



# MM-Narrator: Narrating Long-form Videos with Multimodal In-Context Learning

Chaoyi Zhang<sup>♦\*</sup> Kevin Lin<sup>♦</sup> Zhengyuan Yang<sup>♦</sup> Jianfeng Wang<sup>♦</sup> Linjie Li<sup>♦</sup>  
Chung-Ching Lin<sup>♦</sup> Zicheng Liu<sup>♦</sup> Lijuan Wang<sup>♦</sup>

<sup>♦</sup>University of Sydney <sup>♦</sup>Microsoft <sup>♦</sup>Advanced Micro Devices

<https://MM-Narrator.github.io>

Titanic (1997) | Start: 01:21:53 | End: 01:22:08



**Subtitles**  
> All right ...  
> ... open your eyes.  
> I'm flying.  
> Jack.

**Context AD**  
Jack and Rose stand together on the deck of the boat, enjoying their intimate moment. Jack and Rose embrace, looking out at the ocean, feeling like they're flying.

**AD Prediction** (via MM-Narrator)  
Jack and Rose hold hands, admiring the sunset together.

Friends S1 E1 -The One Where It All Began (1994) | Start: 00:04:06 | End: 00:04:16



**Subtitles**  
> This is everybody. This is Chandler and Phoebe... and Joey.  
> And remember my brother, Ross?  
> - Sure! - Hey.  
> - Hi. - Oh, God.

**Context AD**  
Ross listens to the conversation, occasionally sipping from his blue mug. Rachel stands before the group, introducing herself. She stands in her wedding dress, talking to Ross and the group.

**AD Prediction** (via MM-Narrator)  
Rachel, in her wedding dress, interacts with Ross and the group in the room.

The Matrix (1999) | Start: 00:29:28 | End: 00:29:33



**Subtitles**  
> ...you stay in Wonderland... and I show you how deep the rabbit hole goes.  
> Remember ... all I'm offering is the truth.  
Nothing more.

**Context AD**  
Neo hesitates before choosing between the red and blue pills. He closes his eyes, taking a moment to reflect before making his decision.

**AD Prediction** (via MM-Narrator)  
Neo finally makes his decision and picks up the red pill.

Spider-Man: Into the Spider-Verse (2018) | Start: 00:01:28 | End: 00:01:29



**Subtitles**  
> All right, let's do this one last time.  
> My name is Peter Parker.  
> I was bitten by a radioactive spider.  
> And for 10 years...

**Context AD**  
A close-up reveals a Spider-Man comic book cover. Peter Parker's name tag is shown with a red border and white text. Peter in his Spider-Man costume, jumps onto a metal platform.

**AD Prediction** (via MM-Narrator)  
Spider-Man jumps off a yellow taxi and continues running on the street.

**Figure 1.** We present MM-Narrator, a training-free framework towards automatic audio description (AD) generation for long-form videos via iterations: for each scene, it perceives multimodal inputs (i.e., seeing **visual frames** and hearing **character dialogues**), recalls the **context AD** depicting past scenes, and infers **AD prediction** for the current scene. Zoom in for details.

## Abstract

We present MM-Narrator, a novel system leveraging GPT-4 with multimodal in-context learning for the generation of audio descriptions (AD). Unlike previous methods that primarily focused on downstream fine-tuning with short video clips, MM-Narrator excels in generating precise audio descriptions for videos of extensive lengths, even beyond hours, in an autoregressive manner. This capability is made possible by the proposed memory-augmented generation process, which effectively utilizes both the short-

term textual context and long-term visual memory through an efficient register-and-recall mechanism. These contextual memories compile pertinent past information, including storylines and character identities, ensuring an accurate tracking and depicting of story-coherent and character-centric audio descriptions. Maintaining the training-free design of MM-Narrator, we further propose a complexity-based demonstration selection strategy to largely enhance its multi-step reasoning capability via few-shot multimodal in-context learning (MM-ICL). Experimental results on MAD-eval dataset demonstrate that MM-Narrator consistently outperforms both the existing fine-tuning-based approaches and LLM-based approaches in most scenarios,

\* Work done during internship at Microsoft.

as measured by standard evaluation metrics. Additionally, we introduce the first segment-based evaluator for recurrent text generation. Empowered by GPT-4, this evaluator comprehensively reasons and marks AD generation performance in various extendable dimensions.

## 1. Introduction

Audio Description (AD) is an essential task that transforms visual content into spoken narratives [1], primarily assisting visual impairments in accessing video content. Given its evident importance, the notable expectations for AD to fulfill include complementing the existing audio dialogue, enhancing viewer understanding, and avoiding overlap with the original audio. This process involves identifying not just who is present in the scene and what actions are taking place, but also precisely how and when the actions occur. Additionally, AD should capture subtle nuances and visual cues across different scenes, adding layers of complexity to its generation.

In addition to aiding visually impaired audiences, AD also enhances media comprehension for autistic individuals, supports eyes-free activities, facilitates child language development, and mitigates inattentional blindness for sighted users [25, 47]. However, traditional human-annotated AD, while detailed, incurs significant costs and often suffers from inconsistencies due to low inter-annotator agreement [21], highlighting the need for automatic AD generation systems. Furthermore, AD serves as an emerging testbed for benchmarking the capabilities of LLM/LMM systems in long-form multimodal reasoning [21, 22, 30], towards next-level advanced video understanding.

In this paper, we present MM-Narrator, a multimodal AD narrator, to effectively leverage multimodal clues, including visual, textual, and auditory elements, to enable comprehensive perception and reasoning. In particular, MM-Narrator distinguishes itself by naturally identifying characters through their dialogues, in contrast to existing methods that may underutilize subtitles [21, 22].

Apart from an intricate multimodal understanding of the video content, generating story-coherent AD for long-form videos also relies on an accurate tracking and depicting of character-centric evolving storylines over extended durations, even spanning hours. This differs AD generation from conventional dense video captioning [24, 28, 61, 65]: Unlike mere frame-by-frame scene description, AD should weave a coherent narrative, utilizing characters as pivotal elements to maintain an uninterrupted storytelling flow [1]. To achieve contextual understanding, we propose to leverage both short-term and long-term memories to assist MM-Narrator in its recurrent AD generation process. Specifically, short-term textual memory sets the stage for generating coherent narrations, whereas long-term visual memory

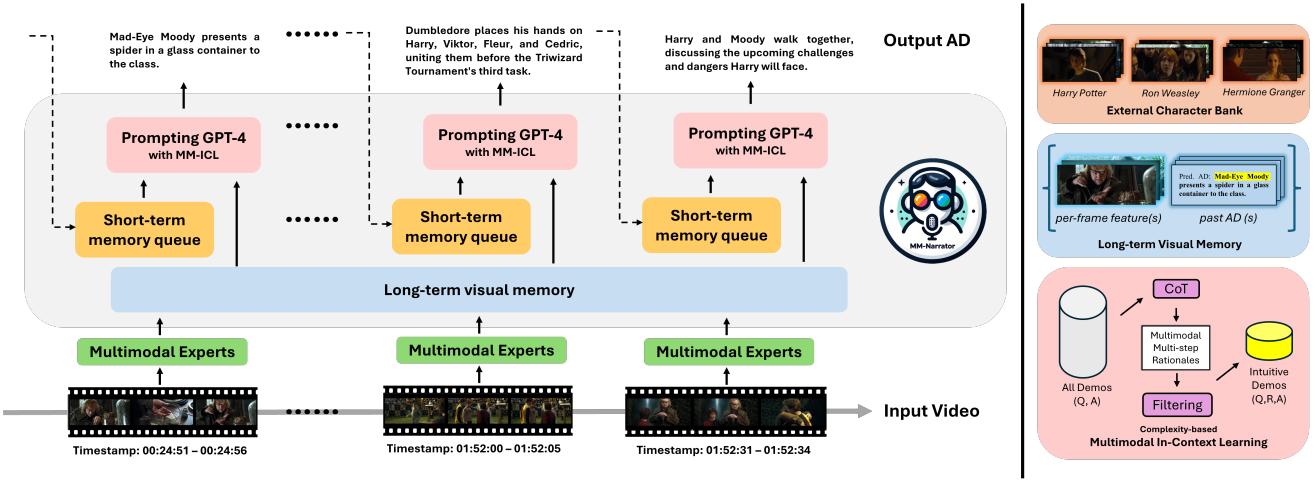
aids in character re-identification during long-form videos, especially for scenes lacking dialogue.

As a GPT-4 empowered multimodal agent, MM-Narrator could further benefit from multimodal in-context learning (MM-ICL) via our proposed complexity-based multimodal demonstration selection. With complexity defined with the chain-of-thought (CoT) technique [63], MM-Narrator could efficiently form and learn from a smaller candidate pool of multimodal demonstrations, effectively improving its multimodal reasoning capability in a few-shot approach. This proposed complexity-based selection surpasses both random sampling and similarity-based retrieval, which are classic ICL solutions in choosing few-shot examples.

In summary, our contributions are four-folds: (1) We present MM-Narrator, an automatic AD narrator for long-form videos that can perceive multimodal inputs, recall past memories, and prompt GPT-4 to produce story-coherent and character-centric AD. (2) We propose a complexity-based multimodal in-context learning (MM-ICL) to further boost its AD generation performance with few-shot examples, offering new insights into the question “*what makes good ICL examples?*” under complex text generation scenarios with multimodal reasoning needed. (3) Our training-free MM-Narrator outperforms both fine-tuning-based SOTAs and LLM/LMM baselines, including GPT-4V, in most classic captioning metrics. (4) Furthermore, we introduce the first GPT-4 based evaluator for recurrent text generation, measuring more comprehensive AD generation qualities at both text-level and sequence-level. Results suggest that MM-Narrator generates AD comparable to human annotations across several considered aspects.

## 2. Related Work

**Audio Description** (AD) offers verbal narration of key visual elements in videos [1], enriching the viewing experience for individuals who are blind or have low vision. AD differs from video captioning [7, 24, 28, 31, 61, 65], which solely describes the visual content of a given video clip. Instead, AD generation considers multiple modalities, aiming to generate coherent narratives of storylines, characters, and actions in a way that complements the regular audio track. Initial studies [50, 51, 53, 57] concentrated on developing audio segmentation and transcription system to collect high-quality video datasets with temporally aligned ADs. These foundational efforts pave the way for more advanced explorations in LSMDC [51]. Recent research [22] has ventured into training transformer models equipped with a frozen LLM. Researchers also incorporate an external character bank [21] to enhance the accuracy of AD generation. Different from prior works [21, 22] that rely on downstream fine-tuning, our proposed MM-Narrator generates accurate ADs in a training-free manner.



**Figure 2.** MM-Narrator generates AD sequence for long-form videos via iterations.

**LLM for Video Understanding.** The remarkable success of Large Language Models (LLMs) [8, 14, 15, 17, 44, 58] has sparked increasing interest in their application to video understanding. Recent works [6, 12, 27, 30, 37, 54] generally fall into two main categories: (i) visual instruction tuning, and (ii) prompting LLMs. The first approach [27, 32, 37, 38] typically fine-tunes an LLM-based model. This involves integrating the pre-trained LLMs and additional trainable networks. The second category [4] involves prompting LLMs to invoke specialized expert tools, transforming the input video into a textual document, which then serves as input to the LLMs for reasoning [6, 12, 30]. However, this strategy may not be effective for processing lengthy or speech-dense videos, as the LLMs often face challenges with excessive token lengths. Different from prior work, we propose to leverage short-term textual memory and long-term visual memory with a register-and-recall mechanism, to effectively generate ADs for long-form videos.

