

Frequency-Modulated Visual Restoration for Matryoshka Large Multimodal Models

Supplementary Material

001	This supplementary material contains several sections	
002	that provide additional details related to our work. Specifi-	
003	cally, it will cover the following topics:	
004	• In Appendix A, we provide details of the experimen-	
005	tal setup, including model architectures and evaluation	
006	benchmarks.	
007	• In Appendix B, we provide addition experiments, in-	
008	cluding FMVR for advanced open-source LMM, FMVR	
009	for larger language model, Ablation study on other Ma-	
010	tryoshka visual token sampling methods, and case results.	
011	A. Details of Experimental Setup	
012	A.1. Model Architectures	
013	LLaVA-1.5 [18]. As a representative line of open-source	
014	multimodal large language models, the LLaVA family is	
015	widely adopted for its efficiency and strong capability. It	
016	mainly has three components, where CLIP [24] is used to	
017	extract visual representations, Vicuna [9] serves as the lan-	
018	guage backbone, and a lightweight projection module maps	
019	visual features into the LLM’s embedding space. With visu-	
020	al instruction tuning, this basic setup enables the model	
021	to perform a wide range of image–text tasks. LLaVA-1.5	
022	significantly enhance multimodal reasoning ability. At the	
023	common resolution of 336×336, the model outputs 576 vi-	
024	sual tokens for each image.	
025	LLaVA-Next [19]. LLaVA-NeXT (also referred to as	
026	LLaVA-1.6) pushes the LLaVA architecture further by in-	
027	troducing a dynamic-resolution strategy aimed at strength-	
028	ening the model’s visual understanding. Instead of relying	
029	on a fixed input size, it adjusts the aspect ratio according	
030	to the native dimensions of the image and allows the effective	
031	resolution to scale by up to four times. Importantly, the vi-	
032	sual encoder itself remains unchanged. The high-resolution	
033	input is partitioned into multiple tiles, each matching the	
034	original input size, and every tile is processed independently	
035	before their visual features are concatenated and passed	
036	to the language model. This tiling-based high-resolution	
037	pipeline substantially boosts performance on tasks requiring	
038	fine-grained perception, such as OCR, detailed reasoning,	
039	and factual recognition. For consistency in our evaluations,	
040	we standardize the resolution to 672×672—four times the	
041	original—yielding a total of 2,880 visual tokens.	
042	Qwen2.5-VL [3]. Qwen2.5-VL represents a member of	
	Qwen-VL family. Its visual backbone has been thoroughly	043
	redesigned, adopting a revised Vision Transformer that in-	044
	corporates windowed attention, SwiGLU activations, and	045
	RMSNorm, following the architectural paradigm of the	046
	Qwen2.5 language model. Beyond static image under-	047
	standing, it embraces a dynamic visual processing pipeline,	048
	supporting flexible input resolutions and adaptive frame	049
	rate handling. Thanks to these advantages, Qwen2.5-VL	050
	achieves comparable performance in detailed visual percep-	051
	tion tasks such as detection, OCR, layout and document	052
	parsing, as well as structured information extraction.	053
	A.2. Evaluation Benchmarks	054
	A.2.1. General Image Benchmarks	055
	VQA_{v2} [10]. It is the version of the VQA benchmark [2],	056
	constructed to assess a model’s integrated understanding	057
	of visual content, natural language, and commonsense rea-	058
	soning. The dataset contains 265,016 images drawn from	059
	COCO [17] as well as synthetic abstract scenes, with each	060
	image paired with an average of 5.4 questions. Every ques-	061
	tion is further provided with 10 reference answers and 3 ad-	062
	ditional plausible distractor responses. Following common	063
	practice, we adopt the test-dev split for evaluation.	064
	GQA [13]. It is a large-scale VQA dataset constructed from	065
	real-world images sourced from the Visual Genome corpus	066
	[11], aimed at evaluating the model’s capacity for composi-	067
	tional reasoning and fine-grained visual comprehension. It	068
	contains more than 22 million question–answer pairs, and	069
	every image is accompanied by a richly annotated scene	070
	graph that captures object categories, attributes, and inter-	071
	object relationships. For our experiments, we report results	072
	on the test-dev balanced split.	073
	VizWiz [11]. This benchmark targets visual question an-	074
	swering in a real-world accessibility scenarios: blind or	075
	low-vision users take photos in everyday environments and	076
	verbally pose questions about them. Each visual ques-	077
	tion is accompanied by ten answers, offering diverse per-	078
	spectives on noisy and incomplete visual information. The	079
	dataset defines two core tasks—providing answers to the vi-	080
	sual questions and detecting when a question cannot be an-	081
	swered from the image alone—thereby emphasizing real-	082
	world challenges such as low-quality imagery, occlusion,	083
	and ambiguous scene content. We evaluate our model on	084
	the official test split.	085

