

From Text to Pixel: Advancing Long-Context Understanding in MLLMs

Yujie Lu¹ Xiujun Li² Tsu-Jui Fu¹ Miguel Eckstein¹ William Yang Wang¹

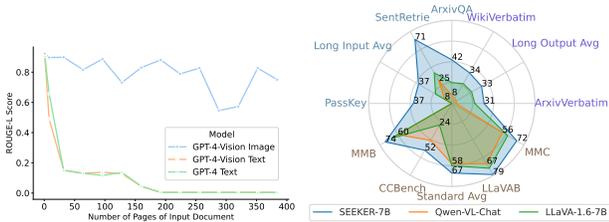


Figure 1: Left: Performance plot on First-Sentence-Retrieval task revealing compact nature of image tokens in representing long content. Right: Radar chart demonstrating the superior performance of the SEEKER (ours) model across both short and long-context multimodal tasks.

Abstract

The rapid progress in Multimodal Large Language Models (MLLMs) has significantly advanced their ability to process and understand complex visual and textual information. However, the integration of multiple images and extensive textual contexts remains a challenge due to the inherent limitation of the models’ capacity to handle long input sequences efficiently. In this paper, we introduce SEEKER, a multimodal large language model designed to tackle this issue. SEEKER aims to optimize the compact encoding of long text by compressing the text sequence into the visual pixel space via images, enabling the model to handle long text within a fixed token-length budget efficiently. Our empirical experiments on six long-context multimodal tasks demonstrate that SEEKER can leverage fewer image tokens to convey the same amount of textual information compared with the OCR-based approach, and is more efficient in understanding long-form multimodal input and generating long-form textual output, outperforming all existing proprietary and open-source MLLMs by large margins.

¹University of California, Santa Barbara, Santa Barbara, CA, USA ²University of Washington, Seattle, WA, USA. Correspondence to: Yujie Lu <yujielu@ucsb.edu>.

1. SEEKER: Long-context Vision and Language Understanding

We propose SEEKER, a multimodal large language model designed to handle long-context images and texts. In Section 1.1, we discuss the innovative use of image tokens to represent lengthy textual data compactly. Then we introduce long-context multimodal task and instruction data in Section 1.2. Finally, in Section 1.3, we illustrate the architecture of our SEEKER to support both long-context and short-context multimodal understanding.

1.1. Using Image Tokens to Encode Text Helps Context Length Extrapolation

We follow the approach outlined in (Xiong et al., 2023) to evaluate model’s extrapolation capability in the First-Sentence-Retrieval task. In this task, models are required to retrieve the first sentence at a specific length. We conduct this synthetic task on various numbers of documents with different page counts. We probe the performance of GPT-4-Vision Image by feeding its images of documents and compare it with GPT-4-Vision Text and GPT-4, which receive extracted text using the OCR model Nougat (Blecher et al., 2023). Nougat achieves over a 90 BLEU score on OCR text from scientific documents. All these models have a context length limit of 128k tokens.

On the left side of Figure 1, we visualize the Rouge-L (Lin, 2004) score in relation to the total number of pages of input documents, which range from 1 (approximately 1k text tokens) to 448 (approximately 500k text tokens). We observe a significant performance degradation in models fed with text input. In contrast, without any additional changes, we see improved extrapolation when representing length text content with visual tokens by feeding images of documents directly to the model.

1.2. Long-Context Multimodal Task

We mainly consider two categories of long-context multimodal capabilities, as outlined in Table 1: 1) Long-form multimodal input: This involves multiple text-rich images interleaved with text as the input context. 2) Long-form text output: This requires generating long text.

Table 1: **Long-Context Multimodal Task.** $\text{Img}/\#\text{In}$: the number of input images, $\text{Text Tok}/\#\text{In}$ and $\#\text{Out}$: the number of input and output text tokens. Full examples are presented in Appendix D.1.

Task	Prompt Example	Img	Text Tok	
		#In.	#In.	#Out.
<i>Long-Form Multi-Image Input</i>				
Index	Which Image contains the given sentence?	6.6	100.4	1.0
SentRetrie	What is the first sentence on the first image?	1.0	23.0	35.5
ArxivQA	What is the main purpose of the article as stated in the abstract?	8.2	13.9	35.0
PassKey	What is the ;PASSKEY _i in the provided images?	4.0	95.4	2.6
<i>Long-Form Text Output</i>				
ArxivVerb	Read the text in the image verbatim.	1.0	10.0	1301.6
WikiVerb	Read the text in the image verbatim.	1.0	16.0	1107.1

Instruction Data for Long-Form Multi-Image Input

First, we combine an arbitrary number of single-image visual instruction data (Liu et al., 2023c) sourced from CC3M into the multi-image format for the intra-image reasoning task. This helps initiate model’s capability of understanding sequences of images (e.g., img_1 *This image depicts a...* img_2 *This image shows a...*). To enable understanding of long-form text-rich image sequences, we collect compiled PDFs from arXiv documents. Each page from these documents is processed as images, ranging from 4 to 24 pages. We use GPT-4V to generate descriptive or conversational instruction data for these scientific documents. To further improve the model’s understanding of each provided image, we create a multi-image text grounding task, requiring the model to ground the question to the referred image (e.g., img_1 img_2 ... img_8 *Which image contains the answer to the question / Which image contains the sentence...*).

Instruction Data for Long-Form Text Output To enhance long-form text generation capabilities related to the given image, we propose a task that involves reading the text in the image verbatim (e.g., img_1 *Quote the text in the image verbatim.*). This challenging task requires the vision backbone to encode character-level image details and the language backbone to attend to the image token while producing very long text without hallucinating on previously generated content.

1.3. Long-Context Multimodal Large Language Model

To enable long-context multimodal reasoning, our model architecture should: 1) encode multiple images interleaved with text, 2) align images and text at a fine-grained level, and 3) decode long texts that attend to extended multimodal contexts. The following paragraphs illustrate the design of our proposed SEEKER for this purpose.

Long-Context Multi-Image Encoding For effective feature integration in scenarios involving multiple images, it is crucial to include image separators to concatenate text and

image sequences as:

$$\text{Query} = \text{Query}_{\text{system}} + \sum_{i=1}^N (\mathbf{Q}_{\text{img},i} + \mathbf{Q}_{\text{txt},i})$$

$$\mathbf{Q}_{\text{img},i} = \text{start}(\text{img}, i) + \text{content}(\text{img}, i) + \text{end}(\text{img}, i) \quad (1)$$

Specifically, we use $\text{start}(\text{img},i)$ and $\text{end}(\text{img},i)$ as special tokens ‘ startofimg_i ’ and ‘ endofimg_i ’ to distinguish the start and end of each image, respectively. We observe this strategy is essential for maintaining model performance, especially when training is limited to a small dataset of long-context multimodal instructions. The encoding process and the concatenation of the feature vectors of the input sequence can be described as:

$$t_i = \text{Enc}_t(T_i), v_i = \text{MLP}_{v \rightarrow t}(\text{Enc}_v(I_i)) \quad (2)$$

$$Q = [t_0; v_1; t_1; v_2; t_2; \dots; v_n; t_n]$$

Here, Enc_v encodes each image i into a feature vector and projects it to the word embedding space. The concatenated vector Q integrates sequences of image and text feature vectors, where $[\cdot; \cdot]$ denotes concatenation along the feature dimension.