**In-Context Learning (ICL)** [13, 34, 36, 39, 40], as a new paradigm, allows LLMs to learn from a few examples without needing parameter updates via downstream fine-tuning. This learning-from-analogy strategy [16] augments original query question with a context formed by natural language demonstrations. Existing studies highlight that the success of ICL largely depends on the selection of effective demonstrations. One common solution [33, 52] is to form the ICL prompt with closest neighbors, which are retrieved with highest similarity to the query embedding. Other query-based metrics are also explored in finding supportive ICL examples on the basis of query content, such as mutual information [55, 56] and perplexity [20]. Although prior works have demonstrated their superiority in text classification tasks or open-domain QA [16], they have not explored ICL on complex text generation tasks under multimodal scenarios. In this work, we propose to quantify the demonstration complexity as the number of reasoning steps in chain-of-thoughts (CoTs) [63, 68], and select the most

intuitive examples to improve AD generation with few-shot MM-ICL.

### 3. Method

Given a long-form video  $\mathcal{V}$ , consisting of multiple video clips  $\{v_t\}$ , MM-Narrator generates an AD sequence  $\{\mathcal{T}_t\}$  in an autoregressive manner, as shown in Figure 2. We first present MM-Narrator, a multimodal narrator that conducts recurrent AD generation via prompting GPT-4 (§3.1). Building upon MM-Narrator, we propose the complexity-based MM-ICL to further enhance its multimodal reasoning capabilities through intuitive few-shot demonstrations (§3.2). Notably, the entire MM-Narrator framework operates in a training-free manner.

#### 3.1. Recurrent AD Narrator

At each iteration of scanning through a specific long-form video, MM-Narrator utilizes multimodal experts for perception, recalls past memories in both short-term and long-term contexts, and prompts LLM to generate an audio description. We describe each step as below.

**Multimodal perceptions.** We employ specialized vision and audio expert models to extract multimodal information from the input video clip. These off-the-shelf multimodal models are employed as integral tools within our MM-Narrator framework. We denote a video clip consisting of  $N$  frames with timestamp  $t$  as  $v_t = \{\mathcal{I}_1, \mathcal{I}_2, \dots, \mathcal{I}_N\}$ . We deploy vision experts [48, 60, 64] to gather visual perceptions, which involves obtaining per-frame visual features and text-formed outputs. Specifically, for each frame  $\mathcal{I}_i$ , we collect CLIP-ViT features  $x_i^{CLIP}$ , image captions  $x_i^{cap}$ , and people detections  $x_i^{det}$ . Alongside these crucial visuals, we observe the spoken dialogues play a profound role, which is underutilized in existing approaches [21, 22]. The spoken dialogues not only offer information complementary to the visuals, but also primarily serve as the

only access to identify characters with their names when no external video metadata is given. To be specific, we concatenate the subtitles within a certain time window  $T_{sub}$  as  $\mathbf{x}_{t \in T_{sub}}^{sub}$ . These subtitles can be sourced from the Internet or generated through automated speech recognition (ASR) as an audio expert [10]. To summarize, for a given video clip  $v_t$ , the multimodal experts produce a comprehensive tuple of perception clues  $\mathcal{X}_t = (\{\mathbf{x}_i^{CLIP}\}, \{\mathbf{x}_i^{cap}\}, \{\mathbf{x}_i^{det}\}, \mathbf{x}_{t \in T_{sub}}^{sub})$ , where  $\{\mathbf{x}_i\}$  denotes the per-frame outputs. Among these,  $\{\mathbf{x}_i^{cap}\}$  and  $\mathbf{x}_{t \in T_{sub}}^{sub}$  are directly used in constructing LLM prompts, while the others facilitate the proposed register-and-recall mechanism for long-term character re-identification.

**Short-term memory queue.** To equip MM-Narrator with contextual understanding for coherent AD generation, we maintain a short-term memory queue  $\mathcal{M}_{short} = \{\mathcal{T}_{t-K}, \dots, \mathcal{T}_{t-1}\}$  to contain the  $K$  most recently predicted ADs with timestamps. The short-term memory queue will be updated over time during inference. This lightweight textual queue is instrumental in creating story-coherent AD narrations, enabling visually impaired audiences to follow the storytelling more intuitively.

**Long-term visual memory.** To endow MM-Narrator with the ability to recall characters identified in previous video clips, we construct a frame-level character re-identification visual bank. This visual bank, designed for long-term use, is operated by a register-and-recall mechanism as follows: (1) we register  $\mathbf{x}_j^{CLIP}$  as the visual signature for each globally-indexed frame  $\mathcal{I}_j$  in all previous video clips  $\mathcal{I}_j \in \{v_1, v_2, \dots, v_{t-1}\}$ , and (2) for each current frame  $\mathcal{I}_i$ , we first filter-out the invalid matches resulting in nonpositive cosine similarity  $\text{Sim}_{cos}(\mathbf{x}_i^{CLIP}, \mathbf{x}_j^{CLIP})$ , and then retrieve the past predicted AD which owns the highest similarity to the current visual signature  $\mathbf{x}_i^{CLIP}$ . For simplicity, this mechanism is activated only when a single individual is detected in a frame (i.e.,  $|\mathbf{x}_i^{det}| = 1$ ), typically in close-up shots of the character, making frame-level CLIP-ViT features [48] compatible for character re-identification. Given any AD that covers multiple frames, this frame-level visual retriever supports the MM-Narrator in re-identifying multiple characters appearing in the video clip. Additionally, the retrieval candidate pool is refined to include only past predicted ADs where person named entities are recognized through a Named Entity Recognition (NER) tool [18]. This strategy focuses MM-Narrator on the main characters who contribute to the past storyline.

**Prompting LLM for AD generation.** Gathering all aforementioned text-formed outputs, MM-Narrator builds prompts to query GPT-4 for recurrent AD generation. Specifically, the input prompt contains the following elements: task introduction, visual captions ( $\mathbf{x}_i^{cap}$ ) with successfully re-identified characters, recent context ADs

( $\mathcal{M}_{short}$ ) and character dialogues ( $\mathbf{x}_{t \in T_{sub}}^{sub}$ ). Noticeably, we also found that adding task-specific hints into the prompt could empirically benefit overall AD generation, which we attribute as an explicit attention guidance via prompt engineering. A breakdown of our AD generation prompt constructed by MM-Narrator, is provided in the supplementary (Figure 7).

### 3.2. Multimodal In-Context Learning

In this section, we further extend MM-Narrator with multimodal in-context learning (MM-ICL) on few-shot examples. Our exploration begins by examining two primary methods of demonstration selection: random and similarity-based approaches. We then critically evaluate the question, “*What makes for effective ICL examples?*” and propose a complexity-based MM-ICL approach to improve the multimodal reasoning capability with the most intuitive multimodal demonstrations.

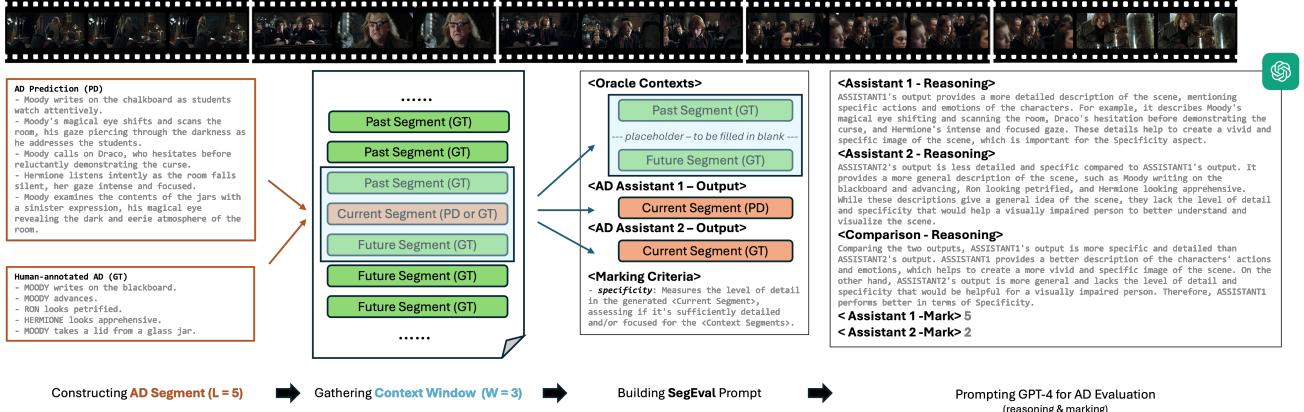
**Random MM-ICL.** Firstly, we build an in-context learning (ICL) demonstration pool, denoted as  $\mathcal{P}$ , from the training dataset. Each demonstration within the pool is composed of a pair  $(\mathcal{Q}, \mathcal{A})$ , where  $\mathcal{Q}$  represents the text-formed *question* created using multimodal experts, and  $\mathcal{A}$  is the corresponding ground-truth AD, serving as the *answer*. Then, for each test query  $q$ , we randomly sample  $C$  demonstrations from  $\mathcal{P}$  to facilitate the ICL process.

Furthermore, we are further interested in two essential questions: “*What makes good examples for AD generation?*” and “*How to find and use them for ICL?*”

**Similarity-based MM-ICL.** A common approach, as suggested in existing literature [33], is to identify “good examples” based on similarity, employing a  $k$ -NN algorithm to select examples that exhibit the highest  $k$  similarity between the embeddings of  $\mathcal{Q}_i$  and the test query  $q$ . This solution expects to find supportive examples to benefit few-shot performance via a “soft-copy” mechanism [23, 43], which is often used in text classification tasks, such as sentiment analysis, or relatively-simple text generation task such as open-domain QA [16].

However, we empirically find that this similarity-based approach does not manage to enhance the ICL capability for AD generation, regardless of whether the retrieved examples are presented in descending order [33] or ascending order [16]. We hypothesize that for complex text generation tasks such as AD generation, which requires multimodal perception and reasoning, similarity or relevance may not be the most suitable criteria for identifying effective ICL examples for improving overall performance.

**Complexity-based MM-ICL.** Our empirical analysis reveals that not all questions are equally challenging, in terms of the complexity of multimodal fusion. Take Figure 1 as an example: when comparing *Titanic* (1997) to *Spider Man* (2018), the latter presents a more complex case. It requires



**Figure 3.** Our proposed SegEval evaluator to measure recurrent text generation quality with GPT-4 under customized marking criteria. Noticeably, GPT-4 is agnostic to the source of each assistant output (i.e., which Seg is GT or PD), and it would measure Seg quality taking oracle contexts into consideration. Take the response shown above as example, its corresponding re-scaled  $r$  is 2.25. Zoom in for details.

the inference that “Peter and Spider Man are the same character”, a deduction drawn from context AD and subtitles, alongside describing his actions from visual frames, enriched by contextual understanding from the context AD.

This observation led us to hypothesize that complexity could be a more suitable metric for identifying effective ICL examples for tasks involving intricate multimodal fusion. To this end, we propose to query LLM to articulate the chain-of-thoughts (CoTs) as reasoning steps, denoted as  $\mathcal{R}$ , that assist in deriving the answer  $\mathcal{A}$  from the question  $\mathcal{Q}$ . This process evolves our demonstration format from simple  $(\mathcal{Q}, \mathcal{A})$  pairs to more comprehensive  $(\mathcal{Q}, \mathcal{R}, \mathcal{A})$  tuples.

Instead of the conventional random sampling from the entire pool  $\mathcal{P}$ , we propose selecting the most straightforward examples, quantified by the shortest number of reasoning steps. These are compiled into a simpler subset pool  $\mathcal{P}_{simple}$ , from which we conduct our demonstration sampling. This method ensures the inclusion of more intuitive and concise examples in our MM-ICL process. We present detailed ablation study in §5.4, validating that complexity serves as a robust measure for selecting effective ICL examples for improving AD generation.

