

# Myopia Rectification: KV Cache Pruning for MLLMs via Dynamic Attention Subsidy and Token Reclamation

## Supplementary Material

### 1. Details of datasets

In Table A3, we provide a comprehensive taxonomy of Milebench, detailing its task composition along with the corresponding number of samples and metrics for each task. The dataset encompasses 6,440 multimodal multi-image samples derived from 21 existing or self-constructed datasets, averaging 15.2 images and 422.3 words per sample. It is divided into two primary subsets: **Realistic Evaluation** and **Diagnostic Evaluation**.

**Realistic Evaluation** component challenges MLLMs to perform tasks within extensive multimodal contexts, highlighting the models’ capabilities in understanding and reasoning across prolonged interactions.

**Diagnostic Evaluation** requires MLLMs to extract relevant information from the provided context, emphasizing their proficiency in long-range information retrieval and the ability to filter out distractors. Figure A1 shows an example in the dataset.

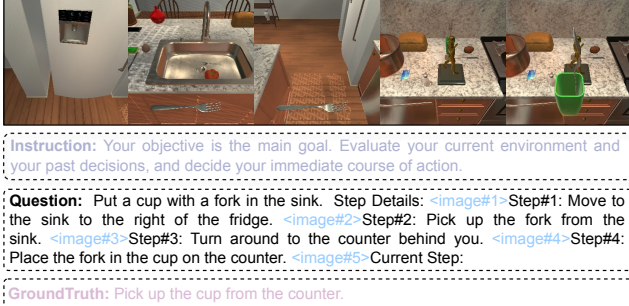


Figure A1. An illustration of a complicated multi-image reasoning sample from MileBench. Such tasks impose requirements on the model’s long contextual understanding and comprehensive spatial awareness abilities.

### 2. Analysis on visual attention

Figure A2 illustrates the attention scores of visual tokens at different positions within the sequence, where a higher value indicates greater importance and a higher retention probability. It can be observed that later images in the sequence are assigned relatively low scores across various layers. These images may be improperly discarded during the pruning process. We consider this an unfair phenomenon and a fundamental cause of attention myopia.

Table A1. Comparison between head-wise and head-aggregate.

Method	T1	T4	S3	NH	IR
Head-wise	<b>39.9</b>	39.8	13.2	5.3	<b>4.9</b>
Head-aggregate	39.7	<b>39.9</b>	13.2	<b>5.5</b>	4.8

Table A2. Image caption task with LLaVA-Next-7B model under 40% cache budget.

Method	COCO	Nocaps	Flickr30K
FastV	93.62	82.15	61.88
Ours ( $\alpha=0.5$ )	<b>95.33</b>	<b>84.28</b>	<b>63.49</b>
Ours ( $\alpha=0.2$ )	94.11	82.80	63.22
Ours ( $\alpha=0.8$ )	94.69	83.42	63.01

### 3. Head-wise information

We try to assign different  $\gamma$  and  $\alpha$  for each attention head, but it does not lead to significant accuracy improvements (as in Table. A1). This indicates that the average head mechanism is sufficient to represent most of the information in attention.

### 4. Generalization of $\gamma$ and $\alpha$ .

For  $\gamma$ , we observe that attention scores decay with position, resembling an exponential change. The parameter  $\gamma$ , serving as compensation (1.0 to 2.5), allows the derivative of this exponential function to smooth out within a specific range. For  $\alpha$ , it is used to adjust the contribution of the cross-attention scores between text and image. We believe that a centered value (0.4 to 0.6) can maximize the gains from text guidance. Table A2 demonstrates the validity of  $\alpha$  in single image tasks.

Table A3. Details of MileBench.

Category	Task	Dataset	Count	Metric
<i>Realistic Evaluation</i>				
<b>Temporal Multi-image</b>	<b>Action Understanding and Prediction (T-1)</b>	Action Localization	200	Accuracy
		Action Prediction	200	Accuracy
		Action Sequence	200	Accuracy
	<b>Object and Scene Understanding (T-2)</b>	Object Existence	200	Accuracy
Object Interaction		200	Accuracy	
Moving Attribute		200	Accuracy	
Object Shuffle		200	Accuracy	
<b>Visual Navigation and Spatial Localization (T-3)</b>	Egocentric Navigation	200	Accuracy	
	Moving Direction	200	Accuracy	
<b>Counterfactual Reasoning and State Change (T-4)</b>	Counterfactual Inference	200	Accuracy	
	State Change	200	Accuracy	
	Character Order	200	Accuracy	
	Scene Transition	200	Accuracy	
<b>Semantic Multi-image</b>	<b>Knowledge Grounded QA (S-1)</b>	Webpage QA	200	Accuracy
		Textbook QA	200	Accuracy
		Complex Multimodal QA	200	Accuracy
		Long Text with Images QA	200	Accuracy
	<b>Text-Rich Images QA (S-2)</b>	Slide QA	200	Accuracy
OCR QA		200	Accuracy	
Document QA		200	Accuracy	
<b>Visual Relation Inference (S-3)</b>	Visual Change Captioning	400	ROUGE-L	
	Visual Relationship Expressing	200	ROUGE-L	
<b>Dialogue (S-4)</b>	Multimodal Dialogue	200	Accuracy	
	Conversational Embodied Dialogue	200	ROUGE-L	
<b>Space Understanding (S-5)</b>	Space Understanding	200	Accuracy	
<i>Diagnostic Evaluation</i>				
<b>Needle In A Haystack</b>	<b>Text Needle (N-1)</b>	Text Needle In A Haystack	320	Accuracy
<b>Image Retrieval</b>	<b>Image Needle (N-2)</b>	Image Needle In A Haystack	320	Accuracy
<b>Image Retrieval</b>	<b>Image Retrieval (I-1)</b>	Image Retrieval	600	Accuracy