086 **ScienceQA** [21]. This dataset is a large-scale multimodal
087 multiple-choice QA benchmark that covers a broad spec-
088 trum of scientific disciplines. It comprises 21,208 questions
089 drawn from natural sciences, language sciences, and social
090 sciences, which are further organized into 26 topics, 127
091 fine-grained categories, and 379 skill types. The questions
092 vary in modality: 48.7% provide visual content, 48.2% in-
093 clude textual context, and 30.8% offer both forms of infor-
094 mation. For evaluation, we follow prior works and adopt the
095 test split that contains image-based questions (ScienceQA-
096 IMG).

097 **POPE** [16]. This benchmark targets the evaluation of object
098 hallucination in large vision–language models. It is con-
099 structed using images from COCO [17], enabling the mea-
100 surement of hallucinated predictions. The performance is
101 summarized using precision, recall, and F1 scores. Follow-
102 ing prior works, we report results on the test split.

103 **MME** [6]. It evaluates both the perceptual and cogni-
104 tive abilities of multimodal large language models across
105 14 subtasks. The perception tasks cover coarse- and fine-
106 grained recognition as well as OCR, ranging from basic ob-
107 ject presence, count, and attributes to identifying specific
108 scenes, landmarks, and artworks. The cognition tasks as-
109 sess commonsense reasoning, numerical calculation, trans-
110 lation, and code understanding.

111 **MMBench** [20]. This benchmark provides a broad eval-
112 uation of vision–language abilities across diverse tasks. It
113 offers a larger varied set of questions and skills than prior
114 benchmarks. MMBench also proposes a CircularEval pro-
115 tocol, which uses ChatGPT to convert open-ended outputs
116 into structured choices for more stable and consistent scor-
117 ing.

118 **MM-Vet** [32]. It examines the integration of multiple multi-
119 modal abilities. It covers six core capabilities—recognition,
120 OCR, knowledge, language generation, spatial reason-
121 ing, and mathematics—through 218 challenging examples.
122 Evaluation is conducted using a ChatGPT-based assessor,
123 which provides unified metrics for answers of different for-
124 mats.

125 A.2.2. Text-oriented benchmarks

126 **TextVQA** [26]. This benchmark examines the model’s ca-
127 pacity to recognize and reason over textual content embed-
128 ded in images. The visual data, largely drawn from Open
129 Images v3 [15], contains real-world settings such as signs,
130 advertisements, and product labels, all featuring substantial
131 scene text. As a result, the benchmark provides an assess-
132 ment of OCR integration and text-sensitive visual reason-
133 ing. We evaluate on the validation split.

134 **ChartQA** [23]. It is a large benchmark for chart-based
135 question answering, emphasizing visual understanding and
136 logical or complex reasoning. It contains 9.6K human-
137 written questions and 23.1K questions generated from
138 chart summaries. Unlike earlier template-driven datasets,

Table 1. Performance comparison of different pruning methods on Qwen2.5-VL-7B.