Dense Image-Text Alignment We inherit the general image-text alignment from the pre-training image-text pairs. To enhance the visual representation of dense text in images, and improve the alignment between image and text representation of rendered text, we curate a visual-embedded task that renders text into visual space.

Supervised Fine-tuning Strategy We aim to leverage sequential data processing to fine-tune models on a combination of textual and visual inputs, enabling them to generate coherent and contextually relevant responses based on both text and image data, using autoregressive training objective.

2. Main Results

2.1. Long Image and Text Context

Long-Form Multi-Image Input In Table 3, SEEKER achieves significantly surpassing larger open-source MLLMs across all four long-form multi-image input tasks. We concatenate the images for models that can not handle image sequences. Additionally, SEEKER-TINY ranks second best. On average, our models also outperform the proprietary GPT-4V model. This indicates our auxiliary tasks, as detailed in Section 1.2, enhance the models’ reasoning across multiple images and grounding content to specific images. Thus our models excel at handling long-context tasks involving long-form multiple text-rich image inputs.

Table 2: Short Image and Text Context. : proprietary models, : the proposed models.

Models	Multi-Image			Single-Image								
	NLVR2	BLINK	Avg	MMB	MMC	SEED	CCBench	AI2D	LLaVAB	ChartQA	TextVQA	Avg
Close-source MLLMs												
GPT-4V (OpenAI, 2023b)	71.7	51.1	61.4	75.1	74.4	71.6	46.5	75.9	93.1	78.5	78.0	60.3
Open-source MLLMs												
Qwen-VL-Chat (Bai et al., 2023b)	30.8	28.1	29.5	60.6	56.3	64.8	41.2	63.0	67.7	49.8	60.7	58.0
LLaVA-1.5-7B (Liu et al., 2023a)	61.7	37.1	49.4	65.2	59.0	65.8	27.5	55.5	61.8	17.8	45.4	49.8
LLaVA-Next-7B (Liu et al., 2024)	58.7	41.2	49.9	67.4	62.3	69.6	24.3	67.0	72.7	55.4	64.4	60.4
LLaVA-Next-7B (Mistral) (Liu et al., 2024)	43.5	37.5	40.5	69.5	61.3	72.4	30.0	<u>69.0</u>	67.8	51.8	65.2	63.1
DeepSeek-VL-7B (Lu et al., 2024)	46.6	40.9	43.7	74.1	71.4	70.4	<u>51.7</u>	65.3	77.8	59.1	64.9	<u>66.8</u>
IDEFICS2-8B (Laurençon et al., 2024)	79.9	46.8	63.4	75.3	67.3	71.9	37.6	72.3	49.1	24.36	68.9	66.3
Monkey-Chat-10B (Li et al., 2023)	66.0	40.5	53.3	71.0	65.8	68.9	48.4	68.5	60.5	<u>59.5</u>	<u>65.5</u>	63.5
LLaVA-1.5-13B (Liu et al., 2023a)	66.2	<u>42.7</u>	54.4	69.2	65.0	68.2	30.4	61.1	66.1	18.2	48.9	53.4
LLaVA-Next-13B (Liu et al., 2024)	64.3	42.6	53.4	70.7	79.0	<u>71.9</u>	28.8	72.2	73.9	61.4	66.9	65.6
Open-source Tiny MLLMs												
DeepSeek-VL-1.3B (Lu et al., 2024)	61.3	38.8	50.1	64.0	62.9	66.0	37.6	51.5	51.1	47.4	57.8	54.8
MiniCPM-V-3B (Hu et al., 2024)	63.1	40.0	51.5	67.9	62.6	65.6	41.4	56.3	51.3	44.2	56.6	55.7
Ours												
SEEKER-TINY-1.3B	69.9	40.5	55.2	64.8	63.7	66.0	37.3	49.0	81.7	45.4	56.3	58.0
SEEKER-7B	<u>72.4</u>	42.1	<u>57.2</u>	<u>74.0</u>	<u>72.6</u>	71.1	52.0	64.6	<u>79.3</u>	58.3	65.3	67.1

Long-Form Text Output In Table 3, our SEEKER achieves the best performance for long-context tasks requiring long-form text output. On average, LLaVA-Next (Liu et al., 2024)-13B also performs well, likely because these tasks usually require a single image. Its feature of splitting images into four tiles as additional 2304 image tokens, combined with the original image, greatly enhances its ability to capture visual details. This is particularly beneficial for verbatim tasks involving Arxiv and Wikipedia content rendered in the image. Meanwhile, DeepSeek-VL (Lu et al., 2024) achieves the best scores among other open-source 7B MLLMs, primarily due to its alignment of image and text by enforcing text reading from a large scale of visual-situated real-world data, such as documents and PDFs. By incorporating our small-scale verbatim task data, which includes images rendered with text of various font sizes, into the instruction-tuning stage, our models achieve a 38.1% performance improvement.

Fix-length Image Tokens are more Expressive than Text Tokens If a model can interpret text within images, it confirms that this method is a valid way to present information. Additionally, if the model requires fewer image tokens than text tokens to understand the text, this indicates that pixels can represent text more compactly. To investigate this, we conduct a probing task involving question-answering using various pages of documents fed into the model, as shown in Table 4. Notably, in this task, we use a version of our SEEKER with the same context length as the compared model, which is 4,096 tokens. Our observations indicate that when the text token count is up to around 4,000, the response accuracy remains within the context length limit of 4,096 tokens without performance degradation for the

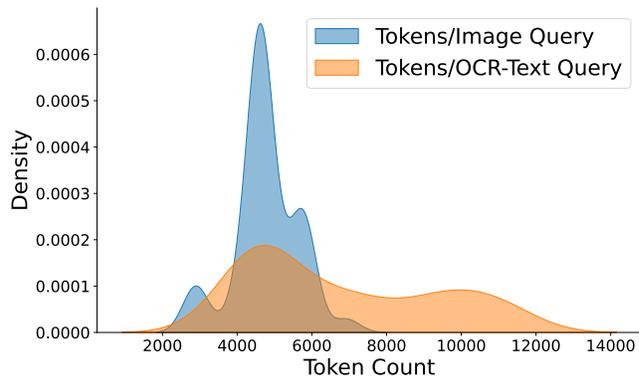


Figure 2: Density plot comparing token counts for image token (blue) and OCR-text (orange) representations. Image tokens are more compact than text, fitting well within 8192 context length.

language model (LLM). When the text token count exceeds 4,000 but the image token count remains below 4,000, the vision-language model (VLM) outperforms the LLM by 4 to 8 percentage points. However, when the image token count exceeds 4,000, the performance of the VLM also declines, though it remains slightly superior to that of the LLM.