## 4. Segment-based GPT-4 Evaluator

The lack of standard AD annotation guidelines, varying cultural background and preferences of human annotators imply that AD is an inherently subjective recurrent text generation process, leading to notable inter-annotator disagreements [21] and challenges in evaluation using traditional reference-based captioning metrics. To this end, inspired by [32, 35], we propose a segment-based GPT-4 evaluator SegEval to measure the recurrent AD generation, in terms of multi-domain qualities.

Suppose  $L$  ADs form one segment Seg. For each Seg, the evaluator takes into consideration an oracle context window Ctx of length  $W$ , to measure its multi-aspect scores. Specifically, we gather  $W - 1$  adjacent segments to form

Ctx, which consists of  $\frac{W-1}{2}$  past and  $\frac{W-1}{2}$  future segments surrounding the targeted Seg. Given a pair of predicted (PD) and ground-truth (GT) AD segments, SegEval would treat them as outputs of two separate AD generation systems, and query GPT-4 to reason and mark their raw marks independently. The final score is calculated as the ratio  $r$  of these raw marks between predicted and human-annotated AD, via post-processing. If the re-scaled  $r$  is higher than 1.0, it indicates that GPT-4 might favour the predicted AD over human annotations under the specific aspect. Besides, this rescaling operation makes it comparable among different approaches, sharing human annotations as the marking standard. Noticeably, although GPT-4 is unaware of the segment source that which Seg is the GT or PD, we always form Ctx from GT annotations to set the oracle for investigating contextual influences.

Overall, as shown in Figure 3, SegEval can measure context-irrelevant, short-context and long-context scores by flexibly changing the value of  $W$ . For example, it could measure *text-level qualities* such as originality and consistency (when  $W = 1$ ), while it could also mark *sequence-level qualities* such as coherence, diversity and specificity (when  $W > 1$ ). The details of each marking criteria are provided in supplementary (§D).

## 5. Experiments

### 5.1. Evaluation Setup

**Datasets.** We conduct experiments on the AD generation benchmark established in AutoAD [22], where MAD-v2-Named and MAD-eval-Named are released as training and testing splits, respectively. **MAD-v2-Named** consists of 334,296 ADs and 628,613 subtitles from 488 movies, while **MAD-eval-Named** is compromised of 6,520 ADs and 10,602 subtitles from 10 movies.

**Metrics.** Following AutoAD [22], we report three traditional captioning metrics to measure the quality of ADs gen-

erated versus human-annotated ones, including ROUGE-L [29] (**R-L**), CIDEr [59] (**C**) and SPICE [9] (**S**). Besides, we follow AutoAD-II to benchmark the text sequence generation over their recall-based metric ‘Recall@k within Neighbours’ (**R@k/N**), where the text similarity is measured by BertScore [67]. We also report Bleu-1 [46] and METEOR [11] for ablation studies. To reduce experimental variability, each experiment of MM-Narrator is repeated *three* times in Tables 1 to 6, as well as Figure 5, with mean (and std) reported.

## 5.2. Comparison with State-of-the-Art Approaches

**Fine-tuning-based SOTAs.** We first compare our training-free framework against the fine-tuning-based SOTAs, including ClipCap [41], ClipDec [42] and AutoAD-I [22]. As shown in Table 1, our training-free approach outperforms its fine-tuning-based counterparts [22, 41, 42], in terms of ROUGE-L, SPICE and R@k/N, especially the AutoAD-I [22] (R-L 12.1 vs 11.9; S 4.5 vs 4.4; R@k/N 48.0 vs 42.1) which is proposed to conduct partial data pretraining over an extra large-scale text-only AV-AD dataset [2, 22] (consisting of 3.3M ADs from over 7k movies) to address the lack of paired training data for AD generation. Unlike [21, 22] who report to struggle with benefiting from character dialogues, our MM-Narrator could better integrate multimodal information and effectively identify characters from appropriate subtitle usage (shown as model D in §5.3).

**Training-free LLM/LMM Baselines.** We next compare our MM-Narrator with LLM and LMM baselines: (a) VLog [6] and (b) VideoChat-Text [27] are two LLM-based methods for multimodal video understanding. They convert multimodal perceptions into natural languages via several pretrained models [26, 49, 62, 64], and then utilizes a LLM to generate texts based on task-specific prompts. To make a fair comparison, we make them query GPT-4 with the same AD generation prompt as we use in MM-Narrator. (c) MM-Vid [30] is a LMM system which generates AD through incorporating external knowledge with clip-level video description generated by GPT-4V [45, 66].

As shown in Table 2, our MM-Narrator (*w/o MM-ICL*) would outperform VLog and VideoChat, which is mainly attributed to the proposed short-term memory queue and long-term visual memory to effectively leverage relevant contextual information recalled from past ADs. In addition, while MM-Narrator is based on GPT-4 (text-only), it also surpasses the GPT-4V(ision) based MM-Vid system in terms of R-L and SPICE. The results suggest that a memory-augmented LLM can be comparably valuable to the perception-enhanced ones. Furthermore, with our proposed MM-ICL, MM-Narrator outperforms these training-free LLM/LMM counterparts by a large margin. Finally, the bottom two rows of Table 2 further validate the effectiveness of the proposed MM-ICL design.

Method	Training-Free	R-L (↑)	C (↑)	S (↑)	R@5/16 (↑)
ClipCap [41]	✗	8.5	4.4	1.1	36.5
ClipDec [42]	✗	8.2	6.7	1.4	-
AutoAD-I [22]	✗	11.9	<b>14.3</b>	4.4	42.1
<b>MM-Narrator</b>	✓	<b>12.1</b>	11.6	<b>4.5</b>	<b>48.0</b>

**Table 1.** Comparisons with fine-tuning-based state-of-the-art methods on MAD-eval-Named benchmark. Note: the random guess will result in a R@5/16 of 31.3%.

Method	LLM/LMM	R-L (↑)	C (↑)	S (↑)	R@5/16 (↑)
VLog [6]	GPT-4	7.5	1.3	2.1	42.3
VideoChat [27]	GPT-4	7.9	2.4	1.8	42.5
MM-Vid [30]	GPT-4V	9.8	6.1	3.8	46.1
<b>MM-Narrator</b>					
<i>w/o MM-ICL</i>	GPT-4	10.3	4.9	3.8	47.1
<i>w/ MM-ICL</i>	GPT-4	<b>12.1</b>	<b>11.6</b>	<b>4.5</b>	<b>48.0</b>

**Table 2.** Comparisons with training-free LLM/LMM baselines on MAD-eval-Named benchmark.

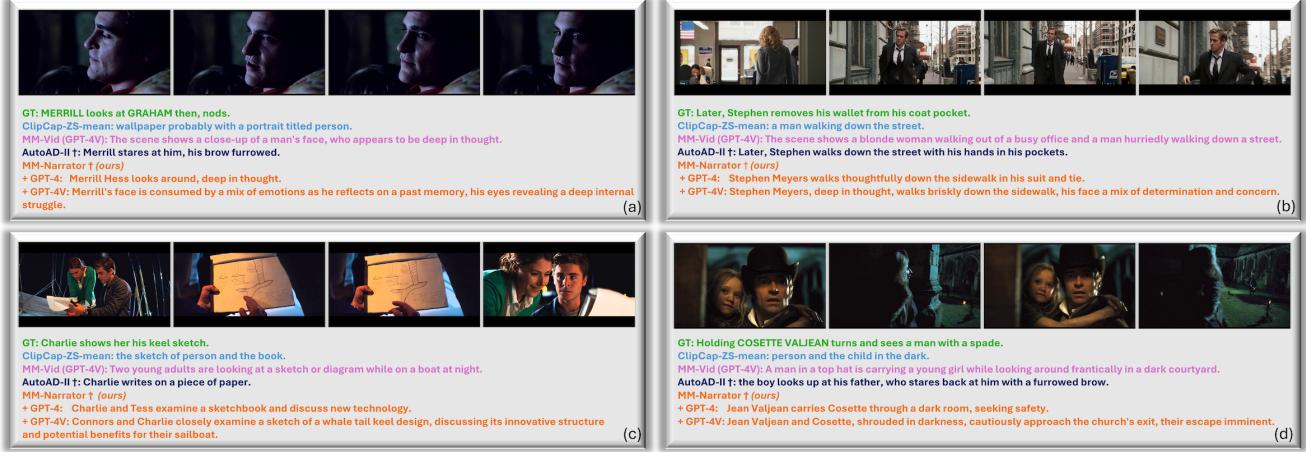
Method	Training-Free	R-L (↑)	C (↑)	S (↑)	R@5/16 (↑)
AutoAD-II † [21]	✗	13.4	19.5	-	50.8
<b>MM-Narrator</b> †	✓	13.4	13.9	5.2	49.0

**Table 3.** Evaluation on MAD-eval-Named benchmark, with an external character bank annotated and utilized for improved character recognition (denoted as †).

**Utilizing External Character Bank.** Previously, all discussed methods share the same and only knowledge source to assist in character recognition. More specifically, they, like us humans, mostly identify characters and infer their names through hearing (i.e., auditory cues) *alone* when watching movies. Given this single source of gaining character information, our MM-Narrator would convey contextual information via retrieving visual and temporal memories. However, these methods suffer from an unavoidable limitation: The character identities would unfortunately remain mystery until their names are being first-time called in dialogues.

To alleviate that, following AutoAD-II [21] we also investigate how our method could benefit from incorporating an external character bank. To construct this character bank, [21] exploits actor portrait images (from an external movie database) to retrieve a few most similar frames for each main character in each movie. Unlike [21] who trains an auxiliary character recognizer from these retrieved frames, we maintain our training-free designs by simply concatenating these frames into short video clips to introduce each character (with ADs as their names). Next, we prepend these video clips to the long-form videos, such that they could work compatibly with our register-and-recall mechanism. As shown in Table 3, our MM-Narrator (*w/ ExtChar-Bank*) could further boost its performance and generate outcomes comparable to the fine-tuning-based AutoAD-II.

**Qualitative Results.** Qualitative comparisons over MAD-



**Figure 4.** Qualitative comparisons between ClipCap, MM-Vid, AutoAD-II, and our MM-Narrator, where the latter two approaches are equipped with the external character bank. The movies are from (a) *Signs* (2002), (b) *Ides of March* (2011), (c) *Charlie St. Cloud* (2010), and (d) *Les Misérables* (2012). Zoom in for details.

eval dataset are shown as Figure 4, while the qualitative demonstrations of applying our MM-Narrator on other long-form videos (external to the MAD-eval dataset) are shown in Figure 1. Additional qualitative results are included in supplementary (Figure 10 and 11).

### 5.3. Building MM-Narrator From Image Captioner

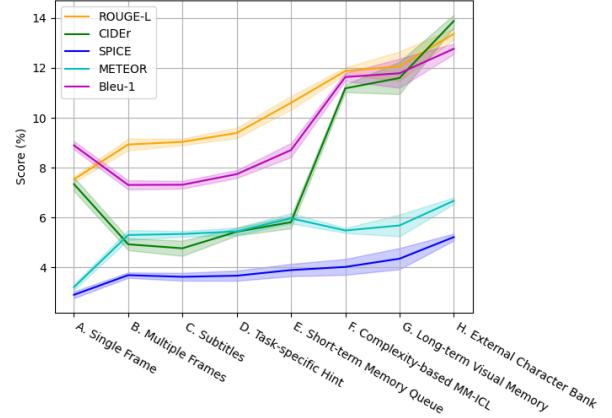
As shown in Figure 5, we quantitatively demonstrate how our training-free MM-Narrator are developed step by step. Starting from (A) an image captioner, we elaborate how multimodal perception benefits MM-Narrator to form an intricate multimodal understanding over video content. Specifically, it includes adding (B) multiple frames, (C) subtitles, and (D) a task-specific hint<sup>1</sup>. Noticeably, simply adding the dialogues (C) might not result in an immediate performance gain. However, with prompt engineering in (D), MM-Narrator pays more attention to effectively leverage multimodal clues for character-centric AD generation.