Methods	#Vision Tokens	TextVQA	ChartQA	AI2D	MMB ^{EN}	Avg. (%)
Qwen2.5-VL	1296	84.8	86.1	80.4	82.8	83.5
FastV	128	73.8	52.2	71.4	72.9	67.6
DivPrune	128	67.0	50.4	72.1	77.8	66.8
CDPruner	128	77.8	59.2	74.0	76.2	71.8
FMVR (Ours)	81	76.3	62.2	77.5	79.6	73.9

ChartQA requires multi-step reasoning over both chart vi-
suals and their underlying data. We use the test split for
evaluation.

AI2D [14]. AI2D is a diagram-based question answering
benchmark containing over 5,000 grade-school science dia-
grams, annotated with more than 150,000 structured labels
and syntactic parses. It also provides over 15,000 multiple-
choice questions paired with the diagrams, supporting re-
search on visual reasoning and scientific diagram under-
standing. We use the masked test split for evaluation.

A.2.3. Video benchmarks

MSVD [7]. MSVD is a video question answering bench-
mark, containing roughly 2k short video clips and about 50k
question–answer pairs. It evaluates a model’s ability to un-
derstand visual content over time, including actions, events,
and temporal relationships, as well as its capacity to ground
and answer natural language queries about the video.

MSRVTT [30]. It is a large-scale video captioning and
video understanding dataset consisting of 10k video clips
collected from real-world web videos. Each video is paired
with 20 captions, covering a broad range of everyday
events, activities, and scenes. The dataset is widely used
for video–language tasks such as video captioning, video re-
trieval, and video question answering, serving as a standard
benchmark for evaluating multimodal temporal understand-
ing.

ActivityNet [4]. It is a large-scale benchmark for human
activity understanding in untrimmed videos. It contains
around 20k videos spanning 200 activity categories, cover-
ing a wide range of daily human actions. Designed for
tasks such as action recognition, temporal action localiza-
tion, and temporal segmentation, ActivityNet provides di-
verse and long-duration videos that evaluate both visual per-
ception and temporal reasoning capabilities of multimodal
models.

Video-based Generative Performance Benchmark [22].
It contains five dimensions: correctness, detail orienta-
tion, contextual understanding, temporal understanding,
and consistency. The evaluation pipelines for both the open-
ended question-answering and the generative performance
benchmarks adhere to Video ChatGPT [22].

Table 2. Performance comparison on LLaVA-1.5-13B across 10 image-based benchmarks. ‘#VisionTokens’ is the number of vision tokens. Δ rows show the change vs. LLaVA-v1.5-13B.

Methods	#Vision Tokens	VQAv2	GQA	VisWiz	SQA ^{IMG}	VQA ^{Text}	POPE	MME	MMB ^{EN}	MMB ^{CN}	MMVet	Avg. (%)
LLaVA-1.5-13B [18]	576	80.0	63.3	53.6	72.8	61.2	86.0	1531.2	68.5	63.5	36.2	66.2
FastV [8]	128	75.3	58.3	54.6	74.2	58.6	75.5	1460.6	66.1	62.3	32.8	63.1
PyramidDrop [29]	128	78.2	61.0	53.8	73.3	60.2	83.6	1489.5	67.5	62.8	32.1	64.7
SparseVLM [35]	128	77.6	59.6	51.4	74.3	59.3	85.0	1487.9	68.4	62.6	35.2	64.8
Prumerge+ [25]	128	76.2	58.3	52.8	73.3	56.1	82.7	1445.9	66.3	61.2	33.6	63.3
TRIM [27]	128	76.4	59.4	49.7	72.4	55.0	86.8	1426.9	67.1	58.4	35.1	63.2
VisionZip [31]	128	76.8	57.9	52.3	73.8	58.9	82.7	1449.2	67.4	62.5	36.0	64.1
DART [28]	128	75.7	57.7	53.0	74.2	58.7	80.4	1395.0	65.4	62.2	34.8	63.2
DivPrune [1]	128	77.1	59.2	53.5	72.8	58.0	86.8	1457.7	66.3	60.7	34.4	64.2
CDPruner [34]	128	77.7	59.7	52.9	73.2	58.4	87.3	1478.0	67.5	61.5	36.2	64.8
VisPruner [33]	128	76.1	58.8	53.4	72.8	57.9	86.0	1452.2	66.4	61.3	36.1	64.1
MQT-LLaVA [12]	144	77.9	59.4	54.1	73.6	60.2	86.1	1471.3	66.5	61.0	31.4	64.4
M3 [5]	144	79.2	60.6	53.5	72.9	60.4	87.7	1478.3	67.2	59.8	32.3	64.8
Ours												
	1	72.2	57.6	50.3	69.4	53.9	82.5	1384.7	62.2	55.8	27.8	60.1
	9	76.3	60.2	51.8	70.9	58.4	84.7	1470.5	65.1	60.4	30.7	63.2
FMVR-LLaVA	36	78.4	61.2	53.4	71.5	60.5	86.6	1477.8	66.3	62.5	32.7	64.7
	144	79.7	63.6	56.7	72.0	61.3	87.2	1489.4	68.3	63.8	34.2	66.1
	576	80.7	64.2	57.2	72.4	60.7	88.1	1528.9	69.5	65.1	35.4	67.0
Δ vs. LLaVA-v1.5-13B	144	-0.3	+0.3	+3.1	-0.8	+0.1	+1.2	-41.8	-0.2	+0.3	-2.0	-0.1

