
501 hui Wang, Jian Zhao, Kai Yu, and Xie Chen. F5-tts: A fairy- 558  
502 taler that fakes fluent and faithful speech with flow matching, 559  
503 2024. 2, 4, 6, 8 560  
504 [10] Zhipeng Chen, Kun Zhou, Wayne Xin Zhao, Jingyuan Wang, 561  
505 and Ji-Rong Wen. Low-redundant optimization for large lan- 562  
506 guage model alignment. *ArXiv*, abs/2406.12606, 2024. 4 563  
507 [11] Zesen Cheng, Sicong Leng, Hang Zhang, Yifei Xin, Xin 564  
508 Li, Guanzheng Chen, Yongxin Zhu, Wenqi Zhang, Ziyang 565  
509 Luo, Deli Zhao, and Lidong Bing. Videollama 2: Advancing 566  
510 spatial-temporal modeling and audio understanding in video- 567  
511 llms, 2024. 1 568  
512 [12] Minkyu Choi, Harsh Goel, Mohammad Osama, Yunhao 569  
513 Yang, Sahil Shah, and Sandeep Chinchali. Towards neuro- 570  
514 symbolic video understanding. In *Computer Vision – ECCV 571  
515 2024: 18th European Conference, Milan, Italy, September 572  
516 29–October 4, 2024, Proceedings, Part LXVIII*, page 573  
517 220–236, Berlin, Heidelberg, 2024. Springer-Verlag. 2 574  
518 [13] Joon Son Chung and Andrew Zisserman. Out of time: auto- 575  
519 mated lip sync in the wild. In *Computer Vision–ACCV 2016 576  
520 Workshops: ACCV 2016 International Workshops, Taipei, 577  
521 Taiwan, November 20–24, 2016, Revised Selected Papers, 578  
522 Part II 13*, pages 251–263. Springer, 2017. 6 579  
523 [14] Gaoxiang Cong, Liang Li, Yuankai Qi, Zhengjun Zha, Qi 580  
524 Wu, Wenyu Wang, Bin Jiang, Ming-Hsuan Yang, and Qing- 581  
525 ming Huang. Learning to dub movies via hierarchical 582  
526 prosody models, 2023. 1, 2, 6, 8 583  
527 [15] Gaoxiang Cong, Liang Li, Yuankai Qi, Zheng-Jun Zha, Qi 584  
528 Wu, Wenyu Wang, Bin Jiang, Ming-Hsuan Yang, and Qing- 585  
529 ming Huang. Learning to dub movies via hierarchical 586  
530 prosody models. In *Proceedings of the IEEE/CVF Confer- 587  
531 ence on Computer Vision and Pattern Recognition*, pages 588  
532 14687–14697, 2023. 1 589  
533 [16] Gaoxiang Cong, Jiadong Pan, Liang Li, Yuankai Qi, Yuxin 590  
534 Peng, Anton van den Hengel, Jian Yang, and Qingming 591  
535 Huang. Emodubber: Towards high quality and emotion con- 592  
536 trollable movie dubbing, 2024. 2

- 593 [33] Junpeng Jiang, Gangyi Hong, Lijun Zhou, Enhui Ma, Heng-  
594 tong Hu, xia zhou, Jie Xiang, Fan Liu, Kaicheng Yu,  
595 Haiyang Sun, Kun Zhan, Peng Jia, and Miao Zhang. DiVE:  
596 Dit-based video generation with enhanced control. In *ECCV  
597 2024 Workshop on Multimodal Perception and Comprehen-  
598 sion of Corner Cases in Autonomous Driving*, 2024. 5
- 599 [34] Justin Johnson, Bharath Hariharan, Laurens van der Maaten,  
600 Li Fei-Fei, C. Lawrence Zitnick, and Ross B. Girshick.  
601 Clevr: A diagnostic dataset for compositional language and  
602 elementary visual reasoning. *2017 IEEE Conference on  
603 Computer Vision and Pattern Recognition (CVPR)*, pages  
604 1988–1997, 2016. 2
- 605 [35] Kangsan Kim, Geon Park, Youngwan Lee, Woongyeong  
606 Yeo, and Sung Ju Hwang. Videoicl: Confidence-based it-  
607 erative in-context learning for out-of-distribution video un-  
608 derstanding, 2024. 1
- 609 [36] Sungwon Kim, Kevin J. Shih, Rohan Badlani, Joao Felipe  
610 Santos, Evelina Bakhturina, Mikyas T. Desta, Rafael Valle,  
611 Sungroh Yoon, and Bryan Catanzaro. P-flow: A fast and  
612 data-efficient zero-shot TTS through speech prompting. In  
613 *Thirty-seventh Conference on Neural Information Process-  
614 ing Systems*, 2023. 2
- 615 [37] Jiyoung Lee, Joon Son Chung, and Soo-Whan Chung. Imag-  
616 inary voice: Face-styled diffusion model for text-to-speech.  
