

Visual Summary Thought of Large Vision-Language Models for Multimodal Recommendation

Yuqing Liu*

yliu363@uic.edu

University of Illinois at Chicago

Chicago, IL, USA

Lichao Sun

lis221@lehigh.edu

Lehigh University

Bethlehem, PA, USA

Yu Wang*

ywang617@uic.edu

University of Illinois at Chicago

Chicago, IL, USA

Philip S. Yu

psyu@uic.edu

University of Illinois at Chicago

Chicago, IL, USA

ABSTRACT

The evolution of large vision-language models (LVLMs) has shed light on the development of many fields, particularly for multimodal recommendation. While LVLMs offer an integrated understanding of textual and visual information of items from user interactions, their deployment in this domain remains limited due to inherent complexities. First, LVLMs are trained from enormous general datasets and lack knowledge of personalized user preferences. Second, LVLMs struggle with multiple image processing, especially with discrete, noisy, and redundant images in recommendation scenarios. To address these issues, we introduce a new reasoning strategy called Visual-Summary Thought (VST) for Multimodal Recommendation. This approach begins by prompting LVLMs to generate textual summaries of item images, which serve as contextual information. These summaries are then combined with item titles to enhance the representation of sequential interactions and improve the ranking of candidates. Our experiments, conducted across four datasets using three different LVLMs: GPT4-V, LLaVA-7b, and LLaVA-13b validate the effectiveness of VST.

KEYWORDS

Large Vision-Language Models, Multimodal Recommendation, Reasoning Strategy

1 INTRODUCTION

To address the cold-start issues that recommender systems lack sufficient records of new items/users, multimodal recommender systems (MMRSs) [5, 8, 14, 27, 29, 35, 37, 41, 42] are proposed by involving the complementary content of items from multiple perspectives, e.g., textual description and visual illustration, thus enriching the recommender system’s knowledge. However, the knowledge of MMRSs is primarily learned from scratch using a limited user-item

*Both authors contributed equally to this research.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).
CIKM '24, October 21–25, 2024, Boise, ID, USA.
© 2024 Copyright held by the owner/author(s).
ACM ISBN 979-8-4007-0436-9/24/10
<https://doi.org/10.1145/XXXXXX.XXXXXX>

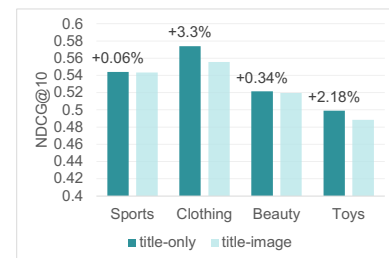


Figure 1: The performance of GPT4-V on four representative Amazon datasets with title-only and title-image concatenation inputs.

interaction dataset that is often biased and noisy [2, 15, 32, 40]. Additionally, the product image provided by the seller contains critical marketing highlights that attract buyers, e.g., the game’s duration and thematic ambiance, elements that traditional embedding-based MMRSs may struggle to effectively capture. Moreover, traditional MMRSs encounter challenges in fusing multimodal knowledge, where inefficient integration can further degrade the recommender system’s performance [11, 13, 36, 39].

Meanwhile, the remarkable success of large vision-language models (LVLMs) [7, 12, 17, 25, 28, 33, 34, 43] offers encouraging solutions to the above issues encountered by traditional MMRSs. LVLMs are proficient in comprehending both textual and visual information about an item, owing to their training on enormous datasets. Their ability to distill and adapt item information across modalities into natural language space exhibits an opportunity for effective knowledge fusion. Despite these strengths, the incorporation of pretrained LVLMs into MMRSs remains an under-explored area. Two possible obstacles may hinder the widespread adoption of LVLMs in MMRSs:

First, *LVLMs are trained from vast general knowledge* and, as such, lack domain-specific knowledge for understanding user preferences revealed through their interactions. This gap results in the under-exploration of LVLMs’ capacity in recommendation scenarios. To bridge this gap, it is essential to integrate additional knowledge to inform LVLMs in the context necessary for making appropriate recommendations. This approach, however, introduces the second challenge: *LVLMs’ inefficiency in processing multiple images*. Although models like GPT4-V have been evaluated in video understanding scenarios to examine their capacity in capturing dynamic content across frames [1, 16, 20, 22, 28, 31, 33], the scenario

with MMRSs involves handling multiple, discrete, and noisy images. This complexity can pose a significant challenge even from a human perspective, making it difficult to extract meaningful knowledge from such diverse interactions. Our preliminary experiments as shown in Figure 1 indicate this issue, showing that a simple concatenation of multiple images with item titles performs worse than methods relying solely on item titles for recommendations even with powerful GPT4-V. Furthermore, current reasoning algorithms, e.g., in-context learning (ICL) [3, 9, 18, 19, 30] and chain-of-thought (CoT) [10, 21, 23, 26, 38], are primarily designed for NLP tasks ignoring visual modality. However, the principal challenge in multimodal recommendation is how to effectively leverage image-based knowledge and integrate it into the recommendation process. Thus, effective LVLM-based MMRS requires the design of specific prompting strategies that can utilize their visual comprehension strength without caving to the complexities associated with processing multiple images simultaneously.