3. Analysis

3.1. Context Length Extrapolation

We analyze the effectiveness of using image tokens versus OCR text tokens for image representation. The density plot in Figure 2 illustrates the distribution of token counts for both methods. The Image token representation is notably

Table 3: **Long Image and Text Context.** : proprietary models, : the proposed models, #Tok/Img: the number of tokens per image. We report accuracy on multiple-choice task Index, and Rouge-L score for other tasks.

Models	Params	#Tok/Img	Long-Form Multi-Image Input					Long-Form Text Output		
			Index	SentRetrie	ArxivQA	PassKey	Avg	ArxivVerb	WikiVerb	Avg
Close-source MLLMs										
GPT-4V (OpenAI, 2023b)	–	85	32.50	71.10	45.19	27.16	43.98	32.58	5.96	19.27
Open-source MLLMs										
Qwen-VL-Chat (Bai et al., 2023b)	7B	–	2.49	25.05	8.24	0.00	8.94	4.90	5.41	5.15
LLaVA-1.5 (Liu et al., 2023b)	7B	576	23.74	30.61	35.60	0.00	22.48	4.14	3.80	3.97
LLaVA-Next (Liu et al., 2024)	7B	2880	17.49	34.35	20.50	0.00	18.08	22.33	22.94	22.63
LLaVA-Next (Mistral) (Liu et al., 2024)	7B	2880	17.49	34.45	21.39	0.00	18.33	20.11	20.92	20.51
DeepSeek-VL (Lu et al., 2024)	7B	576	13.74	10.37	19.83	0.17	11.02	<u>31.59</u>	16.48	24.03
IDEFICS2 (Laurençon et al., 2024)	8B	64	10.83	63.46	9.68	0.13	21.02	12.12	5.93	9.02
Monkey-Chat (Li et al., 2023)	10B	–	16.24	23.65	17.90	0.00	14.44	5.82	2.08	3.95
LLaVA-1.5 (Liu et al., 2023a)	13B	576	22.49	41.02	32.31	0.00	23.95	9.57	7.12	8.34
LLaVA-Next (Liu et al., 2024)	13B	2880	11.24	37.55	15.60	0.00	16.09	27.14	<u>31.05</u>	<u>29.09</u>
Open-source Tiny MLLMs										
DeepSeek-VL (Lu et al., 2024)	1.3B	576	14.99	10.46	21.29	0.15	11.72	20.06	10.43	15.24
MiniCPM-V (Hu et al., 2024)	3B	–	8.74	12.01	31.42	0.00	13.04	1.50	2.98	2.24
Ours										
SEEKER-TINY	1.3B	576	33.74	<u>66.99</u>	42.68	<u>24.99</u>	<u>42.10</u>	23.52	25.33	24.42
SEEKER	7B	576	<u>27.49</u>	71.33	<u>42.35</u>	37.91	44.77	31.85	34.98	33.41

Table 4: Probing Question Answering with Varying Page Context: Our SEEKER model seeks more accurate text answers within compact image tokens of image sequences compared to OCR-based approaches with the same context length.

Models	Input Type	ArxivQA				
		p=4:6	p=6:8	p=8:10	p=10:12	Avg
<i>LLM</i>						
DeepSeek-LLM	OCR Txt	35.79	35.74	36.00	29.99	34.38
SEEKER-LLM	OCR Txt	45.26	46.17	50.57	39.18	45.29
<i>MLLM</i>						
DeepSeek-VL	Seq Img	29.30	37.97	36.67	28.38	33.08
SEEKER	Seq Img+OCR Txt	35.30	41.22	40.73	33.49	37.68
SEEKER	Seq Img	44.43	50.81	58.10	39.95	48.32

more compact, with a significant peak at lower token counts, whereas the OCR-text displays a broader distribution with higher counts. This variation shows that OCR-text length can be vulnerable and uncontrollable in images rich in text, often leading to wide-ranging token counts. In contrast, image tokens maintain a consistent token length regardless of textual density. With a model context length set to 8192 tokens, image tokens are handled 100% of the time without truncation, whereas OCR-text frequently exceeds this limit, achieving only 66.25% execution success without truncation. Meanwhile, truncating OCR text compromises performance as shown in Table 4. This highlights the advantages of image tokens for predictable and efficient encoding of long multimodal contexts.

3.2. Inference Efficiency

In addition to its context length extrapolation capability, our model SEEKER solves long-context multimodal tasks

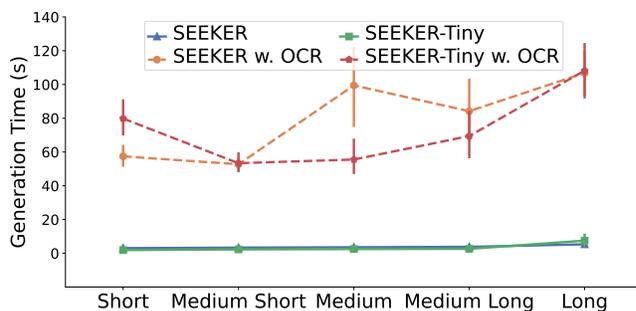


Figure 3: Generation times for SEEKER and SEEKER-TINY with and without OCR.

more efficiently compared to the OCR-based approach. For example, when comparing the inference time cost of SEEKER with and without OCR, the latter first extracts long text from multiple images and then feeds text into SEEKER. By eliminating the time-consuming OCR step, our model achieves a significant reduction in inference time. Specifically, in the longest context scenario, SEEKER is approximately three times faster than OCR-based approach, showcasing the substantial time efficiency.

4. Conclusion

Our SEEKER advances the field of long-context comprehension in multimodal large language models. By enhancing the processing of lengthy texts presented in visual formats and continual instruction-tuning on extended context tasks, SEEKER surpasses existing multimodal large language models in handling extensive multimodal contexts.