Table 3. Ablation study on Matryoshka visual token sampling, including average pooling, sequential sampling, spatial sampling, and max pooling.

Num of Vis Tokens	TextVQA				MMBench			
	Avg Pooling	Sequential	Spatial	Max Pooling	Avg Pooling	Sequential	Spatial	Max Pooling
1	49.2	46.3	45.9	47.1	60.7	57.9	58.6	59.8
9	50.8	46.7	46.4	48.2	64.2	60.4	61.5	63.9
36	55.3	51.2	49.0	52.8	65.2	61.8	63.1	64.5
144	55.5	53.7	50.6	53.3	65.8	64.2	63.7	64.9
576	57.8	55.6	52.6	54.7	65.9	64.7	64.0	65.3

180 B. Additional Experimental Results

B.2. FMVR for larger language model

194

181 B.1. FMVR for advanced open-source LMM

182 In addition to LLaVA, we further apply FMVR to one of
 183 the most advanced open-source LMM, i.e., Qwen2.5-VL.
 184 The input image is resized to 1008×1008 and the visual en-
 185 coder produces a total of 1,296 tokens. We directly insert
 186 the trained FMVR into the back of the Qwen2.5-VL’s visual
 187 encoder. We compare our FMVR with FastV DivPrune,
 188 and CDPruner. As shown in Table 1, FMVR consistently
 189 achieves the best results across on four benchmarks. FMVR
 190 with 81 visual tokens outperforms the next-best method,
 191 CDPruner with 128 visual tokens, by 2.1% in average. It
 192 demonstrates the strong generalizability of FMVR on ad-
 193 vanced LLM architectures.

To further examine whether our approach scales to stronger
 language backbones, we conduct experiments on LLaVA-
 1.5-13B. As shown in Table 2, Increasing the size of LLM
 yields substantial performance gains. Across all token prun-
 ing methods with 128 visual tokens except for Sparse-
 VLM, CDPruner, and M3, FMVR with only 36 visual to-
 kens consistently achieves the highest performance 64.7%.
 It is worth noting that our FMVR with 144 visual tokens
 achieved comparable performance to LLaVA-1.5-13B with
 576 visual tokens, i.e., 66.1% vs. 66.2%, demonstrating its
 effectiveness on larger language models.

195
196
197
198
199
200
201
202
203
204
205

Table 4. Baseline results of token numbers on LLaVA-1.5-7B.