Accordingly, we propose a novel **Visual-Summary Thought (VST)** reasoning principle of LVLMs for MMRSs. Our approach includes two primary components: First, we utilize user historical interactions as contextual data for the LVLMs' personalized recommendations. This involves using sequences of both item titles and images as inputs to the LVLMs. Second, to overcome the shortage of handling multiple images, we prompt the LVLMs with one static image to obtain a corresponding textual summary. Then, we construct user history sequences by substituting the images with their textual comprehensions one by one, serving as an intermediate representation for LVLMs during the reasoning phase. This strategy allows for the recommendation based on a more manageable comprehension of user preferences, transitioning from the complex and noisy image sequences to a simpler task of understanding visual-summary enhanced preference dynamics. To validate the efficacy of our proposed reasoning algorithm, we conduct experiments using GPT4-V, LLaVA-7b, and LLaVA-13b as reasoning backbones. We observe consistent improvements over other existing reasoning strategies, such as concatenation, ICL, and CoT. Our contributions can be summarized as follows:

- To the best of our knowledge, this is the first attempt to investigate the reasoning strategies for LVLMs in multimodal recommendation scenarios.
- We introduce a novel Visual-Summary Thought (VST) reasoning strategy, specifically designed for the multimodal recommendation context, to harness the proficiency of LVLMs' visual understanding and remedy their deficiency in handling multiple images simultaneously.
- We conduct comprehensive experiments to evaluate VST, utilizing both API-based LVLMs like GPT4-V, and open-source models such as LLaVA-7b and LLaVA-13b. The consistent improvements observed across these models demonstrate the effectiveness of VST for LVLM-based MMRSs.

2 METHODOLOGY

2.1 Problem definition

In this paper, we follow the problem settings in [6, 24] that use the pretrained LVLMs as reranker to make recommendations to user u via reranking the given n candidate item titles $v = \{v_1, v_2, \dots, v_n\}$.

For each user, we have their historical interactions, which is the sequence of title and image pair of items: $u = \{(t_1, i_1), (t_2, i_2), \dots, (t_m, i_m)\}$.

2.2 Preliminary

LVLMs exhibit limitations in handling multiple images. We evaluated the LVLMs' ability to handle multimodal inputs by concatenating the item titles and images of user histories. Surprisingly, leveraging complementary visual information led to poorer results compared to only using item titles as shown in Figure 1. (An example can be found in section 3.4.) This underscores a critical insight: adding more information to the LVLMs' prompt context without a thoughtful design can lead to confusion, especially with discrete and noisy images full of redundancy. To address this challenge, we introduce a novel visual-summary thought of prompting strategy (VST) as shown in Figure 2.

2.3 Visual-Summary Generation

Existing LVLMs, e.g., GPT4-V and LLaVA, primarily focus on static image understanding scenarios, where LVLMs generate textual descriptions of a given image. However, this paradigm is inefficient for handling multiple images [?]. Existing strategies include concatenating images for LVLM reasoning [?], or adapting LVLMs to video comprehension scenarios via finetuning on video datasets [17, 20, 25, 34]. Yet, neither approach is suitable for the unique demands of MMRSs, where the image sequence of a user history is discrete and noisy, lacking the continuous nature of video frames and making sequential correlations difficult to discern. To deal with these issues, we propose leveraging LVLMs' strengths in temporal understanding within natural language modality and their capacity for static image interpretation. Instead of processing a sequence of images, we focus on distilling critical marketing highlights from individual image. The prompt can be formalized as: $s_i = \text{summary}(i) = \text{"What's in this image?"}$ For each item, we use one image and get the summarization of each image independently. In this way, we can not only obtain marketing highlights of items via distilling image comprehension from LVLMs but also simplify the temporal user preference understanding from the visual modality to the textual modality, where the LVLMs demonstrate proficiency.