References

- Bai, J., Bai, S., Yang, S., Wang, S., Tan, S., Wang, P., Lin, J., Zhou, C., and Zhou, J. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*, 2023a.
- Bai, J., Bai, S., Yang, S., Wang, S., Tan, S., Wang, P., Lin, J., Zhou, C., and Zhou, J. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond, 2023b.
- Blecher, L., Cucurull, G., Scialom, T., and Stojnic, R. Nougat: Neural optical understanding for academic documents, 2023.
- Dai, W., Li, J., Li, D., Tiong, A. M. H., Zhao, J., Wang, W., Li, B., Fung, P., and Hoi, S. Instructblip: Towards general-purpose vision-language models with instruction tuning. *arXiv preprint arXiv:2305.06500*, 2023.
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., and Houlsby, N. An image is worth 16x16 words: Transformers for image recognition at scale, 2021.
- Gao, T., Wang, Z., Bhaskar, A., and Chen, D. Improving language understanding from screenshots, 2024.
- Hu, S., Tu, Y., Han, X., He, C., Cui, G., Long, X., Zheng, Z., Fang, Y., Huang, Y., Zhao, W., Zhang, X., Thai, Z. L., Zhang, K., Wang, C., Yao, Y., Zhao, C., Zhou, J., Cai, J., Zhai, Z., Ding, N., Jia, C., Zeng, G., Li, D., Liu, Z., and Sun, M. Minicpm: Unveiling the potential of small language models with scalable training strategies, 2024.
- Jiang, D., He, X., Zeng, H., Wei, C., Ku, M., Liu, Q., and Chen, W. Mantis: Interleaved multi-image instruction tuning, 2024.
- Jin, H., Han, X., Yang, J., Jiang, Z., Liu, Z., Chang, C.-Y., Chen, H., and Hu, X. Llm maybe longlm: Self-extend llm context window without tuning, 2024.
- Laurençon, H., Tronchon, L., Cord, M., and Sanh, V. What matters when building vision-language models?, 2024.
- Li, Z., Yang, B., Liu, Q., Ma, Z., Zhang, S., Yang, J., Sun, Y., Liu, Y., and Bai, X. Monkey: Image resolution and text label are important things for large multi-modal models. *arXiv preprint arXiv:2311.06607*, 2023.
- Lin, C.-Y. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pp. 74–81, Barcelona, Spain, July 2004. Association for Computational Linguistics. URL <https://aclanthology.org/W04-1013>.
- Liu, H., Li, C., Li, Y., and Lee, Y. J. Improved baselines with visual instruction tuning. *arXiv preprint arXiv:2310.03744*, 2023a.
- Liu, H., Li, C., Li, Y., and Lee, Y. J. Improved baselines with visual instruction tuning, 2023b.
- Liu, H., Li, C., Wu, Q., and Lee, Y. J. Visual instruction tuning. In *NeurIPS*, 2023c.
- Liu, H., Li, C., Wu, Q., and Lee, Y. J. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*, 2023d.
- Liu, H., Li, C., Li, Y., Li, B., Zhang, Y., Shen, S., and Lee, Y. J. Llava-next: Improved reasoning, ocr, and world knowledge, January 2024. URL <https://llava-vl.github.io/blog/2024-01-30-llava-next/>.
- Liu, N. F., Lin, K., Hewitt, J., Paranjape, A., Bevilacqua, M., Petroni, F., and Liang, P. Lost in the middle: How language models use long contexts, 2023e.
- Lu, H., Liu, W., Zhang, B., Wang, B., Dong, K., Liu, B., Sun, J., Ren, T., Li, Z., Yang, H., Sun, Y., Deng, C., Xu, H., Xie, Z., and Ruan, C. Deepseek-vl: Towards real-world vision-language understanding, 2024.
- Lu, Y., Li, X., Wang, W. Y., and Choi, Y. Vim: Probing multimodal large language models for visual embedded instruction following, 2023.
- McKinzie, B., Gan, Z., Fauconnier, J.-P., Dodge, S., Zhang, B., Dufter, P., Shah, D., Du, X., Peng, F., Weers, F., Belyi, A., Zhang, H., Singh, K., Kang, D., Jain, A., Hè, H., Schwarzer, M., Gunter, T., Kong, X., Zhang, A., Wang, J., Wang, C., Du, N., Lei, T., Wiseman, S., Yin, G., Lee, M., Wang, Z., Pang, R., Grasch, P., Toshev, A., and Yang, Y. Mml: Methods, analysis insights from multimodal llm pre-training, 2024.
- OpenAI. Gpt-4: Technical report. *arXiv preprint arXiv:2303.08774*, 2023a.
- OpenAI. Gpt-4v(ision) system card. <https://openai.com/research/gpt-4v-system-card>, 2023b.
- Rust, P., Lotz, J. F., Bugliarello, E., Salesky, E., de Lhoneux, M., and Elliott, D. Language modelling with pixels. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=FkSp8VW8RjH>.
- Team, G., Anil, R., Borgeaud, S., Wu, Y., Alayrac, J.-B., Yu, J., Soricut, R., Schalkwyk, J., Dai, A. M., Hauth, A., et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.

- Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.
- Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.
- Tworowski, S., Staniszewski, K., Pacek, M., Wu, Y., Michalewski, H., and Miłoś, P. Focused transformer: Contrastive training for context scaling, 2023.
- Wang, Y., Huang, R., Song, S., Huang, Z., and Huang, G. Not all images are worth 16x16 words: Dynamic transformers for efficient image recognition, 2021.
- Wu, T., Ma, K., Liang, J., Yang, Y., and Zhang, L. A comprehensive study of multimodal large language models for image quality assessment, 2024.
- Xiong, W., Liu, J., Molybog, I., Zhang, H., Bhargava, P., Hou, R., Martin, L., Rungta, R., Sankararaman, K. A., Oguz, B., Khabsa, M., Fang, H., Mehdad, Y., Narang, S., Malik, K., Fan, A., Bhosale, S., Edunov, S., Lewis, M., Wang, S., and Ma, H. Effective long-context scaling of foundation models, 2023.
- Zong, Z., Li, K., Song, G., Wang, Y., Qiao, Y., Leng, B., and Liu, Y. Self-slimmed vision transformer, 2022.