Tokens	Method	VQAv2	GQA	VisWiz	SQA ^{IMG}	VQA ^{Text}	POPE	MME	MMB ^{EN}	MMB ^{CN}	MMVet
576	baseline	78.5	62.0	50.0	66.8	58.2	85.9	1510.7	64.3	58.3	30.5
	Ours	79.2	63.0	56.5	68.9	57.8	87.5	1510.1	65.9	58.0	34.3
144	baseline	77.8	61.1	48.3	67.5	57.3	85.0	1480.4	65.1	57.1	30.2
	Ours	78.6	62.3	55.1	69.7	55.5	86.4	1473.9	65.8	57.6	33.4
36	baseline	76.1	58.4	46.8	67.0	56.7	85.2	1448.3	64.5	57.7	29.4
	Ours	76.5	60.9	52.9	69.5	55.3	85.9	1452.5	65.2	58.3	32.2
9	baseline	73.4	56.3	46.0	67.3	52.2	83.6	1422.9	63.7	56.8	28.5
	Ours	74.5	59.1	50.7	69.9	50.8	84.1	1415.0	64.2	57.5	29.0
1	baseline	67.4	54.2	45.5	66.4	49.6	79.5	1291.0	58.2	53.9	26.2
	Ours	68.3	55.2	49.7	68.6	49.2	81.1	1284.8	60.7	53.4	26.4

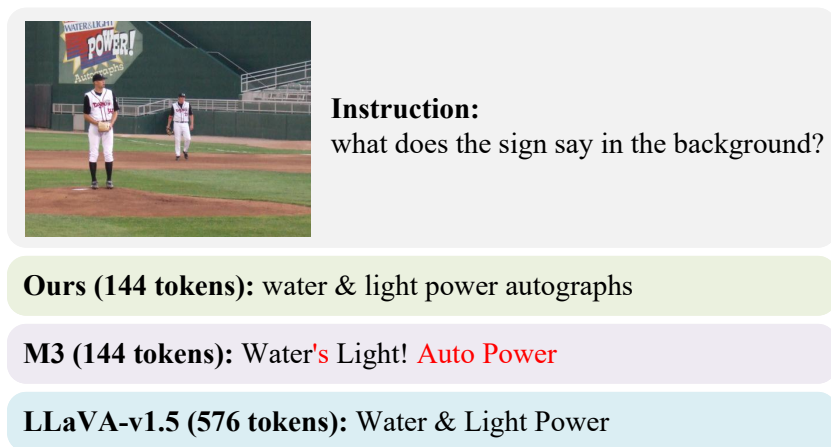


Figure 1. Example demonstrating FMVR’s image understanding capability on more challenging OCR task. Response marked in red indicates errors.

206 B.3. Ablation study on other Matryoshka visual to- 207 ken sampling methods

208 We evaluate four strategies for selecting visual tokens in
209 Matryoshka visual token construction, including average
210 pooling, spatial sampling, sequential sampling, and max
211 pooling. As illustrated in Table 3, average pooling con-
212 sistentlly outperforms the other three approaches because it
213 preserves the complete local semantic information by ag-
214 gregating all tokens within each spatial region. In contrast,
215 spatial sampling and sequential sampling select only one to-
216 ken from each region or sequence, which discards important
217 information and disrupts spatial consistency. Max pooling
218 keeps only the strongest activation, ignoring other seman-
219 tic information. As a result, average pooling provides the
220 most informative representation, leading to superior perfor-
221 mance.

222 B.4. Additional case results

223 Fig. 1 presents an example that highlights FMVR’s strength
224 on a challenging OCR task. The text in the image varies in

size and is somewhat obscured, which typically degrades 225
recognition performance in LLMs. Despite this difficulty, 226
our FMVR correctly read all the text using only 144 visu- 227
al tokens. In contrast, M3 misidentify key words, such 228
as mistakenly identifying ‘autographs’ as ‘auto power’ and 229
introduce an incorrect ‘s’. Furthermore, LLaVA-v1.5 fails 230
to recognize ‘autographs’. These results demonstrate that 231
FMVR possesses more detailed visual perception capabili- 232
ties. 233

B.5. Baseline results of token numbers 234

Table 4 reports the baseline and our method under dif- 235
ferent token budgets on LLaVA-1.5-7B. As token num- 236
bers decrease from 576 to 1, performance gradually drops. 237
However, our method consistently outperforms the base- 238
line on most benchmarks, demonstrating stronger ro- 239
bustness and effectiveness under aggressive token reduc- 240
tion. 241