2.4 Visual-Summary Thought for MMRSs

After summarizing each item image, we concat the history item titles with their visual summary to construct the prompt for querying user preferences among candidates. The prompt is structured in two parts: the first outlines the user's purchase history in chronological order, demonstrated by each item's title and visual summary. The second segment directs the LVLMs to rerank the candidates represented by their titles. An illustrative prompt might be:

"[Here is a chronological list of my purchase history for some products including the title and the description of each product. $\{(t_1, s_{i_1}), \dots, (t_m, s_{i_m})\}$][There are $|n|$ candidate products I am considering to buy: $\{v_1, \dots, v_n\}$. Please rank these $|n|$ candidate products based on the likelihood that I would like to purchase next most according to the given purchase history. You cannot generate products that are not in the given candidate list.]"

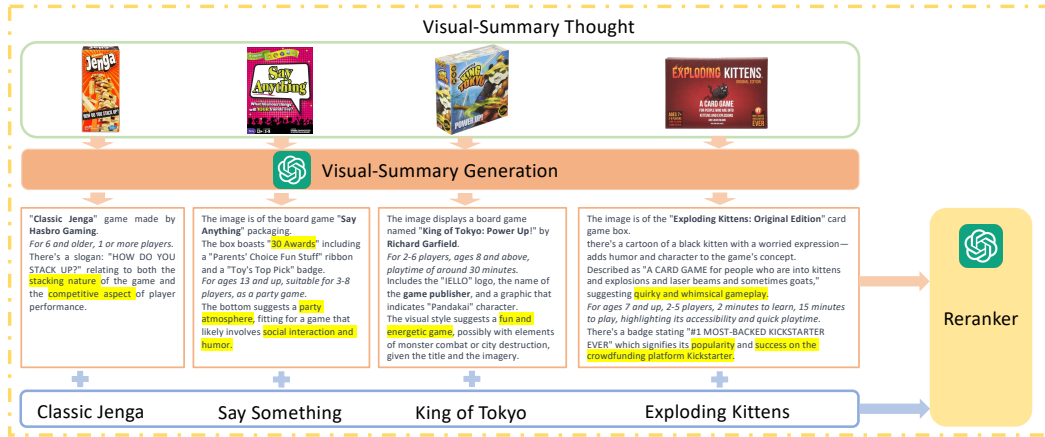


Figure 2: Framework of Visual-Summary Thought of LVLMS for Multimodal Recommendation. Text in yellow highlights some key features obtained through visual-summary generation.

Table 1: Performance comparison of different prompt strategies. Target items are guaranteed to be included in the candidate sets. We highlight the best and the second-best results.