Next, we illustrate how we transform MM-Narrator into recurrent AD narrator to produce story-coherent AD, with incorporation of past memories and complexity-based MM-ICL. Specifically, MM-Narrator maintains (E) a short-term memory queue, learns from (F) multimodal demonstrations via MM-ICL, and retrieves (G) long-term visual memory for character re-identification, which could be further boosted with (H) an external character bank.

### 5.4. Ablations on Multimodal In-Context Learning

We investigated three groups of MM-ICL proposed to augment the baseline. Specifically, we built random R1 and similarity-based S1, by adapting classic ICL techniques [13, 33] from conventional NLP tasks into multi-

<sup>1</sup>“Hint: try to infer character names from subtitles for AD generation.”



**Figure 5.** Ablations on each component for MM-Narrator.

Model Pool Size	Demo. Format	CoT	R-L (↑)	C (↑)	B-1 (↑)
<b>Baseline w/o MM-ICL</b>					
B1	-	-	✗	11.8 <sub>±0.1</sub>	8.6 <sub>±0.1</sub>
<b>Random MM-ICL</b>					
R1	100% ( $\mathcal{Q}, \mathcal{A}$ )	✗	13.2 <sub>±0.1</sub>	12.9 <sub>±0.2</sub>	12.2 <sub>±0.1</sub>
R2	100% ( $\mathcal{Q}, \mathcal{R}, \mathcal{A}$ )	✓	13.4 <sub>±0.1</sub>	13.4 <sub>±0.2</sub>	12.7 <sub>±0.1</sub>
R3	10% random ( $\mathcal{Q}, \mathcal{A}$ )	✗	13.3 <sub>±0.1</sub>	13.0 <sub>±0.1</sub>	12.3 <sub>±0.0</sub>
R4	10% random ( $\mathcal{Q}, \mathcal{R}, \mathcal{A}$ )	✓	13.3 <sub>±0.1</sub>	13.4 <sub>±0.1</sub>	12.6 <sub>±0.0</sub>
<b>Similarity-based MM-ICL</b>					
S1	100% ( $\mathcal{Q}, \mathcal{A}$ )	✗	13.5 <sub>±0.0</sub>	13.1 <sub>±0.0</sub>	12.6 <sub>±0.1</sub>
<b>Complexity-based MM-ICL</b>					
C1	10% shortest ( $\mathcal{Q}, \mathcal{A}$ )	✗	13.2 <sub>±0.1</sub>	13.3 <sub>±0.3</sub>	12.3 <sub>±0.1</sub>
C2	10% shortest ( $\mathcal{Q}, \mathcal{R}, \mathcal{A}$ )	✓	13.4 <sub>±0.0</sub>	13.9 <sub>±0.1</sub>	12.8 <sub>±0.0</sub>
C3	10% longest ( $\mathcal{Q}, \mathcal{R}, \mathcal{A}$ )	✓	13.3 <sub>±0.1</sub>	12.7 <sub>±0.2</sub>	12.4 <sub>±0.1</sub>

**Table 4.** Our different MM-ICL designs for MM-Narrator †. The baseline (B1) and each representative MM-ICL implementation (R1, S1 and C2) are highlighted.

modal AD generation. Next, we presented our complexity-based design as C2.

Method	LLM/LMM	Text-level Quality		Sequence-level Quality					
		Context-irrelevant Scores		Short-context Scores			Long-context Scores		
		Orig. ±0.02	Cons. ±0.02	Cohes. ±0.01	Dive. ±0.06	Spec. ±0.04	Cohes. ±0.01	Dive. ±0.01	Spec. ±0.03
GT	-	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00
ClipCap [41]	GPT-2	0.43	0.42	0.26	0.35	0.35	0.26	0.42	0.33
VLog [6]	GPT-4	1.03	0.88	0.34	0.55	0.52	0.32	0.57	0.43
MM-Vid [30]	GPT-4V	0.85	0.78	0.51	0.81	0.66	0.53	0.84	0.62
<b>MM-Narrator</b>	GPT-4	$1.05 \pm 0.10$	$1.03 \pm 0.05$	$0.52 \pm 0.06$	$0.70 \pm 0.06$	$0.66 \pm 0.04$	$0.57 \pm 0.05$	$0.70 \pm 0.02$	$0.61 \pm 0.05$
<b>MM-Narrator</b>	GPT-4V	$1.49 \pm 0.10$	$1.45 \pm 0.05$	$0.94 \pm 0.07$	$1.01 \pm 0.04$	$1.13 \pm 0.08$	$0.87 \pm 0.04$	$1.05 \pm 0.04$	$1.14 \pm 0.05$
<b>MM-Narrator</b> †	GPT-4	$0.95 \pm 0.02$	$1.06 \pm 0.01$	$0.62 \pm 0.04$	$0.75 \pm 0.01$	$0.76 \pm 0.01$	$0.62 \pm 0.04$	$0.80 \pm 0.03$	$0.71 \pm 0.03$
<b>MM-Narrator</b> †	GPT-4V	$1.45 \pm 0.14$	$1.46 \pm 0.04$	$0.98 \pm 0.03$	$1.06 \pm 0.04$	$1.24 \pm 0.09$	$0.94 \pm 0.02$	$1.09 \pm 0.05$	$1.12 \pm 0.03$

**Table 5.** Evaluating AD generation with SegEval on MAD-eval-Named benchmark, with segment size  $L$  set to 5. The context window sizes  $W$  are set as 1 / 3 / 11 to compute context-irrelevant / short-context / long-context scores, respectively. Orig., Cons., Cohe., Dive., and Spec. stand for *originality*, *consistency*, *coherence*, *diversity*, and *specificity*, respectively. The scoring variances of these GPT-4 evaluators are denoted below for references, which are estimated by three repeated evaluations over the same inference outputs. These re-scaled scores measure the corresponding AD prediction (PD) qualities of each specific method, compared to the shared marking standards set by ground-truth (GT) ADs. For example, given a pair of PD and GT segments, without revealing to the evaluator which segment is GT or PD, if it reasons and marks the raw qualities (R.Q.) as 8 and 5 for PD and GT segments, respectively, we derive the re-scaled score  $r$  as  $\frac{R.Q._{PD}}{R.Q._{GT}} = \frac{8}{5} = 1.6$ . † indicates our incorporation with ExtCharBank.

Method	R-L (↑)	C (↑)	M (↑)	B-1 (↑)
<b>MM-Narrator</b>				
+ GPT-4	$12.1 \pm 0.4$	$11.6 \pm 0.4$	$5.7 \pm 0.2$	$11.8 \pm 0.3$
+ GPT-4V	$11.8 \pm 0.1$	$7.0 \pm 0.2$	$6.5 \pm 0.1$	$9.3 \pm 0.1$
<b>MM-Narrator</b> †				
+ GPT-4	$13.4 \pm 0.0$	$13.9 \pm 0.1$	$6.7 \pm 0.0$	$12.8 \pm 0.0$
+ GPT-4V	$12.8 \pm 0.0$	$9.8 \pm 0.2$	$7.1 \pm 0.0$	$10.9 \pm 0.0$

**Table 6.** Comparisons over classic reference-based captioning scores, when incorporating our MM-Narrator with GPT-4V.

The results as shown in Table 4, verify our hypothesis that complexity serves as an appropriate measure for selecting effective ICL demonstrations for improving AD generation. It also indicates that our proposed complexity-based design (C2) is more preferable than classic ones (R1, S1) for AD generation, especially the CIDEr score. In supplementary (§C), we further discuss three sub-questions to elaborate an in-depth analysis, including 1) *Does CoT help?* 2) *Are more intuitive examples helpful for AD Generation?* and 3) *Does complexity-based MM-ICL work effectively?*

### 5.5. Evaluating AD Generation with GPT-4

In Table 6, we observe a few performance drop on classic reference-based captioning scores when incorporating MM-Narrator with GPT-4V [45]. As shown in Figure 4, the decrease in performance can be primarily attributed to the more detailed and much richer ADs generated by our method, which diverge from the typically shorter human-annotated ADs in MAD-eval-Named. This suggests that taking human annotated AD as oracles to measure AD-level captioning scores might be unsuitable for advanced LMM approaches, which further motivates our proposal of evaluating recurrent text generation with GPT-4.

Adjusting  $W$ , our proposed SegEval could flexibly measure both *text-level* and *sequence-level* qualities. As

shown in Table 5, the performance ranking order observed in SegEval aligns with our other experimental results, validating the reliability of SegEval as an evaluation tool, except for GPT-4V based MM-Vid where ours falls short on *diversity*. Furthermore, when employing GPT-4V as our vision expert, MM-Narrator not only outperforms others by a large margin, but also closely mirrors the quality of human annotated ADs in multiple aspects, gaining more favor from the source-agnostic GPT-4 evaluator.

Compared to classic reference-based captioning scores, SegEval could better reflect the recurrent text generation qualities with GPT-4. One human validation on SegEval is shown in Figure 3, and more examples can be found in supplementary (Figure 9). Moreover, SegEval could be easily extended to support more comprehensive evaluation perspectives by querying it with extra customized marking criteria.

## 6. Conclusion

MM-Narrator represents a significant leap in automatic audio description (AD) generation for long-form videos, leveraging the power of GPT-4 and innovative multimodal in-context learning (MM-ICL). This recurrent AD narrator excels in generating story-coherent and character-centric AD by combining immediate textual context with long-term visual memory. Its training-free design, coupled with our proposed complexity-based MM-ICL demonstration selection strategy, outperforms both existing fine-tuning-based and LLM-based approaches in most scenarios, as measured by traditional captioning metrics. Furthermore, we introduce a GPT-4 empowered evaluator for a more comprehensive measurement of recurrent text generation qualities. Its results suggest that MM-Narrator generates AD comparable to human annotations across several considered aspects.