A. Background

Multimodal Large Language Model Recent advancements of proprietary Large Language Models, GPT-4 (OpenAI, 2023a), Gemini (Team et al., 2023), Claude, QWen (Bai et al., 2023a), and open-source ones, LLaMA (Touvron et al., 2023a;b), Mistral, have shown groundbreaking applications. Their counterparts in the visual domain are followed up, including GPT-4V (OpenAI, 2023b), Gemini-Vision (Team et al., 2023), Claude3-Opus-VL, Qwen-VL (Bai et al., 2023b), InstructBLIP (Dai et al., 2023), LLaVA (Liu et al., 2023d). Some work (Lu et al., 2023; Wu et al., 2024) reveals the deficit of these MLLMs in multiple images reasoning, and recent models (McKinzie et al., 2024; Laurençon et al., 2024; Jiang et al., 2024) improve such capabilities. Other work (Rust et al., 2023; Gao et al., 2024) explore to process both text and images within pixels via task-specific finetuning. However, the long-context capabilities of these MLLMs are underexplored. Our proposed SEEKER advances the long-context multimodal understanding of MLLMs from two aspects, long-form image inputs and long-form text outputs.

Long Context Transformer The Transformer-dominated LLMs have struggled with long context length as studied in (Liu et al., 2023e). LongLLaMA (Tworkowski et al., 2023), Self-Extend (Jin et al., 2024) have been proposed to increase the effective context length by either fine-tuning or training-free approach based on pre-trained LLMs. When it comes to MLLMs, additional long-context issues are introduced from Vision Transformers (ViTs) (Dosovitskiy et al., 2021) for image processing, and connecting with the LLMs. The concept of Dynamic Tokens (Wang et al., 2021) introduces a novel approach where the allocation of computational resources is adapted dynamically, emphasizing that not all image parts equally contribute to the recognition task. Additionally, the development of the Self-slimmed Vision Transformer (Zong et al., 2022) introduces a mechanism for model slimming during the inference phase, reducing computational overhead without significant loss in accuracy. In contrast, our proposed SEEKER utilizes image tokens as compact representations for image and text, alleviating the context length required for the same amount of semantic information in the language model backbone when processing multimodal content.

B. Implementation Details of SEEKER

B.1. Training Loss Curve

In Figure 6, we show the training loss curve of our SEEKER and SEEKER-TINY. Though both model have a quick loss drop initially, we observe a smoother and more consistent decrease of SEEKER than SEEKER-TINY. In the end, SEEKER stabilizes at a lower loss value, suggesting its potentially better generalization capabilities than SEEKER-TINY.

B.2. Evaluation Benchmarks and Metrics

We consider four long-form multi-image input tasks: 1) `Index`: the multiple-choice image indexing task, given a sequence of images and a question, the model selects the option with the index of the image that contains the answer, 2) `SentRetrie`: the sentence retrieval task, given a sequence of images of rendered text sampled from Wikipedia, the model is required to retrieve the first sentence from the first image, 3) `ArxivQA`: the question answering on arxiv documents, the model is required to answer the question according to visual image of arxiv documents. 4) `PassKey`: the passkey retrieval task slightly modified for multimodal model, given the sentence with a masked word, the model need to answer what is the masked word by reading the visually-situated text content from arxiv document. We consider two long-form text output tasks: 1) `ArxivVerb`: extract text from the image of arxiv documents verbatim, 2) `WikiVerb`: extract text from the image of rendered text from Wikipedia verbatim.

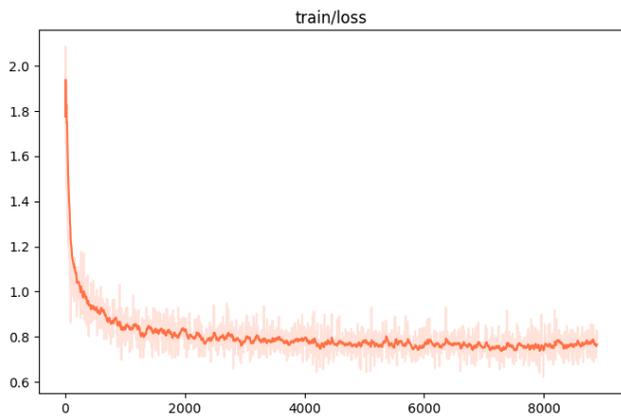


Figure 4: SEEKER

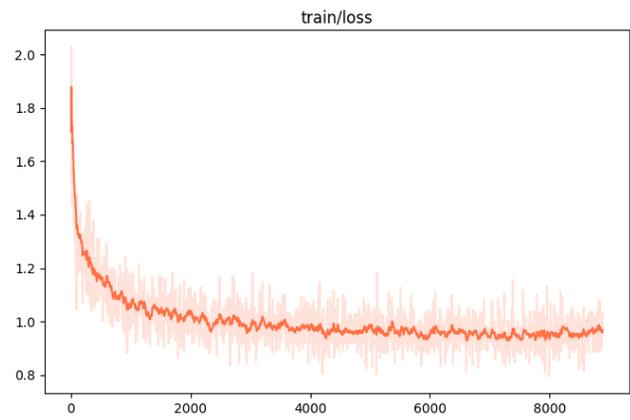


Figure 5: SEEKER-TINY

Figure 6: Training Loss Curve.

C. More Analysis

C.1. Tradeoff of Compact Context Length and High Resolution

In Figure 7, we show GPT-4-Vision with low and high resolution setting on first-sentence-retrieval. With high-resolution mode, more tokens will be used to represent the same image. Although high-resolution usually brings more details and better performance, we can see it tradeoffs capability of extrapolating long page document understanding. And thus only GPT-4-Vision low-resolution model preserves the performance in this probing task. On the right we can see that high-resolution usually take more image tokens to represent text-rich image than text tokens of OCR-extracted content, and thus even drops more quickly than feeding text.

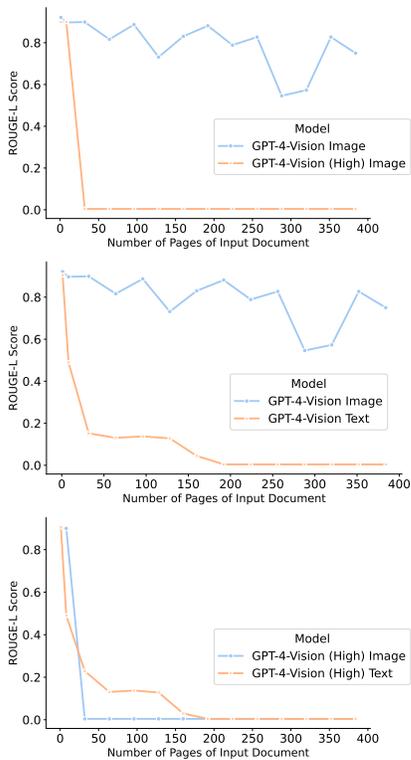


Figure 7: Performance plot on First-Sentence-Retrieval task.