Dataset	Metric	GPT4-V				LLaVA-7b				LLaVA-13b			
		MM	MM-ICL	MM-CoT	VST	MM	MM-ICL	MM-CoT	VST	MM	MM-ICL	MM-CoT	VST
Sports	R@5	0.6900	<u>0.6950</u>	0.5750	0.7250	0.1300	<u>0.1900</u>	0.1800	0.3283	0.2250	<u>0.3300</u>	0.2300	0.3750
	R@10	<u>0.8600</u>	0.8600	0.8150	0.9000	0.2950	<u>0.3400</u>	0.3250	0.5067	0.3200	<u>0.4850</u>	0.3250	0.6250
	R@20	<u>0.8700</u>	0.8650	0.8300	0.9050	0.3100	<u>0.3500</u>	0.3550	0.5117	0.3400	<u>0.5000</u>	0.3450	0.6350
	N@5	0.4880	<u>0.5126</u>	0.4186	0.5263	0.0703	<u>0.1138</u>	0.1043	0.1769	0.1395	<u>0.2087</u>	0.1393	0.2244
	N@10	0.5435	<u>0.5666</u>	0.4961	0.5834	0.1243	<u>0.1619</u>	0.1506	0.2345	0.1706	<u>0.2598</u>	0.1701	0.3063
	N@20	0.5461	<u>0.5678</u>	0.4999	0.5846	0.1281	<u>0.1646</u>	0.1580	0.2357	0.1755	<u>0.2637</u>	0.1752	0.3086
Clothing	R@5	0.6550	0.7100	0.6300	<u>0.6950</u>	0.1400	0.1650	<u>0.1700</u>	0.2800	<u>0.3650</u>	0.3200	0.2550	0.3950
	R@10	0.8950	<u>0.9050</u>	0.8150	0.9300	0.2750	<u>0.3100</u>	0.2600	0.3250	0.6700	0.5450	0.4200	0.6200
	R@20	0.9000	<u>0.9050</u>	0.8200	0.9350	0.2900	<u>0.3150</u>	0.2600	0.3250	0.6950	0.5450	0.4200	0.6250
	N@5	0.4781	0.5580	0.4631	<u>0.5322</u>	0.0851	<u>0.1156</u>	0.1086	0.1875	<u>0.2248</u>	0.2062	0.1554	0.2594
	N@10	0.5555	0.6205	0.5238	<u>0.6085</u>	0.1287	<u>0.1633</u>	0.1386	0.2025	<u>0.3234</u>	0.2787	0.2058	0.3329
	N@20	0.5569	0.6205	0.5252	<u>0.6098</u>	0.1326	<u>0.1646</u>	0.1386	0.2025	<u>0.3301</u>	0.2787	0.2085	0.3343
Beauty	R@5	0.6300	0.6300	0.5500	<u>0.6200</u>	<u>0.2450</u>	0.1800	0.1450	0.2750	0.2650	<u>0.2900</u>	0.2300	0.3200
	R@10	0.8450	<u>0.8700</u>	0.6400	0.9000	0.4050	0.3150	0.1700	<u>0.4000</u>	0.3750	<u>0.4200</u>	0.3200	0.5500
	R@20	0.8500	<u>0.8750</u>	0.6500	0.9000	0.4200	0.3200	0.1750	<u>0.4000</u>	0.3850	<u>0.4200</u>	0.3250	0.5600
	N@5	<u>0.4503</u>	0.4395	0.3964	0.4536	<u>0.1484</u>	0.1202	0.1006	0.1769	0.1641	<u>0.1928</u>	0.1398	0.2183
	N@10	<u>0.5197</u>	0.5183	0.4264	0.5439	<u>0.1996</u>	0.1641	0.1087	0.2179	0.2008	<u>0.2361</u>	0.1692	0.2942
	N@20	<u>0.5211</u>	0.5195	0.4290	0.5439	<u>0.2035</u>	0.1655	0.1101	0.2179	0.2033	<u>0.2361</u>	0.1706	0.2970
Toys	R@5	0.5500	0.6450	0.4950	<u>0.6300</u>	<u>0.1450</u>	0.1150	0.1300	0.3000	0.1875	<u>0.3400</u>	0.2600	0.3617
	R@10	0.7650	<u>0.7800</u>	0.6950	0.8000	<u>0.2750</u>	0.1450	0.1700	0.3800	0.2550	<u>0.4250</u>	0.3800	0.5150
	R@20	0.7750	<u>0.7800</u>	0.7050	0.8000	<u>0.2850</u>	0.1550	0.1850	0.3950	0.2663	<u>0.4350</u>	0.3800	0.5200
	N@5	0.4184	0.4789	0.3967	<u>0.4399</u>	<u>0.0857</u>	0.0842	0.0835	0.2035	0.1389	<u>0.2373</u>	0.1832	0.2412
	N@10	0.4883	0.5227	0.4349	<u>0.4958</u>	<u>0.1281</u>	0.0941	0.0977	0.2299	0.1614	<u>0.2648</u>	0.2228	0.2919
	N@20	0.4911	0.5227	0.4376	<u>0.4958</u>	<u>0.1305</u>	0.0966	0.1015	0.2336	0.1642	<u>0.2672</u>	0.2228	0.2932

3 EXPERIMENTS

In this section, we provide the performance comparison between the proposed VST and three representative reasoning strategies on four public datasets, using GPT4-V, LLaVA-7b, and LLaVA-13b as pretrained LVLMS.

3.1 Experimental Settings

Dataset. In this paper, we adopt the same dataset as in [4] that uses the Amazon Review datasets for evaluation. Due to the limitation of the inference rate, following the common practice [6], we

Table 2: Statistics of the datasets after sampling.

Datasets	#Users	#Items	#Interactions	Sparsity
Sports	200	1750	2333	99.33%
Clothing	200	1291	1362	99.47%
Beauty	200	2024	2797	99.31%
Toys	200	1684	1967	99.42%

only sample 200 users for evaluation. We report the statistics of such datasets in Table 2.

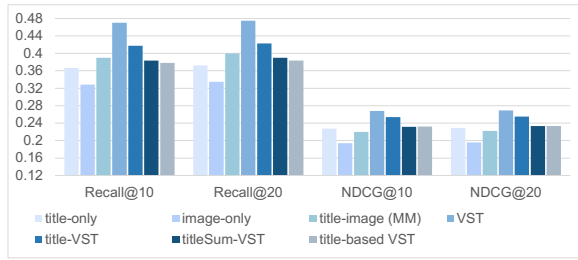


Figure 3: Ablation study. Performance of LLaVA-13b with different prompts on Toys dataset.