## References

- [1] American council of the blind. <https://adp.acb.org/>. 2
- [2] Audiovault. <https://audiovault.net/>. 6
- [3] Google text-to-speech: a python library and cli tool to interface with google translate’s text-to-speech api. <https://gtts.readthedocs.io/en/latest/>. 13
- [4] Langchain. <https://langchain.readthedocs.io/>. 3
- [5] Pyscenedetect: Video scene cut detection and analysis tool. <https://www.scenedetect.com/>. 13
- [6] VLog: Video as a Long document. <https://github.com/showlab/VLog>, 2023. GitHub repository. 3, 6, 8
- [7] Nayyer Aafaq, Ajmal Mian, Wei Liu, Syed Zulqarnain Gilani, and Mubarak Shah. Video description: A survey of methods, datasets, and evaluation metrics. *ACM Computing Surveys*, 52(6):1–37, 2019. 2
- [8] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems*, 35:23716–23736, 2022. 3
- [9] Peter Anderson, Basura Fernando, Mark Johnson, and Stephen Gould. SPICE: Semantic propositional image caption evaluation. In *ECCV*, 2016. 6
- [10] Max Bain, Jaesung Huh, Tengda Han, and Andrew Zisserman. WhisperX: Time-accurate speech transcription of long-form audio. *INTERSPEECH 2023*, 2023. 4, 13
- [11] Satanjeev Banerjee and Alon Lavie. METEOR: An automatic metric for mt evaluation with improved correlation with human judgments. In *ACL workshop*, 2005. 6
- [12] Aanisha Bhattacharya, Yaman K Singla, Balaji Krishnamurthy, Rajiv Ratn Shah, and Changyou Chen. A video is worth 4096 tokens: Verbalize story videos to understand them in zero shot. *EMNLP*, 2023. 3
- [13] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *NeurIPS*, 33:1877–1901, 2020. 3, 7, 13
- [14] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhang-hao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yong-hao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source chatbot impressing GPT-4 with 90% ChatGPT quality, March 2023. 3
- [15] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. PaLM: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022. 3
- [16] Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, and Zhifang Sui. A survey for in-context learning. *arXiv preprint arXiv:2301.00234*, 2022. 3, 4
- [17] Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, et al. PaLM-E: An embodied multimodal language model. *arXiv preprint arXiv:2303.03378*, 2023. 3
- [18] Jenny Rose Finkel, Trond Grenager, and Christopher D Manning. Incorporating non-local information into information extraction systems by gibbs sampling. In *ACL*, pages 363–370, 2005. 4
- [19] Yao Fu, Hao Peng, Ashish Sabharwal, Peter Clark, and Tushar Khot. Complexity-based prompting for multi-step reasoning. *arXiv preprint arXiv:2210.00720*, 2022. 12
- [20] Hila Gonen, Srini Iyer, Terra Blevins, Noah A Smith, and Luke Zettlemoyer. Demystifying prompts in language models via perplexity estimation. *arXiv preprint arXiv:2212.04037*, 2022. 3
- [21] Tengda Han, Max Bain, Arsha Nagrani, Gü̈l Varol, Weidi Xie, and Andrew Zisserman. AutoAD II: The Sequel - who, when, and what in movie audio description. In *ICCV*, 2023. 2, 3, 5, 6
- [22] Tengda Han, Max Bain, Arsha Nagrani, Gü̈l Varol, Weidi Xie, and Andrew Zisserman. AutoAD: Movie description in context. In *CVPR*, 2023. 2, 3, 5, 6, 13
- [23] Xiaochuang Han, Daniel Simig, Todor Mihaylov, Yulia Tsvetkov, Asli Celikyilmaz, and Tianlu Wang. Understanding in-context learning via supportive pretraining data. *arXiv preprint arXiv:2306.15091*, 2023. 4
- [24] Ranjay Krishna, Kenji Hata, Frederic Ren, Li Fei-Fei, and Juan Carlos Niebles. Dense-captioning events in videos. In *ICCV*, pages 706–715, 2017. 2
- [25] Elisa Lewis. Deep dive: How audio description benefits everyone, 2021. Accessed on: 2023-11-13. 2
- [26] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. BLIP-2: bootstrapping language-image pre-training with frozen image encoders and large language models. In *ICML*, 2023. 6
- [27] Kunchang Li, Yinan He, Yi Wang, Yizhuo Li, Wenhui Wang, Ping Luo, Yali Wang, Limin Wang, and Yu Qiao. Videochat: Chat-Centric video understanding. *arXiv preprint arXiv:2305.06355*, 2023. 3, 6
- [28] Linjie Li, Jie Lei, Zhe Gan, Licheng Yu, Yen-Chun Chen, Rohit Pillai, Yu Cheng, Luowei Zhou, Xin Eric Wang, William Yang Wang, et al. VALUE: A multi-task benchmark for video-and-language understanding evaluation. In *NeurIPS Datasets and Benchmarks Track*, 2021. 2
- [29] Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. *ACL*, 2004. 6
- [30] Kevin Lin, Faisal Ahmed, Linjie Li, Chung-Ching Lin, Ehsan Azarnasab, Zhengyuan Yang, Jianfeng Wang, Lin Liang, Zicheng Liu, Yumao Lu, et al. MM-Vid: Advancing video understanding with GPT-4V (ision). *arXiv preprint arXiv:2310.19773*, 2023. 2, 3, 6, 8
- [31] Kevin Lin, Linjie Li, Chung-Ching Lin, Faisal Ahmed, Zhe Gan, Zicheng Liu, Yumao Lu, and Lijuan Wang. SwinBert: End-to-end transformers with sparse attention for video captioning. In *CVPR*, 2022. 2
- [32] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *NeurIPS*, 2023. 3, 5

- [33] Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. What makes good in-context examples for GPT-3? *arXiv preprint arXiv:2101.06804*, 2021. 3, 4, 7, 13
- [34] Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroyaki Hayashi, and Graham Neubig. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Computing Surveys*, 55(9):1–35, 2023. 3
- [35] Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. GPTEval: NLG evaluation using GPT-4 with better human alignment. *arXiv preprint arXiv:2303.16634*, 2023. 5
- [36] Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for science question answering. *NeurIPS*, 35:2507–2521, 2022. 3
- [37] Ruipu Luo, Ziwang Zhao, Min Yang, Junwei Dong, Minghui Qiu, Pengcheng Lu, Tao Wang, and Zhongyu Wei. Valley: Video assistant with large language model enhanced ability. *arXiv preprint arXiv:2306.07207*, 2023. 3
- [38] Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. Video-ChatGPT: Towards detailed video understanding via large vision and language models. *arXiv preprint arXiv:2306.05424*, 2023. 3
- [39] Sewon Min, Mike Lewis, Luke Zettlemoyer, and Hannaneh Hajishirzi. MetaICL: Learning to learn in context. *arXiv preprint arXiv:2110.15943*, 2021. 3
- [40] Sewon Min, Xinxi Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. Rethinking the role of demonstrations: What makes in-context learning work? *arXiv preprint arXiv:2202.12837*, 2022. 3
- [41] Ron Mokady, Amir Hertz, and Amit H Bermano. Clip-Cap: Clip prefix for image captioning. *arXiv preprint arXiv:2111.09734*, 2021. 6, 8, 13
- [42] David Nukrai, Ron Mokady, and Amir Globerson. Text-only training for image captioning using noise-injected clip. *EMNLP Findings*, 2022. 6
- [43] Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, et al. In-context learning and induction heads. *arXiv preprint arXiv:2209.11895*, 2022. 4
- [44] OpenAI. GPT-4 technical report. 2023. 3, 14
- [45] OpenAI. GPT-4V(ision) system card. 2023. 6, 8, 14
- [46] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *ACL*, 2002. 6
- [47] Elisa Perego. Gains and losses of watching audio described films for sighted viewers. *Target*, 28(3):424–444, 2016. 2
- [48] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, pages 8748–8763. PMLR, 2021. 3, 4, 13
- [49] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision. In *ICML*, 2023. 6
- [50] Anna Rohrbach, Marcus Rohrbach, Niket Tandon, and Bernt Schiele. A dataset for movie description. In *CVPR*, pages 3202–3212, 2015. 2
- [51] Anna Rohrbach, Atousa Torabi, Marcus Rohrbach, Niket Tandon, Christopher Pal, Hugo Larochelle, Aaron Courville, and Bernt Schiele. Movie description. *Int. J. Comput. Vis.*, 123:94–120, 2017. 2
- [52] Ohad Rubin, Jonathan Herzig, and Jonathan Berant. Learning to retrieve prompts for in-context learning. *arXiv preprint arXiv:2112.08633*, 2021. 3
- [53] Mattia Soldan, Alejandro Pardo, Juan León Alcázar, Fabian Caba, Chen Zhao, Silvio Giancola, and Bernard Ghanem. MAD: A scalable dataset for language grounding in videos from movie audio descriptions. In *CVPR*, pages 5026–5035, 2022. 2
- [54] Enxin Song, Wenhao Chai, Guanhong Wang, Yucheng Zhang, Haoyang Zhou, Feiyang Wu, Xun Guo, Tian Ye, Yan Lu, Jenq-Neng Hwang, et al. MovieChat: From dense token to sparse memory for long video understanding. *arXiv preprint arXiv:2307.16449*, 2023. 3
- [55] Taylor Sorensen, Joshua Robinson, Christopher Michael Rytting, Alexander Glenn Shaw, Kyle Jeffrey Rogers, Alexia Pauline Delorey, Mahmoud Khalil, Nancy Fulda, and David Wingate. An information-theoretic approach to prompt engineering without ground truth labels. *arXiv preprint arXiv:2203.11364*, 2022. 3
- [56] Eshaan Tanwar, Manish Borthakur, Subhabrata Dutta, and Tanmoy Chakraborty. Multilingual llms are better cross-lingual in-context learners with alignment. *arXiv preprint arXiv:2305.05940*, 2023. 3
- [57] Atousa Torabi, Christopher Pal, Hugo Larochelle, and Aaron Courville. Using descriptive video services to create a large data source for video annotation research. *arXiv preprint arXiv:1503.01070*, 2015. 2
- [58] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. LLaMA 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023. 3
- [59] Ramakrishna Vedantam, C. Lawrence Zitnick, and Devi Parikh. CIDEr: Consensus-based image description evaluation. In *CVPR*, 2015. 6
- [60] Chien-Yao Wang, Alexey Bochkovskiy, and Hong-Yuan Mark Liao. YOLOv7: Trainable bag-of-freebies sets new state-of-the-art for real-time object detectors. In *CVPR*, pages 7464–7475, 2023. 3
- [61] Xin Wang, Jiawei Wu, Junkun Chen, Lei Li, Yuan-Fang Wang, and William Yang Wang. VaTeX: A large-scale, high-quality multilingual dataset for video-and-language research. In *ICCV*, October 2019. 2
- [62] Yi Wang, Kunchang Li, Yizhuo Li, Yinan He, Bingkun Huang, Zhiyu Zhao, Hongjie Zhang, Jilan Xu, Yi Liu, Zun Wang, Sen Xing, Guo Chen, Junting Pan, Jiashuo Yu, Yali

- Wang, Limin Wang, and Yu Qiao. Internvideo: General video foundation models via generative and discriminative learning. 2022. 6
- [63] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *NeurIPS*, 35:24824–24837, 2022. 2, 3
- [64] Jialian Wu, Jianfeng Wang, Zhengyuan Yang, Zhe Gan, Zicheng Liu, Junsong Yuan, and Lijuan Wang. GReT: A generative region-to-text transformer for object understanding. *arXiv preprint arXiv:2212.00280*, 2022. 3, 6
- [65] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. MSR-VTT: A large video description dataset for bridging video and language. In *CVPR*, pages 5288–5296, 2016. 2
- [66] Zhengyuan Yang, Linjie Li, Kevin Lin, Jianfeng Wang, Chung-Ching Lin, Zicheng Liu, and Lijuan Wang. The Dawn of LMMs: Preliminary explorations with GPT-4V (ision). *arXiv preprint arXiv:2309.17421*, 2023. 6
- [67] Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*, 2019. 6
- [68] Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. Automatic chain of thought prompting in large language models. *arXiv preprint arXiv:2210.03493*, 2022. 3

## Appendix

In this supplementary, we present more details and discussions of AD generation with MM-Narrator (§A), our proposed complexity-based MM-ICL (§B), ablation studies of MM-ICL (§C) and AD evaluation with SegEval (§D). Next, we elaborate our implementation details (§E) and discuss the future work on both AD generation and evaluation (§F).

### A. AD Generation

MM-Narrator builds prompts to query GPT-4 for recurrent AD generation, including the following elements: task-specific introduction  $I_{\text{task}}$  and hint  $H_{\text{task}}$ , main query  $q_{\text{main}}$ , as well as a set of few-shot multimodal demonstrations  $\mathcal{D}_{\text{ICL}}$  to conduct in-context learning. With a breakdown shown in Figure 7, we present the details as follows.