D. Long-Context Multimodal Tasks

D.1. Task Examples

In Section 1.2, we first introduce multimodal long-context tasks categorized in long-form multi-image input and long-form text output. And in Figure 8-13, we visualize full task examples.



arXiv:1912.08239v1 [math.CA] 17 Dec 2019

A NOTE ON A TAUBERIAN THEOREM FOR ARITHMETIC FUNCTIONS

ALEXANDER HED PATAWOSKI

ABSTRACT. We offer new Tauberian theorems for a generalized partition function of power series with conditions on its coefficients. The well-known Hardy-Littlewood Tauberian result [H, pg.137] states that if a_n are non-negative, constant $C > 0$, and

$$\sum_{n=1}^{\infty} a_n x^n = O\left(\frac{1}{1-x}\right)$$

as $x \rightarrow 1^-$ in \mathbb{R} , then

$$A(x) = \sum_{n=0}^{\infty} a_n = O\left(\frac{1}{1-x}\right)$$

as $x \rightarrow \infty$. Here we use the same definition in the sense that $f(x) = O(g(x))$ means there exists a constant $C > 0$ such that $|f(x)| \leq C|g(x)|$. Additionally we write $f(x) \sim g(x)$ to mean that $\lim_{x \rightarrow \infty} f(x)/g(x) = 1$. The main class of asymptotic formulae we will present in this note follow directly from the following well-known result [H, (also found in [L, Lemma 1])].

Lemma 1.1. Suppose that the sequence (a_n) is real, and $\sum_{n=0}^{\infty} a_n x^n = O(x^{-\epsilon})$, $\epsilon > 0$, as $x \rightarrow \infty$. Here the limit is taken to be positive or infinite. Then we have that

$$\sum_{n=0}^{\infty} a_n x^n = O\left(\frac{1}{(1-x)^{\epsilon+1}}\right)$$

as $x \rightarrow 1^-$, where $\Gamma(\epsilon)$ is the classical Gamma function.

4

ALEXANDER HED PATAWOSKI

Theorem 2.2. Let $R(x)$ denote the k th Bernoulli polynomial. For positive ϵ and $k \geq 2$,

$$\sum_{n=0}^{\infty} p_n(x) = O\left(\frac{1}{(1-x)^k}\right) + O\left(\frac{1}{(1-x)^{k-1}}\right) + O(x^{\epsilon-k})$$

as $x \rightarrow 1^-$. We need the formula [2, pg.81, Proposition 9.2.12]

$$p_n(x) = \sum_{j=0}^n \binom{n}{j} B_j(x) \frac{(-1)^{n-j}}{(n-j)!}$$

valid for each $k \geq 2$. Summing (2) over the interval $1-\epsilon \leq x \leq \epsilon$ and applying (2.2) gives the theorem after using standard properties of O for polynomials. \square

A nice corollary to this result is

$$\sum_{n=0}^{\infty} p_n(x) = O\left(\frac{1}{(1-x)^k}\right) + O\left(\frac{1}{(1-x)^{k-1}}\right)$$

as $x \rightarrow 0^+$.

3

REMARKS RELATED TO THE VON MANGoldt FUNCTION

As we previously noted, the sum $\sum_{n=0}^{\infty} a_n A(x)$ was examined in [H]. It was shown [H, pg.139] heuristically that

$$\sum_{n=0}^{\infty} a_n A(x) \sim S A(x)$$

as $x \rightarrow \infty$, where $S = \sum_{n=0}^{\infty} p_n(x) \log(p_n(x))$, and p is a suitable multiplicative function such that S is finite when $M \rightarrow \infty$. If (B_j) holds true, and $R(x) \sim x^k$, $\epsilon > 0$, it would be the case then that we have

$$\sum_{n=0}^{\infty} a_n A(x) \sim C \frac{1}{(1-x)^{k-1}}$$

A simple special case of this would be if $a_n = 1$ for all n , which would imply

$$\sum_{n=0}^{\infty} A(x) \sim C \frac{1}{(1-x)^{k-1}}$$

This we know to be true by the Prime Number Theorem $\sum_{n \leq x} \Lambda(n) \sim x$ [H, pg.81, eq.(2.1)], and the Hardy-Littlewood result noted in the introduction [H, pg.137]. We give a general result which we believe to be of some interest.

5

ALEXANDER HED PATAWOSKI

Admittedly some of the form $\sum_{n=0}^{\infty} a_n x^n$, where k is chosen to be an arbitrary function $k: \mathbb{N} \rightarrow \mathbb{C}$, have been studied for some time. In the case of $k = \lambda(x)$, the von Mangoldt function, estimates were proved in [H]. A recent study on sums of the form $\sum_{n=0}^{\infty} a_n \lambda(x)$, where $\lambda(x)$ denotes the x -th divisor function, was produced in [6].

The purpose of this note is to offer asymptotic results for power series of the form

$$\sum_{n=0}^{\infty} a_n x^n$$

as $x \rightarrow 1^-$. Our main results are centered on partitions in the next section, and some additional results follow in the last section concerning the case $k = \lambda(x) = \sum_{d|x} \Lambda(d) \log(x/d)$, where $\Lambda(x)$ is the Möbius function [8,8].

2

THE PARTITION FUNCTION $p_k(x)$

First let us consider $p_k(x)$, the number of partitions of k into at most x parts. It is an elementary fact [H, pg.218] that $p_k(x) \leq (x+1)^k$. Since $(x+1)^k = O(x^k)$, we have the trivial estimate

$$\sum_{n=0}^{\infty} a_n p_k(x) = O(x^k A(x))$$

which suggests an interesting application of Lemma 1.1 would be of interest in understanding

$$\sum_{n=0}^{\infty} a_n p_k(x) x^n$$

However, it seems we can do better by appealing to a result found in [7].

Theorem 2.1. Let H denote the set consisting of k positive integers which have a greatest common divisor of 1, i.e., are relatively prime. If $p_k(x)$ is the number of partitions of k into parts taken from the set H , we have that

$$\sum_{n=0}^{\infty} a_n p_k(x) x^n = O\left(\frac{1}{(1-x)^k}\right) + O\left(\frac{1}{(1-x)^{k-1}}\right)$$

as $x \rightarrow 1^-$, where the implied constant depends on the set H .

Proof. First we note [H]

$$p_k(x) = O\left(\frac{1}{(1-x)^k}\right) + O(x^{k-1})$$

6

A NOTE ON A TAUBERIAN THEOREM FOR ARITHMETIC FUNCTIONS

Theorem 3.1. Suppose that $A(x) = O(x^{-\epsilon})$, $\epsilon > 0$, as $x \rightarrow \infty$. We have that

$$\sum_{n=0}^{\infty} a_n A(x) x^n = O\left(\frac{1}{(1-x)^k}\right) + O\left(\frac{1}{(1-x)^{k-1}}\right)$$

as $x \rightarrow 1^-$.