617 In *ICASSP 2023-2023 IEEE International Conference on  
618 Acoustics, Speech and Signal Processing (ICASSP)*, pages  
619 1–5. IEEE, 2023. 1, 2
- 620 [38] Keon Lee, Dong Won Kim, Jaehyeon Kim, Seungjun Chung,  
621 and Jaewoong Cho. DiTTo-TTS: Diffusion transformers for  
622 scalable text-to-speech without domain-specific factors. In  
623 *The Thirteenth International Conference on Learning Rep-  
624 resentations*, 2025. 2
- 625 [39] Chengzu Li, Wenshan Wu, Huanyu Zhang, Yan Xia,  
626 Shaoguang Mao, Li Dong, Ivan Vulić, and Furu Wei. Imag-  
627 ine while reasoning in space: Multimodal visualization-of-  
628 thought, 2025. 1
- 629 [40] Hao Li, Xu Li, Belhal Karimi, Jie Chen, and Mingming Sun.  
630 Joint learning of object graph and relation graph for visual  
631 question answering. In *2022 IEEE International Conference  
632 on Multimedia and Expo (ICME)*, pages 01–06, 2022. 2
- 633 [41] Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maxim-  
634 ilian Nickel, and Matthew Le. Flow matching for genera-  
635 tive modeling. In *The Eleventh International Conference on  
636 Learning Representations*, 2023. 2, 4
- 637 [42] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee.  
638 Visual instruction tuning. In *Thirty-seventh Conference on  
639 Neural Information Processing Systems*, 2023. 2, 3
- 640 [43] Jingming Liu, Yumeng Li, Boyuan Xiao, Yichang Jian,  
641 Ziang Qin, Tianjia Shao, Yao-Xiang Ding, and Kun Zhou.  
642 Enhancing visual reasoning with autonomous imagination in  
643 multimodal large language models, 2024. 1
- 644 [44] Mikołaj Małkiński and Jacek Mańdziuk. A review of emerg-  
645 ing research directions in abstract visual reasoning. *Informa-  
646 tion Fusion*, 91:713–736, 2023. 2
- 647 [45] Shivam Mehta, Ruibo Tu, Jonas Beskow, Éva Székely, and  
648 Gustav Eje Henter. Matcha-tts: A fast tts architecture with  
649 conditional flow matching. In *ICASSP 2024 - 2024 IEEE  
650 International Conference on Acoustics, Speech and Signal  
651 Processing (ICASSP)*, pages 11341–11345, 2024. 2
- [46] OpenAI. Openai gpt-4o, 2024. 6, 7
- [47] OpenAI. Openai o1, 2024. 1, 2
- [48] OpenAI. Openai o1-min, 2024.
- [49] OpenAI. Openai o3 mini, 2025. 1, 2
- [50] Abhirama Subramanyam Penamakuri, Kiran Chhatre, and  
Akshat Jain. Audiopedia: Audio qa with knowledge, 2024.  
1
- [51] Qwen and An Yang et.al. Qwen2.5 technical report, 2025. 6,  
7
- [52] QwenLM. Qwq-32b, 2025. 1, 2
- [53] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya  
Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry,  
Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen  
Krueger, and Ilya Sutskever. Learning transferable visual  
models from natural language supervision. In *Proceedings  
of the 38th International Conference on Machine Learning,  
ICML 2021, 18-24 July 2021, Virtual Event*, pages 8748–  
8763. PMLR, 2021. 4
- [54] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman,  
Christine McLeavey, and Ilya Sutskever. Robust speech  
recognition via large-scale weak supervision, 2022. 5
- [55] Rafael Rafaïlov, Archit Sharma, Eric Mitchell, Christo-  
pher D Manning, Stefano Ermon, and Chelsea Finn. Direct  
preference optimization: Your language model is secretly a  
reward model. In *Thirty-seventh Conference on Neural In-  
formation Processing Systems*, 2023. 4
- [56] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee,  
Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and  
Peter J. Liu. Exploring the limits of transfer learning with  
a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21:  
140:1–140:67, 2020. 4
- [57] Zahraa Al Sahili, Ioannis Patras, and Matthew Purver. Fair-  
cot: Enhancing fairness in diffusion models via chain of  
thought reasoning of multimodal language models, 2024. 1
- [58] Haoyu Song, Li Dong, Weinan Zhang, Ting Liu, and Furu  
Wei. Clip models are few-shot learners: Empirical studies  
on vqa and visual entailment. In *Annual Meeting of the As-  
sociation for Computational Linguistics*, 2022. 2
- [59] Yang Song, Conor Durkan, Iain Murray, and Stefano Ermon.  