Metrics. We adopt Recall@K (R@K) and NDCG@K (N@K) to evaluate the ranking performance of the LVLMS over candidate items, which consist of the title of the ground-truth (target) item and the 9 random sampled items following [6].

Implementation Details. For open-source LVLMS, we use Fastchat to launch models and conduct the model inference on a single GeForce RTX 4090.

Baseline Models. As there is no previous work that only utilizes the inference capacity of LVLMS for multimodal recommendation, we adopt the commonly chosen prompting strategies used in NLP tasks: in-context-learning and chain-of-thought for comparison. **MM:** The plain prompt, using the simple concatenation of the historical item titles and images as the first segment. The second part keeps the same as VST. **MM-ICL:** For ICL, we match each prefix of the user’s historical interaction sequence with its corresponding successor as demonstration examples. For example: “[Here is a chronological list of my purchase history: $\{(t_1, i_1), \dots, (t_{m-1}, i_{m-1})\}$] [Then if I ask you to recommend a new product, you should recommend t_m . Now I’ve just purchased t_m , I want to buy a new product...]”. The remaining part is the same as the second part of VST. **MM-CoT:** For CoT, we adopt zero-shot CoT by adding “Please think step by step.” to the second part of the prompt, while the first part is the same as MM. For example: “[Here is a chronological list of my purchase history: $\{(t_1, i_1), \dots, (t_m, i_m)\}$] [There are $|n|$ candidate products I am considering to buy ... Please think step by step by considering my preferences based on the given titles and image sequence of the purchased products...]”.

3.2 Overall Performance

To demonstrate the effectiveness of our proposed VST strategy, we employ GPT4-V, LLaVA-7b, and LLaVA-13b as pretrained LVLMS and conduct experiments with four different prompt strategies across four datasets. The complete experimental results are shown in Table 1. From the table, we can observe that our proposed VST reasoning strategy achieves the best or comparable performances across all datasets, demonstrating the effectiveness of our approach. Notably, our approach has a better performance on Sports dataset than others. This might be due to the titles of this category of products containing much more noise, making the alignment between textual and visual information more challenging for the employed LVLMS. In contrast, through visual-summary generation, VST can

better leverage visual modality and capture more relevant information from the image, reducing the impact of the noise from different modalities to some extent.

3.3 Ablation Study

To analyze the effectiveness of the VST reasoning principle, we conduct an ablation study on six variants of the proposed strategy. The results on Toys dataset using LLaVA-13b are shown in Figure 3. The reported results are the average of a minimum of three repeated runs, aimed at minimizing the impact of randomness. **titleSum-VST** refers to the prompt that also lets LVLMS distill information from the title of an item: $s_t = \text{summary}(t) = \text{“What information can you get from the title?”}$, then appended by the summary distilled from the corresponding image. **title-based VST** refers to instructing LVLMS to distill information from an image by taking item title into consideration, where $s_i = \text{summary}(i) = \text{“This is an image related to } t. \text{ Please provide a detailed description of the given image.”}$

From the results, we have the following observations: (1) VST can capture more meaningful information from both textual and visual modalities. The results show that VST has the capability to significantly enhance the ranking performance compared to non-VST-based strategies. The improvement stems from VST’s proficiency in multimodal understanding and serves better in sequential scenarios, where information from different sources needs to be integrated effectively. (2) Information from the title can boost performance, but it depends on the quality of the title and the alignment between the title and the image. Compared to the results among VST, title-VST, titleSum-VST, and title-based VST, we can observe that adding the title information doesn’t yield improvement. This lack of improvement is likely due to the visibility of toy titles in images or the easy identification of entities mentioned in titles from the images themselves. Therefore, combining title information with VST does not provide substantial additional benefits. Whether to include titles during reasoning remains a hyperparameter decision dependent on the quality of titles in each dataset.

3.4 Case Study

In this section, we compare the ranking lists generated by LLaVA-13b using VST with title-only and title-image concatenation prompts. The results are shown in Figure 4. Here are our observations from comparing the outputs: Both title-only and VST strategies successfully rank the target item as the first position, while the naive concatenation of title and image places it fourth. This discrepancy suggests that raw images may contain an excess of information, which could be perceived as redundant and introduce additional noise into our ranking task. On the other hand, the VST strategy offers a more refined approach. By utilizing VST, we not only incorporate information from the title but also extract richer and more relevant details from the image itself. Such details also align closely with the marketing selling points of the product. Consequently, the VST strategy emerges as a more effective prompt for multimodal recommendation, as it combines textual and visual cues to provide a comprehensive understanding of the item, thereby enhancing the performance of the ranking results.