**Querying with multimodal clues.** Both the main query  $q_{\text{main}}$  and the demonstration queries in  $\mathcal{D}_{\text{ICL}}$  are formatted with the same query builder, which outputs AD query from multiple text-formed multimodal clues. These multimodal clues include visual captions ( $x_i^{\text{cap}}$ ) with successfully re-identified characters, recent context ADs ( $\mathcal{M}_{\text{short}}$ ) and character dialogues ( $x_{t \in T_{\text{sub}}}^{\text{sub}}$ ).

**Prompting with MM-ICL.** Each MM-ICL demonstration within  $\mathcal{D}_{\text{ICL}}$ , is composed of a pair  $(Q, A)$  or a tuple  $(Q, \mathcal{R}, A)$  when chain-of-thought (CoT) is adopted to generate the multimodal multi-step reasoning  $\mathcal{R}$  that derives answer  $A$  from question  $Q$ .

**More qualitative results.** Apart from Figure 1 and 4 in main paper, we show additional qualitative demonstrations of MM-Narrator on both MAD-eval-Named benchmark and other long-form videos (external to the MAD-eval dataset) as Figure 10 and Figure 11, respectively, in this supplementary.

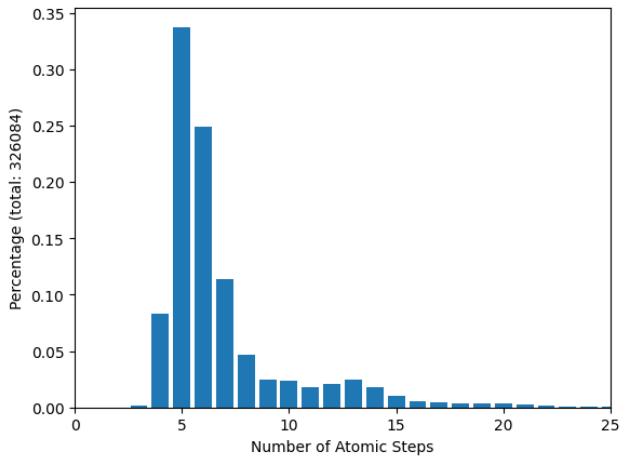
### B. Details of Complexity-based MM-ICL

Combining CoT with complexity-based ranking, our proposed complexity-based MM-ICL performs more favorably than classic ICL solutions. We reveal their details as follows.

**Reasoning with CoT.** We first employ GPT-4 to articulate the chain-of-thoughts (CoTs) as reasoning steps, denoted as  $\mathcal{R}$ , that assist in deriving the answer  $A$  from the question  $Q$ . Practically, we found a CoT-specific constraint<sup>2</sup> helpful to derive reliable CoTs, ensuring a closed-loop reasoning to be inferred. Without this constraint, LLM might unexpectedly generate  $\mathcal{R}$  followed by its own AD prediction, which are different from the human annotated  $A$ .

<sup>2</sup>CoT-specific constraint: “lets fill-in the REASONING process which derives the ANSWER from QUESTION.”

**Quantifying on atomic steps.** Practically, we observe that raw steps decided by LLM itself, might not be a considerably consistent measurement among various examples. Take two demonstrations shown in Figure 8 as example: Steps 3 to 7 in *left example*, conduct reasoning over per-frame captions individually, which are equivalent to step 2 in *right example*, including several sub-steps in analysing the per-frame captions. To this end, following [19], we split  $\mathcal{R}$  into *atomic steps* by newline char “\n”, and propose using the number of atomic steps  $N_{\text{atomic}}$  as our measurement of reasoning complexity.



**Figure 6.** Distributions of multimodal MAD-v2-Named demonstrations over reasoning complexity, quantified by  $N_{\text{atomic}}$ .

**Ranking by complexity.** We propose to select the most intuitive examples to perform few-shot MM-ICL for improving AD generation. Here, we show the distributions of multimodal demonstrations over the complexity in Figure 6. Specifically, the 10% shortest examples lead to a simple demonstration pool  $\mathcal{P}_{\text{simple}}$  with its maximum  $N_{\text{atomic}}$  as 5, while the 10% longest ones result in another pool  $\mathcal{P}_{\text{hard}}$  whose minimum  $N_{\text{atomic}}$  equals to 12.

### C. More Ablations on MM-ICL

Table 4 in the main manuscript implies that complexity is a suitable criterion for selecting efficient ICL demonstrations to enhance AD generation. Here, we further discuss *three sub-questions* to elaborate a few in-depth ablation studies, as following:

**Does CoT help?** We propose to adopt CoT technique to obtain the intermediate reasoning steps  $\mathcal{R}$  that help derive answer  $A$  from question  $Q$ . This automatic process extends demonstration format from  $(Q, A)$  pairs to  $(Q, \mathcal{R}, A)$  tuples. As its consistent gains could be observed multiple times (R1 vs R2; R3 vs R4; C1 vs C2), adding multimodal multi-step reasoning  $\mathcal{R}$  during MM-ICL could help MM-Narrator improve its multimodal reasoning capability

to better incorporate multimodal inputs. Qualitative demonstrations of  $\mathcal{R}$  are shown as Figure 8 in this supplementary.

**Does complexity-based ranking help?** We observed that conducting MM-ICL with the most intuitive examples benefits the overall performance (R4 vs C2), however, switching with the hardest ones which own the longest reasoning steps, MM-ICL actually leads to a decline in performance (R4 vs C3). These results indicate that more straightforward examples, quantified by the shortest number of reasoning steps, compile to a simpler yet more powerful subset MM-ICL demonstration pool for effective AD generation.

#### Does complexity-based MM-ICL work effectively?

Combining CoT with complexity-based ranking, our proposed complexity-based MM-ICL (C2) performs more favorably than the random and similarity-based sampling approaches (R1 [13] and S1 [33]), which are classic solutions in choosing few-shot ICL examples. Besides, ours is easy-to-implement and explainable-to-human, avoiding the computation overhead of retrieval-based selection.

## D. AD Evaluation with GPT-4

Suppose a few ADs form one segment Seg. For each Seg, our proposed SegEval evaluator takes into consideration an oracle context window Ctx of length  $W$ , to measure its AD quality with GPT-4. The details of SegEval prompt are shown as Figure 9. We elaborate each individual marking criteria as follows:

- **originality:** Evaluates if the Seg is novel and non-repetitive, to enrich the watching experience of the visually impaired.
- **consistency:** Checks if the generated Seg maintains a consistent tone or content throughout.
- **coherence:** Determines whether Seg logically connects to the given Ctx. A coherent text flows smoothly and deepen the movie understanding for the visually impaired.
- **diversity:** Focuses on the variety of Seg generated. A good model should produce varied outputs rather than repetitive or highly similar ones against the given Ctx.
- **specificity:** Measures the level of detail in the generated Seg, assessing if it is sufficiently detailed and/or focused for the Ctx.

Noticeably, the first two marking aspects focus on text-level AD quality, which are context-free ( $W = 0$ ) evaluation metrics, while the rest three metrics measure sequence-level AD generation, taking oracle context into consideration.

## E. Implementation Details

**Multimodal Experts.** To obtain framewise image caption and people detection, we utilize vision experts publicly available via the Azure Cognitive Services APIs<sup>3</sup>. For speech recognition, we choose WhisperX [10] as our audio expert. To register and recall long-term visual memory for character re-identification purpose, we adopt CLIP-ViT-L14 [48] as our visual feature extractor, and query GPT-4 as our Person-NER tool with the following prefix: “Extract the people names in the following text as a string splitted by ‘|’ (return ‘none’ if none of names are recognized): ”.

**Building MM-ICL Pool.** We build the MM-ICL demonstrations for each sample in MAD-v2-Named split [22]. As the raw frames are not publicly available, we derive per-frame captions by inferring ClipCap [41] on the released CLIP-ViT features. Differing from the main query  $q_{\text{main}}$ , whose recent context ADs in  $\mathcal{M}_{\text{short}}$  are recurrently generated by MM-Narrator, the queries in MM-ICL demonstrations  $\mathcal{D}_{\text{ICL}}$  are instead built with human annotations as their recent context ADs. Additionally, we omit long-term visual memory retrievals when constructing MM-ICL demonstrations.

**GPT-4 Error Handler.** GPT-4 might inevitably return errors when the content filtering policies<sup>4</sup> are occasionally triggered in Azure OpenAI Service. Such cases account for a very small proportion (less than 0.1%), thus they would not largely affect the overall performance. To address them, we utilize ClipCap [41] as the error handler to output video caption as AD. Specifically, we inference ClipCap on the mean pooled feature among frames in each video clip.

**Deployment on Long-form Videos.** We utilize PySceneDetect [5] for scene detection, and based on that, we cut long-form videos into video clips for recurrent AD generation with MM-Narrator. We utilize Google Text-to-Speech (gTTS) [3] for voice-over audio creation, which narrates AD for each video clip.

**Hyper-parameter Settings.** Following [22], the number of frames  $N$  to be sampled per video clip is set to 8, while we utilize subtitles within a time window  $T_{\text{sub}}$  set to 0.25 minutes. Our short-term memory queue is maintained to contain  $K$  most recently predicted ADs with timestamps, where  $K = 7$ . The number of demonstrations  $C$  equals to 5, which are sampled for conducting MM-ICL. The API versions of GPT-4 and GPT-4V used in our experiments are ‘gpt4-2023-03-15’ and ‘gpt4v-2023-08-01’, respectively.

<sup>3</sup><https://azure.microsoft.com/en-us/products/cognitive-services/vision-services>

<sup>4</sup><https://learn.microsoft.com/en-us/azure/ai-services/openai/concepts/content-filter>

## F. Future Works

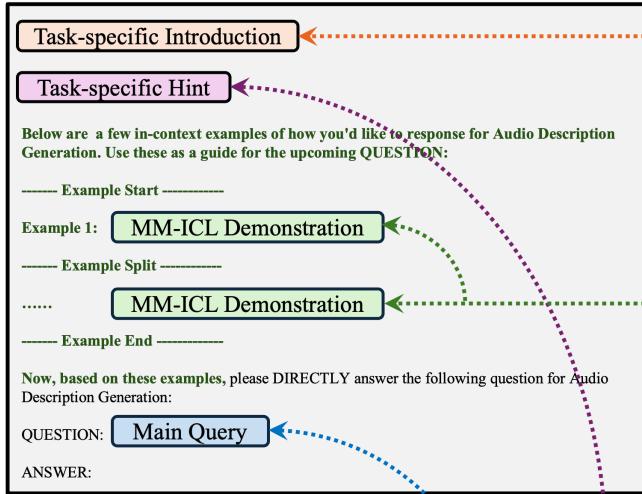
**AD Generation.** In future developments in Audio Description (AD), a critical enhancement will be the integration of advanced audio-visual speaker and character identification, coupled with strategic timing for AD delivery. This direction involves not only recognizing who is speaking or present in a scene but also determining the most opportune moments to provide descriptions without interrupting critical dialogue or action. Additionally, the establishment of a much more comprehensive and reliable external character bank, facilitating retrieval-augmented generation, will further refine AD content, ensuring it is both contextually relevant and timely. These advances are poised to transform AD into a more coherent, immersive experience, significantly improving accessibility for visually impaired audiences.

**AD Evaluation.** In future work for AD evaluation, a crucial focus should be on enhancing the measurement of factuality, an aspect not adequately addressed by current evaluation criteria like SegEval. Given the limitation of traditional reference-based scores in precisely assessing the factual accuracy of AD content, employing AI models such as GPT-4V emerges as a promising solution. GPT-4V’s advanced capabilities in understanding and contextualizing multimedia content could offer a more nuanced and accurate evaluation of AD factuality. This shift towards AI-driven, factuality-focused evaluation methods would not only provide a more comprehensive assessment of AD quality but also ensure that the generated descriptions are reliably accurate, ultimately benefiting visually impaired individuals with a more authentic storytelling experience.