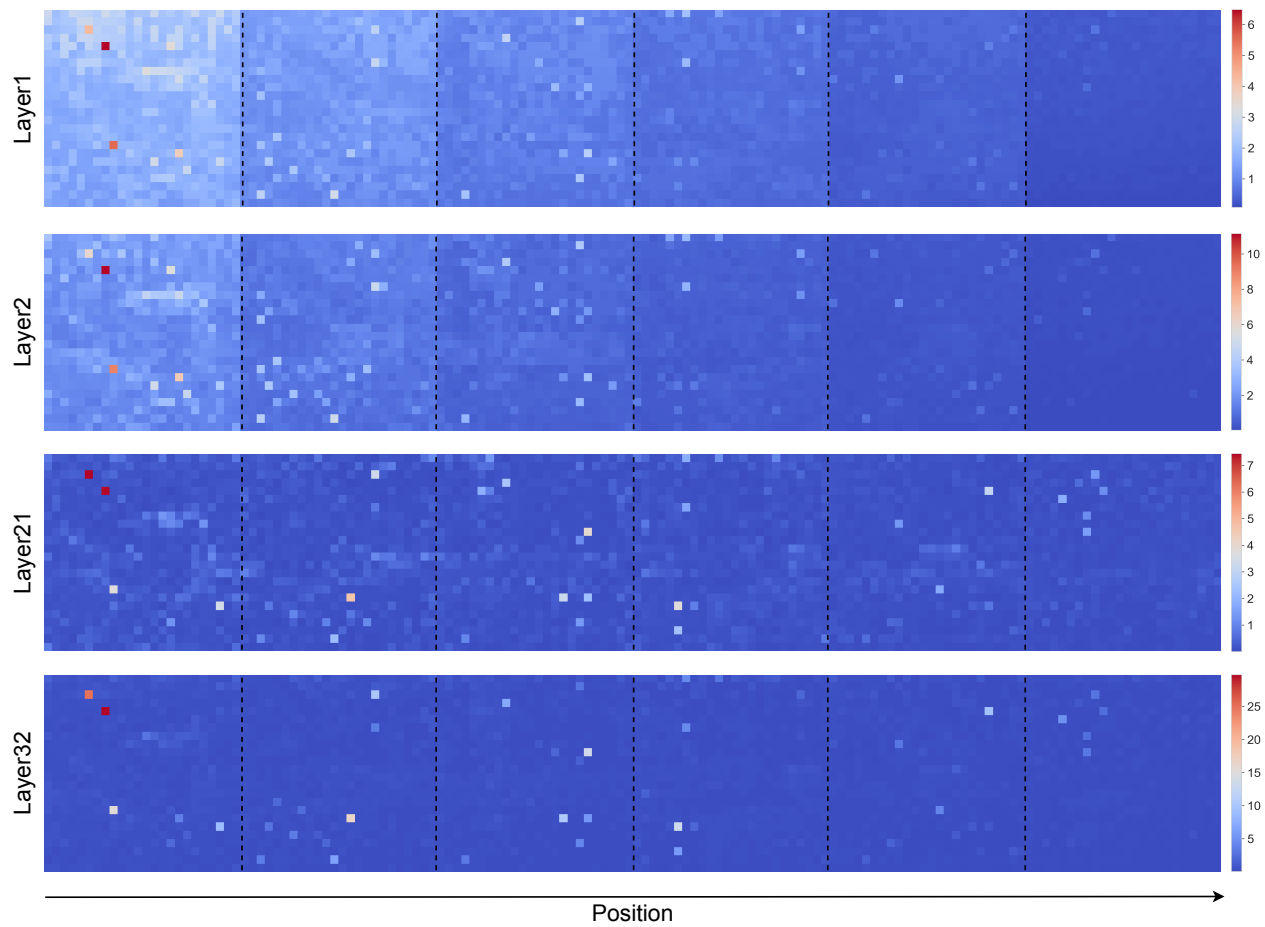


Figure A2. Visualization of attention score on each image.