User's Historical Interaction Instruction (UHI)		
Here is a chronological list of my purchase history for some toys-related products including the <i>title / title and image / image description</i> of each product.		
Title	Title-Image	Image Description
1. Mastermind 2. Say Anything 3. My First Lab Duo-Scope Microscope 4. King of Tokyo Power Up Expansion Game 5. Foam Maverick Pogo Stick 6. Helicopter with Gyro 7. Volcano Making Kit	1.  Mastermind 2.  Say Anything 7.  Volcano Making Kit	1. The image shows the game "Mastermind," a well-known board game . The packaging indicates that this game is meant for two players , who are typically aged 8 and above . It's a logic game . 2. The image appears to be the cover of a board game , called "Say Anything." It's a party game designed for 3-8 players who are 13 years of age or older . The bottom indicates the party atmosphere of the game. The cover also boasts that the game has won 30 awards , signaling its popularity and recognition in the gaming community. 7. The image displays a science kit , including a segmented dish,, a plastic volcanic structure. It is designed for educational purposes to model volcanic eruption , potentially for school-age children as a learning tool .
Candidate Reranking Instruction (CRI)		
There are 10 candidate products I am considering to buy: {... Flyer Scooter, Paint Cups with Color-Coded Lids, Don't Let the Pigeon Drive the Bus Game , Wear Charms Spectacular Spinner ...} Please rank these 10 candidate products that I would like to purchase next most according to the given purchase history.		
CASE STUDY		
Title-only Input: UHI + Title + CRI Output: 1. Don't Let the Pigeon Drive the Bus Game 2. TabletTopics Family: Questions to Start Great Conversations 3. Kid Chuck Bumper Cars 10. Wikki Stix Big Count Box	Title-Image Input: UHI + Title-Image + CRI Output: 1. Paint Cups with Color-Coded Lids 4. Don't Let the Pigeon Drive the Bus Game 10. Wikki Stix Big Count Box	VST Input: UHI + Image Description + CRI Output: 1. Don't Let the Pigeon Drive the Bus Game 2. TabletTopics Family: Questions to Start Great Conversations 3. Paint Cups with Color-Coded Lids 10. Flyer Scooter

Figure 4: Case study. Text in red indicates the target item. Text in orange, purple, or blue indicates the pattern to describe the item for the corresponding prompt. Text in yellow highlights some key features obtained through visual-summary generation.

4 CONCLUSION

In this work, we investigate the performance of different reasoning strategies for LVLMs in multimodal recommendation scenarios and identify a notable limitation in LVLMs' capability to effectively handle multiple images. To bridge this gap, we propose a Visual-Summary Thought (VST) strategy to distill information from images. By leveraging LVLMs' visual understanding, VST aims to harness their strengths while rectifying deficiencies in handling multiple images. Extensive experiments conducted on four real-world datasets using both API-based LVLMs such as GPT4-V and open-source models like LLaVA-7b and LLaVA-13b, consistently demonstrate the effectiveness of our proposed VST.