Proof. Using the fact that $A(x) \leq \log(x)^k$, [H, eq.(1.45)], we have that

$$\sum_{n=0}^{\infty} a_n A(x) \leq \log(x)^k A(x)$$

This together with the growth assumption on $A(x)$ and Lemma 1.1 gives the theorem. \square

A nice consequence of applying the Prime Number Theorem $\sum_{n \leq x} \Lambda(n) \sim x$, to Theorem 3.1 gives us the formula

$$\sum_{n=0}^{\infty} \Lambda(n) A(x) x^n = O\left(\frac{1}{(1-x)^k}\right) + O\left(\frac{1}{(1-x)^{k-1}}\right)$$

as $x \rightarrow 1^-$. Setting $k = 1$ in (3.1) gives us the degree

$$\sum_{n=0}^{\infty} \Lambda(n) x^n = O\left(\frac{1}{(1-x)^k}\right) + O\left(\frac{1}{(1-x)^{k-1}}\right)$$

as $x \rightarrow 1^-$.

REFERENCES

[1] G. H. Hardy, *Number Theory*, W. B. Saunders, Philadelphia, 1975 (Reprinted Dover, New York 1965).

[2] H. Cohen, *Number Theory of U -adic and M -adic Fields*, Graduate Texts in Math. 249, Springer-Verlag 1987.

[3] J. Erdős and H. Halász, *Asymptotic laws for primes*, Ann. of Math. (2), 1983:1049-1065, 1984.

[4] S. Gethner, *Asymptotic estimates for some number-theoretic power series*, Acta Arith. 142 (2004), pp. 147-168.

[5] G. H. Hardy and E. Littlewood, *Partition numbers concerning power series and Dirichlet's series whose coefficients are positive*, Proc. London Math. Soc., 15, (1916) pp. 178-191.

[6] H. Isai and E. Kowalski, *Arithmetic number theory*, American Mathematical Society Colloquium Publications, vol. 59, American Mathematical Society, Providence, RI, 2004.

[7] H. Yoshitane, *Partitions with parts in a finite set*, Proc. Amer. Math. Soc. 130 (1980) 1380-1379.

Question: In this task, please reply with the option letter of which Image contains the given Sentence. Sentence: 'Next we consider a direct corollary of this result by applying to prime number theorem' Instruction: Which Image contains the above Sentence? Select from these options: (A) Image 1 (B) Image 2 (C) Image 3 (D) Image 4 (E) Image 5 (F) Image 6.

Answer: (C) Image 3

Figure 8: Task Index.

12



Quantum-Classical Transitions in Complex Networks 11

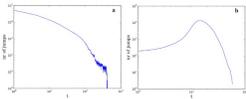


Figure 5. The number of particles that change their energy level (indicated as number of jumps) with time (considering a network with $n = 3000$). **a** During the cooling process, the number of jumps rapidly decreases. **b** During the heating process, at the beginning, all particles are constrained to low-energy levels. After a few time steps particles can find more available states in the upper bundles and the number of jumps increases. This curve reaches its maximum when all particles have available upper energy levels to reach, until these top levels become full and the number of jumps begins to decrease. At the end, all the particles are mainly arranged in the higher energy levels.

free, but as the temperature increases the network loses its metric structure and its hierarchical heterogeneous organization, becoming a classical random network.

Considering that many real complex networks are scale-free while others are not (see for example [16]), we deem that the proposed fermionic model can be considered a good candidate for representing their evolution, at low and high temperatures. As shown in Figure 5, we analyzed also the dynamics of particles during both processes. In each simulation we observed that the cooling process takes more time to allow the particles to get to their final positions. During the cooling process, the number of particles changing position is very high from the first time step. In contrast, during the heating process, we found that, at the beginning, this number is small and rapidly increases after a few (usually about ten) time steps. Then, this number of jumps reaches a maximum and begins to decrease until all particles stop moving. We deem that this behavior is an effect of the Eq. (16), since it has been defined to prevent particles from occupying high-energy levels densely at high temperatures.

6. Conclusions

In this paper, we have defined a fermionic network model that allows us to represent complex networks as quantum gases. Using this model, we have shown that network evolution is a temperature-dependent process characterized by three main phases: classical random, scale-free and *WTA*. On performing a cooling process, a transition from a classical random to a scale-free network takes place. Notably, the system achieves equilibrium when a *WTA* structure is reached, despite the non-equilibrium nature of the network evolution. On the other hand, on performing a heating process which starts from a *WTA* structure, the network evolution follows a slightly different path. In

Question: Read the text in the image verbatim.



Answer: free, but as the temperature increases the network loses its metric structure and its hierarchical heterogeneous organization, becoming a classical random network. Considering that many real complex networks are scale-free while others are not (see for example [16]), we deem that the proposed fermionic model can be considered a good candidate for representing their evolution, at low and high temperatures. As shown in Figure 5, we analyzed also the dynamics of particles during both processes. In each simulation we observed that the cooling process takes more time to allow the particles to get to their final positions. During the cooling process, the number of particles changing position is very high from the first time step. In contrast, during the heating process, we found that...

Figure 11: Task ArxivVerbatim.