Maximum likelihood training of score-based diffusion mod-  
els. In *Advances in Neural Information Processing Systems*,  
2021. 2
- [60] Li Wan, Quan Wang, Alan Papir, and Ignacio Lopez Moreno.  
Generalized end-to-end loss for speaker verification. In *2018  
IEEE International Conference on Acoustics, Speech and  
Signal Processing (ICASSP)*, pages 4879–4883. IEEE, 2018.  
1
- [61] Weiyun Wang and Zhe Chen et.al. Enhancing the reason-  
ing ability of multimodal large language models via mixed  
preference optimization, 2024. 7
- [62] Weiyun Wang, Zhe Chen, Wenhui Wang, Yue Cao,  
Yangzhou Liu, Zhangwei Gao, Jinguo Zhu, Xizhou Zhu,  
Lewei Lu, Yu Qiao, and Jifeng Dai. Enhancing the reason-  
ing ability of multimodal large language models via mixed  
preference optimization. *ArXiv*, abs/2411.10442, 2024. 3, 4

- 707 [63] Sanghyun Woo, Shoubhik Debnath, Ronghang Hu, Xinlei  
708 Chen, Zhuang Liu, In So Kweon, and Saining Xie. Convnext  
709 v2: Co-designing and scaling convnets with masked autoen-  
710 coders. In *2023 IEEE/CVF Conference on Computer Vision*  
711 and Pattern Recognition (CVPR), pages 16133–16142, 2023.  
712 2
- 713 [64] xAI. grok-3, 2025. 1, 2
- 714 [65] Guowei Xu, Peng Jin, Hao Li, Yibing Song, Lichao Sun, and  
715 Li Yuan. Llava-cot: Let vision language models reason step-  
716 by-step, 2025. 3, 5
- 717 [66] Jiaxin Ye, Xin-Cheng Wen, Yujie Wei, Yong Xu, Kunhong  
718 Liu, and Hongming Shan. Temporal modeling matters:  
719 A novel temporal emotional modeling approach for speech  
720 emotion recognition. In *ICASSP 2023-2023 IEEE Interna-*  
721 *tional Conference on Acoustics, Speech and Signal Process-*  
722 *ing (ICASSP)*, pages 1–5. IEEE, 2023. 6
- 723 [67] Wangbo Yu, Chaoran Feng, Jiye Tang, Xu Jia, Li Yuan, and  
724 Yonghong Tian. Evagaussians: Event stream assisted gaus-  
725 sian splatting from blurry images. *CoRR*, abs/2405.20224,  
726 2024. 2
- 727 [68] J.D. Zamfirescu-Pereira, Richmond Y. Wong, Bjoern Hart-  
728 mann, and Qian Yang. Why johnny can't prompt: How non-  
729 ai experts try (and fail) to design llm prompts. In *Proceed-  
730 ings of the 2023 CHI Conference on Human Factors in Com-  
731 puting Systems*, New York, NY, USA, 2023. Association for  
732 Computing Machinery. 2
- 733 [69] Zhedong Zhang, Liang Li, Gaoxiang Cong, Haibing YIN,  
734 Yuhan Gao, Chenggang Yan, Anton van den Hengel, and  
735 Yuankai Qi. From speaker to dubber: Movie dubbing with  
736 prosody and duration consistency learning. In *ACM Multi-  
737 media 2024*, 2024. 1, 2, 6, 7, 8
- 738 [70] Yuan Zhao, Zhenqi Jia, Rui Liu, De Hu, Feilong Bao, and  
739 Guanglai Gao. Mcdubber: Multimodal context-aware ex-  
740 pressive video dubbing. In *National Conference on Man-  
741 Machine Speech Communication*, pages 168–182. Springer,  
742 2024. 1, 2
- 743 [71] Haojie Zheng, Tianyang Xu, Hanchi Sun, Shu Pu, Ruoxi  
744 Chen, and Lichao Sun. Thinking before looking: Improv-  
745 ing multimodal llm reasoning via mitigating visual halluci-  
746 nation, 2024. 1
- 747 [72] Kunyang Zhou. LVP: Language-guide visual projector for  
748 efficient multimodal LLM, 2025. 2