REFERENCES

- [1] Yunkang Cao, Xiaohao Xu, Chen Sun, Xiaonan Huang, and Weiming Shen. 2023. Towards generic anomaly detection and understanding: Large-scale visual-linguistic model (gpt-4v) takes the lead. *arXiv preprint arXiv:2311.02782* (2023).
- [2] Xu Chen, Hanxiong Chen, Hongteng Xu, Yongfeng Zhang, Yixin Cao, Zheng Qin, and Hongyuan Zha. 2019. Personalized fashion recommendation with visual explanations based on multimodal attention network: Towards visually explainable recommendation. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval*. 765–774.
- [3] Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, and Zhifang Sui. 2022. A survey for in-context learning. *arXiv preprint arXiv:2301.00234* (2022).
- [4] Shijie Geng, Juntao Tan, Shuchang Liu, Zuohui Fu, and Yongfeng Zhang. 2023. VIP5: Towards Multimodal Foundation Models for Recommendation. In *Findings of the Association for Computational Linguistics: EMNLP 2023*. 9606–9620.
- [5] Ruining He and Julian McAuley. 2016. VBPR: visual bayesian personalized ranking from implicit feedback. In *AAAI*, Vol. 30.
- [6] Yupeng Hou, Junjie Zhang, Zihan Lin, Hongyu Lu, Ruobing Xie, Julian McAuley, and Wayne Xin Zhao. 2024. Large language models are zero-shot rankers for recommender systems. In *ECIR*. Springer, 364–381.
- [7] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597* (2023).
- [8] Shuaiyang Li, Dan Guo, Kang Liu, Richang Hong, and Feng Xue. 2023. Multimodal Counterfactual Learning Network for Multimedia-based Recommendation. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1539–1548.
- [9] Xiaonan Li and Xipeng Qiu. 2023. Finding supporting examples for in-context learning. *arXiv preprint arXiv:2302.13539* (2023).
- [10] Zhan Ling, Yunhao Fang, Xuanlin Li, Zhao Huang, Mingyu Lee, Roland Memisevic, and Hao Su. 2023. Deductive Verification of Chain-of-Thought Reasoning. *arXiv preprint arXiv:2306.03872* (2023).
- [11] Fan Liu, Huilin Chen, Zhiyong Cheng, Anan Liu, Liqiang Nie, and Mohan Kankanhalli. 2022. Disentangled multimodal representation learning for recommendation. *IEEE Transactions on Multimedia* (2022).
- [12] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2024. Visual instruction tuning. *Advances in Neural Information Processing Systems* 36 (2024).
- [13] Han Liu, Yinwei Wei, Fan Liu, Wenjie Wang, Liqiang Nie, and Tat-Seng Chua. 2023. Dynamic Multimodal Fusion via Meta-Learning Towards Micro-Video Recommendation. *TOIS* 42, 2 (2023), 1–26.
- [14] Kang Liu, Feng Xue, Dan Guo, Le Wu, Shujie Li, and Richang Hong. 2023. Megcf: Multimodal entity graph collaborative filtering for personalized recommendation. *ACM Trans. Recomm. Syst.* 41, 2 (2023), 1–27.
- [15] Xiaohao Liu, Zhulin Tao, Jiahong Shao, Lifang Yang, and Xianglin Huang. 2022. EliMRec: Eliminating Single-modal Bias in Multimedia Recommendation. In *Proceedings of the 30th ACM International Conference on Multimedia*. 687–695.
- [16] Ming Y Lu, Bowen Chen, Drew FK Williamson, Richard J Chen, Kenji Ikamura, Georg Gerber, Ivy Liang, Long Phi Le, Tong Ding, Anil V Parwani, et al. 2023. A Foundational Multimodal Vision Language AI Assistant for Human Pathology. *arXiv preprint arXiv:2312.07814* (2023).
- [17] Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. 2024. Video-ChatGPT: Towards Detailed Video Understanding via Large Vision and Language Models. In *ACL*. 12585–12602.
- [18] Sewon Min, Xinxin Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2022. Rethinking the role of demonstrations: What makes in-context learning work?. In *EMNLP*. 11048–11064.
- [19] Ohad Rubin, Jonathan Herzig, and Jonathan Berant. 2021. Learning to retrieve prompts for in-context learning. *arXiv preprint arXiv:2112.08633* (2021).
- [20] Yunlong Tang, Jing Bi, Siting Xu, Luchuan Song, Susan Liang, Teng Wang, Daoan Zhang, Jie An, Jingyang Lin, Rongyi Zhu, et al. 2023. Video Understanding with Large Language Models: A Survey. *arXiv preprint arXiv:2312.17432* (2023).
- [21] Harsh Trivedi, Niranjana Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022. Interleaving retrieval with chain-of-thought reasoning for knowledge-intensive multi-step questions. *arXiv preprint arXiv:2212.10509* (2022).
- [22] Rohan Wadhawan, Hritik Bansal, Kai-Wei Chang, and Nanyun Peng. 2024. Contextual: Evaluating Context-Sensitive Text-Rich Visual Reasoning in Large Multimodal Models. *arXiv preprint arXiv:2401.13311* (2024).
- [23] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. Self-Consistency Improves Chain of Thought Reasoning in Language Models. In *11th International Conference on Learning Representations*.
- [24] Yu Wang, Zhiwei Liu, Jianguo Zhang, Weiran Yao, Shelby Heinecke, and Philip S Yu. 2023. DRDT: Dynamic Reflection with Divergent Thinking for LLM-based Sequential Recommendation. *arXiv preprint arXiv:2312.11336* (2023).
- [25] Zhanyu Wang, Longyue Wang, Zhen Zhao, Minghao Wu, Chenyang Lyu, Huayang Li, Deng Cai, Luping Zhou, Shuming Shi, and Zhaopeng Tu. 2023. Gpt4video: A unified multimodal large language model for instruction-followed understanding and safety-aware generation. *arXiv preprint arXiv:2311.16511* (2023).