## Acknowledgment

We are deeply grateful to OpenAI for providing access to their exceptional tool [44, 45]. We also extend heartfelt thanks to our Microsoft colleagues for their insights, with special acknowledgment to Faisal Ahmed, Ehsan Azarnasab, and Lin Liang for their constructive feedback.

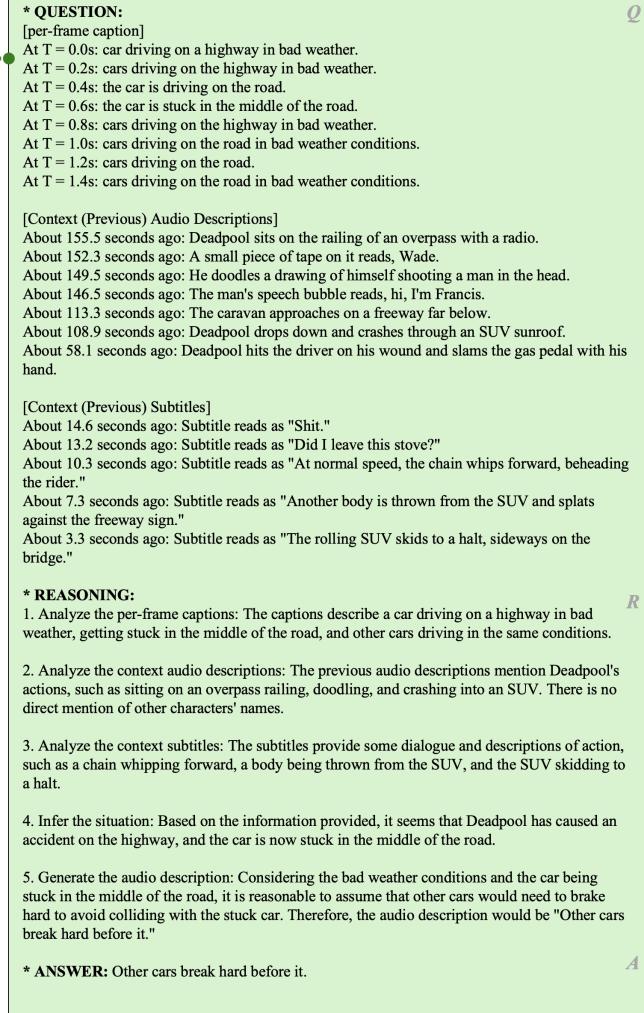
### (A) AD Generation Prompt (Overview)



### (B) Task-specific Introduction

Suppose you are an audio narrator who generates the audio descriptions (ADs) for blind people. However, instead of watching the videoclip, you will read the per-frame captions of the given videoclip, as well as the context ADs and/or subtitles before this videoclip. You must not narrate frame by frame. Instead, you should generate one-sentence brief AD covering the following videoclip. Note: the AD should be brief and concise and should not contain any redundant information. The length of the AD should match the length of the videoclip. Reminder: AD is not a video caption; thus, it MUST NOT contain words like "videoclip", "frames", "video". Specifically, the AD must NOT start with some common phrases like "This is a videoclip about ..." or "The videoclip describes/transitions/displays that ...".

### (D) MM-ICL Demonstration – $(Q, A)$ pair or $(Q, R, A)$ tuple



### (C) Task-specific Hint

(HINT: try to infer the character names from context ADs and subtitles, as well as the per-frame captions - where past memory might be recalled - for AD generation)

### (E) Main Query

**[per-frame caption]**  
At T = 0.0s: a hand holding a yellow plastic container  
At T = 0.3s: a hand holding a yellow plastic cup  
At T = 0.9s: a hand holding a yellow plastic container  
At T = 1.4s: a hand holding a yellow plastic container  
At T = 1.9s: a hand holding a yellow plastic container - *who could be a character we've seen before, especially around the previous scene about 21.9s ago, whose AD reads as "Lisa excitedly opens a gift from George."*  
At T = 2.4s: a hand holding a yellow plastic container  
At T = 2.9s: a hand holding a yellow container - *who could be a character we've seen before, especially around the previous scene about 22.9s ago, whose AD reads as "Lisa excitedly opens a gift from George."*  
At T = 3.4s: a hand holding a yellow container - *who could be a character we've seen before, especially around the previous scene about 23.4s ago, whose AD reads as "Lisa excitedly opens a gift from George."*

**[Context (Previous) Audio Descriptions]**  
About 25.0 seconds ago: The man in a suit continues his conversation with Lisa as George watches.  
About 21.4 seconds ago: Lisa excitedly opens a gift from George.  
About 17.7 seconds ago: Lisa smiles as she discovers the contents of the gift from George.  
About 12.4 seconds ago: George and Lisa engage in conversation at the table.  
About 9.2 seconds ago: The man in a suit converses with Lisa, while George observes.  
About 6.1 seconds ago: The man in a suit presents a gift to Lisa as George looks on.  
About 4.4 seconds ago: Lisa receives a gift from the man in a suit while George watches, then she looks at a yellow cup.

**[Context (Previous) Subtitles]**  
About 10.7 seconds ago: Subtitle reads as "Thanks for not rushing me."  
About 7.2 seconds ago: Subtitle reads as "Yeah."

### (F) Quantitative and Qualitative Analysis



How do you know (2010)  
Start: 01:48:29 | End: 01:48:33

**GT** (via human annotation)

**As Lisa examines it, her diamond watch glitters on her wrist.**

**PD** (via MM-Narrator)

**Lisa examines the yellow container from George.**

**R-L: 21.4 | C: 78.0 | M: 13.0 | B1: 16.1**

**Figure 7.** A breakdown of the AD generation prompt constructed by MM-Narrator, including an (A) overview with ICL-specific instructions marked in green, (B) task-specific introduction  $I_{task}$  and (C) hint  $H_{task}$ , a few multimodal ICL (MM-ICL) demonstrations  $D_{ICL}$  with an example shown as (D), and (E) the main query  $q_{main}$  to be responded by GPT-4, with long-term visual memory marked in gray. Eventually, we show the corresponding (F) quantitative and qualitative analysis of the AD prediction via MM-Narrator against the human AD annotation. Zoom in for details.

## (A) Prompting LLM to articulate CoTs as reasoning steps

You are an audio narrator who generates the audio descriptions (ADs) for blind people. However, instead of watching the videoclip, you will read the per-frame captions of the given videoclip, as well as the context ADs and/or subtitles before this videoclip.

### Task-specific Hint

Now, given the below (QUESTION, ANSWER) pair example of AD generation, lets fill-in the REASONING process which derives the ANSWER from Question.

\* Question: Question

\* Answer: Answer

\* Reasoning: Lets think of step-by-step...

## (B) Two examples of deriving CoTs with LLM

[per-frame caption]  
 At T = 0.0s: close up of a burning firewood in a fireplace.  
 At T = 0.2s: a fire burns in a fireplace.  
 At T = 0.4s: slow motion of a burning flame in a metal bowl.  
 At T = 0.6s: a fire burns in slow motion.  
 At T = 0.8s: abstract background of a yellowish - orange background.  
 At T = 1.0s: person in a dark forest.  
 At T = 1.2s: a man is standing in front of a burning fireplace and is looking at the flame.  
 At T = 1.4s: slow motion of a girl running in the dark.

[Context (Previous) Audio Descriptions]  
 About 33.2 seconds ago: Another hand grabs Doris's wrist.  
 About 30.9 seconds ago: Lane's eyes clear and Debbie appears beside her.  
 About 27.3 seconds ago: The blonde girl stares intensely at the ghoul.  
 About 22.2 seconds ago: Doris flicks her eyes between the two friends.  
 About 12.2 seconds ago: Sarah emerges and throws the bundle into the raging furnace.  
 About 8.4 seconds ago: A shockwave throws Sarah back.  
 About 6.2 seconds ago: Doris' mouth hangs open and black smoke billows from it.

Lane throws the board into the furnace, surging flames throw her back against a wooden beam.

- Analyze the per-frame captions and context audio descriptions to understand the scene and characters involved.
- From the context audio descriptions, we can identify the characters as Doris, Lane, Debbie, and Sarah.
- Now, let's focus on the per-frame captions to understand the actions happening in the scene.
- At T = 0.0s to T = 0.8s, we see various descriptions of fire and flames, indicating that the scene is focused on a burning fireplace or furnace.
- At T = 1.0s, there's a mention of a person in a dark forest, but it's not clear who this person is or how they relate to the scene.
- At T = 1.2s, a man is standing in front of a burning fireplace, looking at the flame. This could be a character from the context audio descriptions, but we don't have enough information to determine who it is.
- At T = 1.4s, there's a slow-motion scene of a girl running in the dark. This could be one of the female characters from the context audio descriptions, such as Lane, Debbie, or Sarah.
- Considering the previous audio descriptions, we know that Sarah threw a bundle into the furnace, and a shockwave threw her back. This could be related to the current scene with the burning fireplace.
- Based on the information available, we can infer that Lane is the one who throws the board into the furnace, as she is one of the characters mentioned in the context audio descriptions and is likely the girl running in the dark at T = 1.4s.
- The surging flames and the shockwave mentioned in the context audio descriptions suggest that Lane is thrown back against a wooden beam after throwing the board into the furnace.

Thus, the audio description generated is: "Lane throws the board into the furnace, surging flames throw her back against a wooden beam."

[per-frame caption]  
 At T = 0.0s: person in a scene from the movie.  
 At T = 0.2s: a young woman is looking through a magnifying glass in a dark room.  
 At T = 0.4s: person in the dark : person.  
 At T = 0.6s: film character and person in the dark.  
 At T = 0.8s: the man in the dark.  
 At T = 1.0s: the man in the dark.  
 At T = 1.2s: the team in the library.  
 At T = 1.4s: the car is parked in the dark.

[Context (Previous) Audio Descriptions]  
 About 453.2 seconds ago: The handsome sailor sheepishly makes his way to the door.  
 About 448.9 seconds ago: Two dancers sandwich him between their bottoms and he bounces.  
 About 438.1 seconds ago: A director descends on a crane.  
 About 358.1 seconds ago: In the beach house living room...  
 About 128.2 seconds ago: Baird gazes off.  
 About 126.2 seconds ago: In an office, Eddie approaches a man wearing glasses.  
 About 9.5 seconds ago: He enters an editing room.

[Context (Previous) Subtitles]  
 About 13.1 seconds ago: Subtitle reads as "Why?"  
 About 11.9 seconds ago: Subtitle reads as "Never mind."  
 About 11.1 seconds ago: Subtitle reads as "Okay."  
 About 10.3 seconds ago: Subtitle reads as "Tell him I'll be there."  
 About 8.4 seconds ago: Subtitle reads as "Thank you, old lady."  
 About 7.6 seconds ago: Subtitle reads as "Oh, well, what is it?"  
 About 5.9 seconds ago: Subtitle reads as "Hello, Cece."  
 About 5.5 seconds ago: Subtitle reads as "I lost the thought we might go to later."  
 About 2.8 seconds ago: Subtitle reads as "Oh, hi, Eddie."  
 About 1.4 seconds ago: Subtitle reads as "Wanna lace up what you got on Barlowy Dance?"

It's up now.

1. Identify the character names and context from the previous audio descriptions and subtitles: We have Eddie, Cece, Baird, and a man wearing glasses.

2. Analyze the per-frame captions:

- At T = 0.0s: person in a scene from the movie.
- At T = 0.2s: a young woman is looking through a magnifying glass in a dark room.
- At T = 0.4s: person in the dark : person.
- At T = 0.6s: film character and person in the dark.
- At T = 0.8s: the man in the dark.
- At T = 1.0s: the man in the dark.
- At T = 1.2s: the team in the library.
- At T = 1.4s: the car is parked in the dark.