242

References

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

280

281

282

283

284

285

286

287

288

289

290

291

292

293

294

295

296

297

298

- [1] Saeed Ranjbar Alvar, Gursimran Singh, Mohammad Akbari, and Yong Zhang. Divprune: Diversity-based visual token pruning for large multimodal models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 9392–9401, 2025. 3
- [2] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pages 2425–2433, 2015. 1
- [3] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023. 1
- [4] Fabian Caba Heilbron, Victor Escorcia, Bernard Ghanem, and Juan Carlos Niebles. Activitynet: A large-scale video benchmark for human activity understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 961–970, 2015. 2
- [5] Mu Cai, Jianwei Yang, Jianfeng Gao, and Yong Jae Lee. Matryoshka multimodal models. *arXiv preprint arXiv:2405.17430*, 2024. 3
- [6] Fu Chaoyou, Chen Peixian, Shen Yunhang, Qin Yulei, Zhang Mengdan, Lin Xu, Yang Jinrui, Zheng Xiawu, Li Ke, Sun Xing, et al. Mme: A comprehensive evaluation benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394*, 2023. 2
- [7] David Chen and William B Dolan. Collecting highly parallel data for paraphrase evaluation. In *Proceedings of the 49th annual meeting of the association for computational linguistics: human language technologies*, pages 190–200, 2011. 2
- [8] Liang Chen, Haozhe Zhao, Tianyu Liu, Shuai Bai, Junyang Lin, Chang Zhou, and Baobao Chang. An image is worth 1/2 tokens after layer 2: Plug-and-play inference acceleration for large vision-language models. In *European Conference on Computer Vision*, pages 19–35. Springer, 2024. 3
- [9] Wei-Lin Chiang, Zhuohan Li, Ziqing Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. See <https://vicuna.lmsys.org> (accessed 14 April 2023), 2(3):6, 2023. 1
- [10] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6904–6913, 2017. 1
- [11] Danna Gurari, Qing Li, Abigale J Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo, and Jeffrey P Bigham. Vizwiz grand challenge: Answering visual questions from blind people. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3608–3617, 2018. 1
- [12] Wenbo Hu, Zi-Yi Dou, Liunian Li, Amita Kamath, Nanyun Peng, and Kai-Wei Chang. Matryoshka query transformer for large vision-language models. *Advances in Neural Information Processing Systems*, 37:50168–50188, 2024. 3
- [13] Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6700–6709, 2019. 1
- [14] Aniruddha Kembhavi, Mike Salvato, Eric Kolve, Minjoon Seo, Hannaneh Hajishirzi, and Ali Farhadi. A diagram is worth a dozen images. In *European conference on computer vision*, pages 235–251. Springer, 2016. 2
- [15] Ivan Krasin, Tom Duerig, Neil Alldrin, Vittorio Ferrari, Sami Abu-El-Haija, Alina Kuznetsova, Hassan Rom, Jasper Uijlings, Stefan Popov, Andreas Veit, et al. Openimages: A public dataset for large-scale multi-label and multi-class image classification. *Dataset available from <https://github.com/openimages>*, 2(3):18, 2017. 2
- [16] Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating object hallucination in large vision-language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 292–305, 2023. 2
- [17] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. 1, 2
- [18] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36:34892–34916, 2023. 1, 3
- [19] Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next: Improved reasoning, ocr, and world knowledge, january 2024. URL <https://llava-vl.github.io/blog/2024-01-30-llava-next>, 2024. 1
- [20] Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multi-modal model an all-around player? In *European conference on computer vision*, pages 216–233. Springer, 2024. 2
- [21] 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. *Advances in Neural Information Processing Systems*, 35:2507–2521, 2022. 2
- [22] Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Khan. Video-chatgpt: Towards detailed video understanding via large vision and language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12585–12602, 2024. 2
- [23] Ahmed Masry, Xuan Long Do, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. Chartqa: A benchmark for question answering about charts with visual and logical reasoning. In *Findings of the association for computational linguistics: ACL 2022*, pages 2263–2279, 2022. 2