- [26] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems* 35 (2022), 24824–24837.
- [27] Yinwei Wei, Xiang Wang, Liqiang Nie, Xiangnan He, Richang Hong, and Tat-Seng Chua. 2019. MMGCN: Multi-modal graph convolution network for personalized recommendation of micro-video. In *Proceedings of the 27th ACM international conference on multimedia*. 1437–1445.
- [28] Licheng Wen, Xueming Yang, Daocheng Fu, Xiaofeng Wang, Pinlong Cai, Xin Li, Tao MA, Yingxuan Li, Linran XU, Dengke Shang, Zheng Zhu, Shaoyan Sun, Yeqi BAI, Xinyu Cai, Min Dou, Shuanglu Hu, Botian Shi, and Yu Qiao. 2024. On the Road with GPT-4V(ision): Explorations of Utilizing Visual-Language Model as Autonomous Driving Agent. In *ICLR 2024 Workshop on Large Language Model (LLM) Agents*.
- [29] Chuhan Wu, Fangzhao Wu, Tao Qi, Chao Zhang, Yongfeng Huang, and Tong Xu. 2022. Mm-rec: Visiolinguistic model empowered multimodal news recommendation. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 2560–2564.
- [30] Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. 2021. An explanation of in-context learning as implicit bayesian inference. *arXiv preprint arXiv:2111.02080* (2021).
- [31] Zhiling Yan, Kai Zhang, Rong Zhou, Lifang He, Xiang Li, and Lichao Sun. 2023. Multimodal chatgpt for medical applications: an experimental study of gpt-4v. *arXiv preprint arXiv:2310.19061* (2023).
- [32] Wei Yang, Zhengru Fang, Tianle Zhang, Shiguang Wu, and Chi Lu. 2023. Modal-aware Bias Constrained Contrastive Learning for Multimodal Recommendation. In *Proceedings of the 31st ACM International Conference on Multimedia*. 6369–6378.
- [33] Zhengyuan Yang, Linjie Li, Kevin Lin, Jianfeng Wang, Chung-Ching Lin, Zicheng Liu, and Lijuan Wang. 2023. The dawn of llms: Preliminary explorations with gpt-4v (ision). *arXiv preprint arXiv:2309.17421* 9, 1 (2023), 1.
- [34] Hang Zhang, Xin Li, and Lidong Bing. 2023. Video-llama: An instruction-tuned audio-visual language model for video understanding. In *EMNLP*. 543–553.
- [35] Jinghao Zhang, Yanqiao Zhu, Qiang Liu, Shu Wu, Shuhui Wang, and Liang Wang. 2021. Mining latent structures for multimedia recommendation. In *Proceedings of the 29th ACM International Conference on Multimedia*. 3872–3880.
- [36] Jinghao Zhang, Yanqiao Zhu, Qiang Liu, Mengqi Zhang, Shu Wu, and Liang Wang. 2022. Latent structure mining with contrastive modality fusion for multimedia recommendation. *IEEE Transactions on Knowledge and Data Engineering* (2022).
- [37] Lingzi Zhang, Xin Zhou, and Zhiqi Shen. 2023. Multimodal Pre-training Framework for Sequential Recommendation via Contrastive Learning. *arXiv preprint arXiv:2303.11879* (2023).
- [38] Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. 2022. Automatic chain of thought prompting in large language models. *arXiv preprint arXiv:2210.03493* (2022).
- [39] Hongyu Zhou, Xin Zhou, Zhiwei Zeng, Lingzi Zhang, and Zhiqi Shen. 2023. A Comprehensive Survey on Multimodal Recommender Systems: Taxonomy, Evaluation, and Future Directions. *arXiv preprint arXiv:2302.04473* (2023).
- [40] Xin Zhou and Zhiqi Shen. 2023. A tale of two graphs: Freezing and denoising graph structures for multimodal recommendation. In *Proceedings of the 31st ACM International Conference on Multimedia*. 935–943.
- [41] Xin Zhou, Hongyu Zhou, Yong Liu, Zhiwei Zeng, Chunyan Miao, Pengwei Wang, Yuan You, and Feijun Jiang. 2023. Bootstrap latent representations for multi-modal recommendation. In *Proceedings of the ACM Web Conference 2023*. 845–854.
- [42] Yan Zhou, Jie Guo, Hao Sun, Bin Song, and Fei Richard Yu. 2023. Attention-guided multi-step fusion: a hierarchical fusion network for multimodal recommendation. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1816–1820.
- [43] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. 2024. Minigpt-4: Enhancing vision-language understanding with advanced large language models. In *12th International Conference on Learning Representations*.