3. Generate the audio description based on the per-frame captions and context:

- "Eddie and Cece are in a dark room, where a young woman is looking through a magnifying glass. The man in the dark, possibly Baird or the man wearing glasses, is also present. The team gathers in the library, and a car is parked outside in the dark."

4. Final Answer: "It's up now."

**Figure 8.** Reasoning with CoT. (A) An overview of prompting LLM to articulate CoTs as multimodal multi-step reasoning  $\mathcal{R}$  that derives the answer  $\mathcal{A}$  from question  $\mathcal{Q}$ , with (B) two examples shown below. Zoom in for details.

### (A) AD evaluation prompt with SegEval

Suppose you are a visually impaired person, and you will be "watching" a movie videoclip with audio description (AD). Here, you are requested to provide feedback (via reasoning and marking) on the performance of two AI assistants ("ASSISTANT1" and "ASSISTANT2") for automatic AD generation task:

#### Evaluation Steps:

1. you will be given <Context ADs>, <ASSISTANT1-output>, and <ASSISTANT2-output>, where <Context ADs> shows a few contextual human-annotated ADs but leaves "<PRESENT-SEGMENT>" empty to be filled with one or multiple AD(s) generated by two AI assistants (i.e., <ASSISTANT1-output> and <ASSISTANT2-output>).
2. you will read though the <Context ADs> and <ASSISTANT1-output> and <ASSISTANT2-output>, and then measure the AD generation quality of the two AI assistants in terms of [ Metric ] aspect.
3. you will complete the following five sections IN ORDER (namely, <Assistant1-Reasoning>, <Assistant2-Reasoning>, <Comparison-Reasoning>, <Assistant1-Score>, and <Assistant2-Score>).

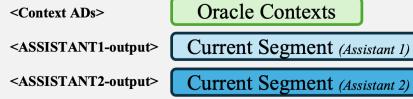
#### HINT:

1. <Context ADs> will be used to provide the context of the movie scene. If it contains no valid ADs, it means that the current evaluation metric ([ Metric ]) will not take contextual information into account.
2. <Assistant1-Reasoning> and <Assistant2-Reasoning> will be used to record your reasoning and comments (with supporting evidence) on the [ Metric ] aspect of the ADs generated by two AI assistants, respectively;
3. <Comparison-Reasoning> will be used to record your feedback (with supporting evidence) for comparisons between the two AI assistants (with respect to the [ Metric ] aspect), which will be used to support the below two marking sections;
4. <Assistant1-Score> and <Assistant2-Score> will be used to record your AD generation [ Metric ] scores (from "1" to "5", where "1" indicates the worst and "10" indicates the excellent) of the two AI assistants, respectively.

#### Evaluation Criteria:

#### Criteria

#### Evaluation Starts



Please make sure you read and understand these instructions carefully, and complete the following five sections IN ORDER:

- (1) firstly reason them individually within "<Assistant1-Reasoning>" and "<Assistant2-Reasoning>";
- (2) secondly compare two assistants within "<Comparison-Reasoning>"; and
- (3) finally mark them within "<Assistant1-Score>" and "<Assistant2-Score>"

### (B) One example of Diversity with L = 5 and W = 3.

- **Diversity**: Focuses on the variety of <PRESENT-SEGMENT> generated. A good model should produce varied outputs rather than repetitive or highly similar ones against the given <Context ADs>.

About 57.0s ago: The Marines share a look.  
About 43.0s ago: The staff sergeant wipes the wounded civilian's brow then sits back and unscrews a water bottle.  
About 14.5s ago: Joe lifts his trembling hand.  
About 11.2s ago: Doc eyes his superior.  
About 8.9s ago: Nantz glances at the Marine and grips the wounded man's hand.  
\* At present: —PLACEHOLDER for "<PRESENT\_SEGMENT>" to be generated by AI assistants below—  
About 62.0s later: On screen, a giant craft hovers.  
About 75.0s later: Name meets her gaze and Santos nods.  
About 81.0s later: They put Joe into the backroom.  
About 83.5s later: Nantz compares a clock with his watch.  
About 87.9s later: Harris checks his watch.

Now: The staff sergeant holds his gaze.  
About 9.3s later: Nantz grabs his helmet.  
About 12.9s later: Doc unfolds a foil blanket.  
About 17.7s later: Now news footage.  
About 34.0s later: Santos watches on.

Now: Cpl. Lee Imlay and Ssgt. Michael Nantz exchange concerned glances as they look down at something, surrounded by their fellow soldiers.  
About 9.3s later: Nantz lies on the ground, tended to by Santos and fellow soldiers.  
About 12.9s later: Santos tends to Nantz's wound as he lies on the floor, surrounded by concerned soldiers.  
About 17.7s later: TSgt. Elena Santos stands silhouetted, raising her arms and pointing at a tree.  
About 34.0s later: Santos and Mottola intently analyze screens, emotions running high.

#### <Assistant1-Reasoning>

Assistant1's output provides diverse descriptions of the scene, covering different actions and characters. The generated ADs include Nantz holding his gaze, grabbing his helmet, Doc unfolding a foil blanket, news footage, and Santos watching on. These descriptions are varied and not repetitive, showing a good level of diversity.

#### <Assistant2-Reasoning>

Assistant2's output also provides diverse descriptions of the scene, covering different actions and characters. The generated ADs include Nantz and Imlay exchanging concerned glances, Nantz lying on the ground being tended to, Santos tending to Nantz's wound, Santos raising her arms and pointing at a tree, and Santos and Mottola analyzing screens. These descriptions are varied and not repetitive, showing a good level of diversity.

#### <Comparison-Reasoning>

Both Assistant1 and Assistant2 provide diverse and varied descriptions of the scene, covering different actions and characters. Assistant1's output focuses more on Nantz's actions and reactions, while Assistant2's output provides more context about the situation and includes more characters. Assistant2's output seems to be slightly more diverse in terms of the actions and characters described, but both assistants perform well in terms of diversity.

<Assistant1-Score> 8

<Assistant2-Score> 9

Post-Processing (note: GPT-4 is agnostic to the output source during above marking)  
Model Score: 1.125

### (C) One example of Originality with L = 5 and W = 1.

- **Originality**: Evaluates if the <PRESENT-SEGMENT> is novel and non-repetitive, to enrich the watching experience of the visually impaired.

\* At present: —PLACEHOLDER for "<PRESENT\_SEGMENT>" to be generated by AI assistants below—

Now: Children play and interact with a cat on the porch of a house.  
About 14.2s later: The children gather on the porch, enjoying their time together.  
About 14.9s later: The children sit together on the porch, petting the dog and chatting.  
About 15.6s later: The children sit together on the porch, petting the dog and chatting.  
About 18.2s later: The children gather around the dog in the yard.

Now: BO carefully pours the water into a metal dog bowl, and slides it across the grass to Houdini.  
About 14.2s later: The Alsatian suddenly jerks his head up towards BO  
About 14.9s later: BO jumps back in shock.  
About 15.6s later: And MORGAN frowns.  
About 18.2s later: BO stands.

#### <Assistant1-Reasoning>

The ADs generated by ASSISTANT1 are not very original. The descriptions are repetitive, especially the two instances of "The children sit together on the porch, petting the dog and chatting." This repetition does not enrich the watching experience for the visually impaired.

#### <Assistant2-Reasoning>

ASSISTANT2's ADs are more original and non-repetitive. Each description provides a unique action or event, such as BO pouring water into a dog bowl, the Alsatian jerking its head up, and BO jumping back in shock. These descriptions provide a richer watching experience for the visually impaired.

#### <Comparison-Reasoning>

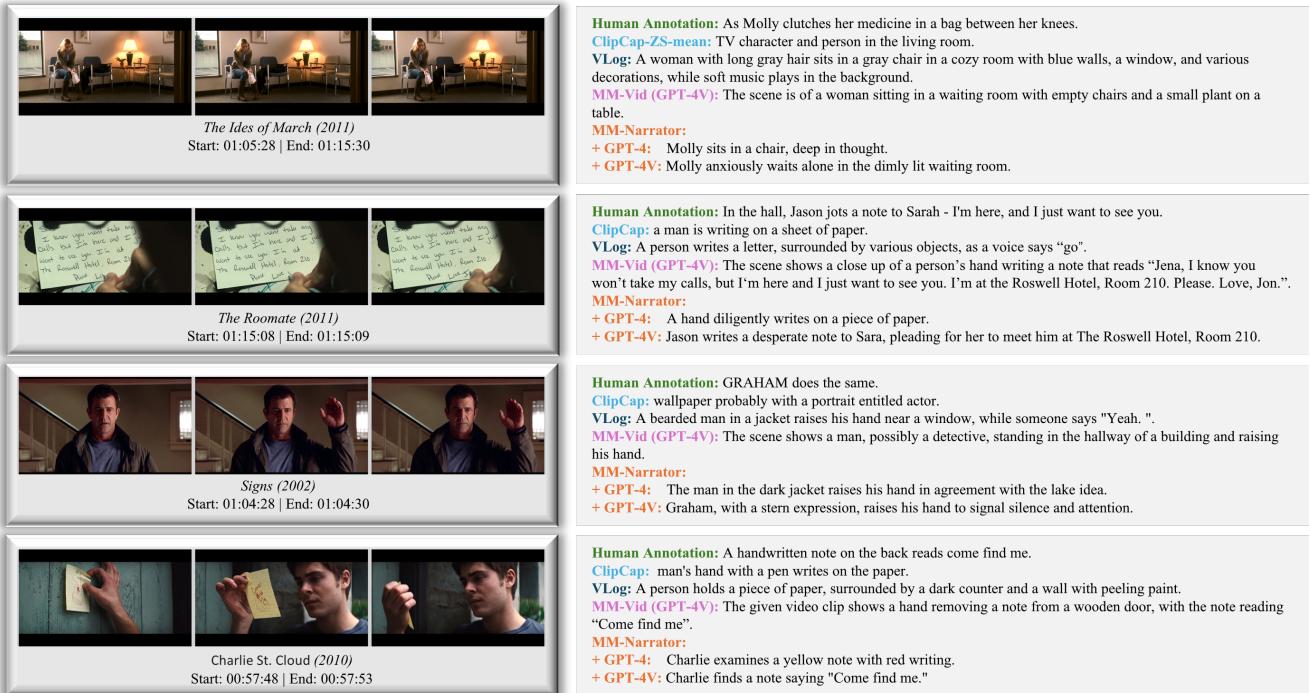
Comparing the two AI assistants, ASSISTANT2's ADs are more original and non-repetitive than ASSISTANT1's. ASSISTANT1's descriptions are repetitive and do not provide a rich watching experience for the visually impaired. On the other hand, ASSISTANT2's descriptions are unique and provide a better understanding of the scene.

<Assistant1-Score> 2

<Assistant2-Score> 8

Post-Processing (note: GPT-4 is agnostic to the output source during above marking)  
Model Score: 0.25

**Figure 9.** AD evaluation with SegEval. (A) An overview of prompting GPT-4 to evaluate AD generation quality, with (B) one *diversity* and (C) one *originality* examples shown below. Zoom in for details.



**Figure 10.** More qualitative comparisons on MAD-eval-Named benchmark. For example, in *The ides of March* (2011), our method generates AD by conditioning on current video clip and the contextual information from timestamp 00:00:00 to 01:05:28. Zoom in for details.



**Figure 11.** More qualitative demonstrations of MM-Narrator on other long-form videos. For example, in *Inception* (2010), our method generates AD by conditioning on current video clip and contextual information from timestamp 00:00:00 to 00:31:52. Zoom in for details.