Beginning of the sequence: Aristotle (; Aristotélēs, ; 384–322 BC) was a Greek philosopher and polymath during the Classical period in Ancient Greece. Taught by Plato, he was the founder of the Lyceum, the Peripatetic school of philosophy, and the Aristotelian tradition. His writings cover many subjects including physics, biology, zoology, metaphysics, logic, ethics, aesthetics, poetry, theatre, music, rhetoric, psychology, linguistics, economics, politics, meteorology, geology and government. Aristotle provided a complex synthesis of the various philosophies existing prior to him. It was above all from his teachings that the West inherited its intellectual lexicon, as well as problems and methods of inquiry. As a result, his philosophy has exerted a unique influence on almost every form of knowledge in the West and it continues to be a subject of contemporary philosophical discussion. Little is known about his life. Aristotle was born in the city of Stagira in Northern Greece. His father, Nicomachus, died when Aristotle was a child, and he was brought up by a guardian. At seventeen or eighteen years of age he joined Plato's Academy in Athens and remained there until the age of thirty-seven (c. 347 BC). Shortly after Plato died, Aristotle left Athens and, at the request of Philip II of Macedon, tutored Alexander the Great beginning in 343 BC. He established a library in the Lyceum which helped him to produce many of his hundreds of books on papyrus scrolls. Though Aristotle wrote many elegant treatises and dialogues for publication, only around a third of his original output has survived, none of it intended for publication. Aristotle's views profoundly shaped medieval scholarship. The influence of physical science extended from Late Antiquity and the Early Middle Ages into the Renaissance, and were not replaced systematically until the Enlightenment and theories such as classical mechanics were developed. Some of Aristotle's zoological observations found in his biology, such as on the hectocotyl (reproductive) arm of the octopus, were disbelieved until the 19th century. He also influenced Judeo-Islamic philosophies (800–1400) during the Middle Ages, as well as Christian theology, especially the Neoplatonism of the Early Church and the scholastic tradition of the Catholic Church. Aristotle was revered among medieval Muslim scholars as "The First Teacher", and among medieval Christians like Thomas Aquinas as simply "The Philosopher", while the poet Dante called him "the master of those who know". His works contain the earliest known formal study of logic, and were studied by medieval scholars such as Peter Abelard and John Buridan. Aristotle's influence on logic continued well into the 19th century. In addition, his ethics, though always influential, gained renewed interest with the modern advent of virtue ethics. Aristotle has been called "the father of logic", "the father of biology", "the father of political science", "the father of zoology", "the father of embryology", "the father of natural law", "the father of scientific method", "the father of rhetoric", "the father of psychology", "the father of realism", "the father of criticism", "the father of individualism", "the father of teleology", and "the father of meteorology". Lifeln general, the details of Aristotle's life are not well-established. The biographies written in ancient times are often speculative and historians only agree on a few salient points. Aristotle, whose name means "the best purpose" in Ancient Greek, was born in 384 BC in Stagira, Chalcidice, about 55 km (34 miles) east of modern-day Thessaloniki. His father, Nicomachus, was the personal physician to King Amyntas of Macedon. While he was young, Aristotle learned about biology and medical information, which was taught by his father. Both of Aristotle's parents died when he was about thirteen, and Proxenus of Atarneus became his guardian. Although little information about Aristotle's childhood

Question: Read the text in the image verbatim.



Answer: Aristotle (; Aristotélēs, ; 384–322 BC) was a Greek philosopher and polymath during the Classical period in Ancient Greece. Taught by Plato, he was the founder of the Lyceum, the Peripatetic school of philosophy, and the Aristotelian tradition. His writings cover many subjects including physics, biology, zoology, metaphysics, logic, ethics, aesthetics, poetry, theatre, music, rhetoric, psychology, linguistics, economics, politics, meteorology, geology and government. Aristotle provided a complex synthesis of the various philosophies existing prior to him.



Beginning of the sequence: An American in Paris is a jazz-influenced orchestral piece by American composer George Gershwin first performed in 1928. It was inspired by the time that Gershwin had spent in Paris and evokes the sights and energy of the French capital during the Années folles. Gershwin scored the piece for the standard instruments of the symphony orchestra plus celesta, saxophones, and automobile horns. He brought back four Parisian taxi horns for the New York premiere of the composition, which took place on December 13, 1928, in Carnegie Hall, with Walter Damrosch conducting the New York Philharmonic. It was Damrosch who had commissioned Gershwin to write his Concerto in F following the earlier success of Rhapsody in Blue (1924). He completed the orchestration on November 18, less than four weeks before the work's premiere. He collaborated on the original program notes with critic and composer Deems Taylor. Background Although the story is likely apocryphal, Gershwin is said to have been attracted by Maurice Ravel's unusual chords, and Gershwin went on his first trip to Paris in 1926 ready to study with Ravel. After his initial student audition with Ravel turned into a sharing of musical theories, Ravel said he could not teach him, saying, "Why be a second-rate Ravel when you can be a first-rate Gershwin?" Gershwin strongly encouraged Ravel to come to the United States for a tour. To this end, upon his return to New York, Gershwin joined the efforts of Ravel's friend Robert Schmitz, a pianist Ravel had met during the war, to urge Ravel to tour the U.S. Schmitz was the head of Pro Musica, promoting Franco-American musical relations, and was able to offer Ravel a \$10,000 fee for the tour, an enticement Gershwin knew would be important to Ravel. Gershwin greeted Ravel in New York in March 1928 during a party held for Ravel's birthday by Éva Gauthier. Ravel's tour reignited Gershwin's desire to return to Paris, which he and his brother Ira did after meeting Ravel. Ravel's high praise of Gershwin in an introductory letter to Nadia Boulanger caused Gershwin to seriously consider taking much more time to study abroad in Paris. Yet after he played for her, she told him she could not teach him. Boulanger gave Gershwin basically the same advice she gave all her accomplished master students: "What could I give you that you haven't already got?" This did not set Gershwin back, as his real intent abroad was to complete a new work based on Paris and perhaps a second rhapsody for piano and orchestra to follow his Rhapsody in Blue. Paris at this time hosted many expatriate writers, among them Ezra Pound, W. B. Yeats, Ernest Hemingway, and artist Pablo Picasso. Composition Gershwin based An American in Paris on a melodic fragment called "Very Parisienne", written in 1926 on his first visit to Paris as a gift to his hosts, Robert and Mabel Schirmer. Gershwin called it "a rhapsodic ballet"; it is written freely and in a much more modern idiom than his prior works. Gershwin explained in Musical America, "My purpose here is to portray the impressions of an American visitor in Paris as he strolls about the city, listens to the various street noises, and absorbs the French atmosphere." The piece is structured into five sections, which culminate in a loose ABA format. Gershwin's first A episode introduces the two main "walking" themes in the "Allegretto grazioso" and develops a third theme in the "Subito con brio". The style of this A section is written in the typical French style of composers Claude Debussy and Les Six. This A section featured duple meter, singsong rhythms, and diatonic melodies

Question: What is the first sentence in the image?



Answer: An American in Paris is a jazz-influenced orchestral piece by American composer George Gershwin first performed in 1928.

Figure 13: Task SentRetrie.

E. Discussion

E.1. Limitations

While our model, SEEKER, has made significant strides in processing extended-context multimodal inputs, it encounters several critical limitations that require deeper investigation. The process of compressing textual information into visual tokens, although efficient, may inadvertently overlook precise textual understanding. Future endeavors should focus on developing hybrid encoding strategies that balance token compression with the preservation of essential information. Additionally, SEEKER could inadvertently learn and perpetuate biases present in its training data. It is imperative that further research is conducted to identify, understand, and address these biases, ensuring the model's equity and inclusiveness.

E.2. Societal Impact

By integrating visual tokens with textual data, SEEKER addresses the limitations of traditional models and supports the handling of longer input sequences. This innovation could transform various sectors, improving information accessibility and retrieval systems across academic research, legal document analysis, and extensive data processing tasks. Particularly beneficial in educational and professional environments, SEEKER enables rapid and accurate extraction of vast informational content, fostering better decision-making and knowledge dissemination. However, this advancement might exacerbate information disparities if not equitably accessible. Steps should be taken to make sure it is both affordable and available to everyone.