- 356 [24] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya
357 Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry,
358 Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning
359 transferable visual models from natural language supervi-
360 sion. In *International conference on machine learning*, pages
361 8748–8763. PmLR, 2021. 1
- 362 [25] Yuzhang Shang, Mu Cai, Bingxin Xu, Yong Jae Lee, and Yan
363 Yan. Llava-prumerge: Adaptive token reduction for efficient
364 large multimodal models. In *Proceedings of the IEEE/CVF
365 International Conference on Computer Vision*, pages 22857–
366 22867, 2025. 3
- 367 [26] Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang,
368 Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus
369 Rohrbach. Towards vqa models that can read. In *Proceedings
370 of the IEEE/CVF conference on computer vision and pattern
371 recognition*, pages 8317–8326, 2019. 2
- 372 [27] Dingjie Song, Wenjun Wang, Shunian Chen, Xidong Wang,
373 Michael X Guan, and Benyou Wang. Less is more: A simple
374 yet effective token reduction method for efficient multi-
375 modal llms. In *Proceedings of the 31st International Confer-
376 ence on Computational Linguistics*, pages 7614–7623, 2025.
377 3
- 378 [28] Zichen Wen, Yifeng Gao, Shaobo Wang, Junyuan Zhang,
379 Qintong Zhang, Weijia Li, Conghui He, and Linfeng
380 Zhang. Stop looking for important tokens in multimodal
381 language models: Duplication matters more. *arXiv preprint
382 arXiv:2502.11494*, 2025. 3
- 383 [29] Long Xing, Qidong Huang, Xiaoyi Dong, Jiajie Lu, Pan
384 Zhang, Yuhang Zang, Yuhang Cao, Conghui He, Jiaqi Wang,
385 Feng Wu, et al. Pyramiddrop: Accelerating your large
386 vision-language models via pyramid visual redundancy re-
387 duction. *arXiv preprint arXiv:2410.17247*, 2024. 3
- 388 [30] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. Msr-vtt: A large
389 video description dataset for bridging video and language. In
390 *Proceedings of the IEEE conference on computer vision and
391 pattern recognition*, pages 5288–5296, 2016. 2
- 392 [31] Senqiao Yang, Yukang Chen, Zhuotao Tian, Chengyao
393 Wang, Jingyao Li, Bei Yu, and Jiaya Jia. Visionzip: Longer
394 is better but not necessary in vision language models. In
395 *Proceedings of the Computer Vision and Pattern Recognition
396 Conference*, pages 19792–19802, 2025. 3
- 397 [32] Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang,
398 Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang.
399 Mm-vet: evaluating large multimodal models for integrated
400 capabilities. In *Proceedings of the 41st International Con-
401 ference on Machine Learning*, pages 57730–57754, 2024. 2
- 402 [33] Qizhe Zhang, Aosong Cheng, Ming Lu, Renrui Zhang, Zhiy-
403 ong Zhuo, Jiajun Cao, Shaobo Guo, Qi She, and Shanghang
404 Zhang. Beyond text-visual attention: Exploiting visual cues
405 for effective token pruning in vlms. In *Proceedings of the
406 IEEE/CVF International Conference on Computer Vision*,
407 pages 20857–20867, 2025. 3
- 408 [34] Qizhe Zhang, Mengzhen Liu, Lichen Li, Ming Lu, Yuan
409 Zhang, Junwen Pan, Qi She, and Shanghang Zhang. Beyond
410 attention or similarity: Maximizing conditional diversity for
411 token pruning in mllms. *arXiv preprint arXiv:2506.10967*,
412 2025. 3
- [35] Yuan Zhang, Chun-Kai Fan, Junpeng Ma, Wenzhao Zheng,
Tao Huang, Kuan Cheng, Denis Gudovskiy, Tomoyuki
Okuno, Yohei Nakata, Kurt Keutzer, et al. SparseVlm: Vi-
sual token sparsification for efficient vision-language model
inference. *arXiv preprint arXiv:2410.04417*, 2024. 3