

Accelerating Multimodal Large Language Models via Dynamic Visual-Token Exit and Empirical Findings

Qiong Wu^{1,2}, Wenhao Lin^{1,2}, Yiyi Zhou^{1,2*}, Weihao Ye¹, Zhanpeng Zeng¹, Xiaoshuai Sun^{1,2}, Rongrong Ji^{1,2}

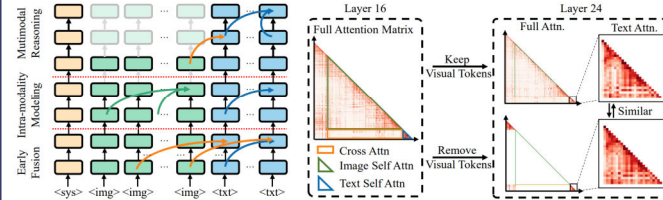
¹Key Laboratory of Multimedia Trusted Perception and Efficient Computing, Ministry of Education of China, School of Informatics, Xiamen University, China.

²Institute of Artificial Intelligence, Xiamen University, China

INTRODUCTION

Multimodal Large Language Models (MLLMs) integrate visual and textual understanding, yet their inference is computationally expensive due to redundant visual tokens. Existing efficiency methods often prune tokens heuristically, lacking in-depth exploration of the intrinsic behaviors of MLLMs. We propose Dynamic Visual Token Exiting (DyVTE), which adaptively removes redundant visual tokens early, improving efficiency with minimal accuracy loss.

MOTIVATION



Left: Illustration of three main stages observed in MLLMs.

Right: The impact of visual tokens on the text self attention at the multimodal reasoning stage.

CONTRIBUTION

1. We study the problem of visual redundancy from the perspective of MLLMs' behaviors, and reveal the dependency between text and visual tokens.
2. Based on the empirical findings, we propose a novel and effective approach to reduce visual redundancy of MLLMs, termed dynamic visual-token exit (DyVTE), which can dynamically evaluate and schedule the contributions of visual tokens to multimodal reasoning.
3. The extensive experiments on a set of MLLMs well validate the motivation and effectiveness of DyVTE, also providing insights into the principle of MLLMs

Pattern Analysis of MLLMs

Fig 2. The averaged attention scores of four MLLMs in terms of cross, visual and text attentions.

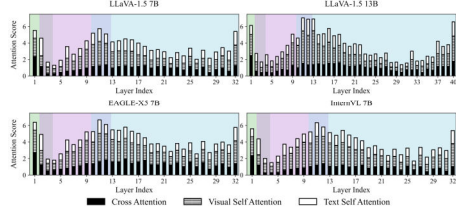


Fig 3. The relationship between manual visual-token exit and performance on LLaVA-1.5-7B and 13B.

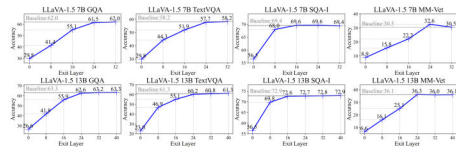
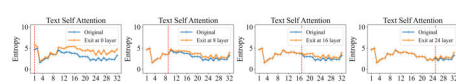


Fig 4. The entropy of text self-attention distribution with different layer to remove all visual tokens.



Dynamic Visual-token Exiting

Method

DyVTE uses lightweight hyper-networks to perceive the text token status of MLLMs and then adaptively judge the right time to remove all visual tokens.

$$p = \text{Softmax}(\text{GELU}(\text{arg}(\mathbf{U}_{i-1}^{(k)} \cdot \mathbf{T}_i^{(k)})) \mathbf{W}_1) \mathbf{W}_2)$$

When the prediction p is exit, DyVTE will remove all visual tokens after this layer, while the text ones are kept in the rest layers of MLLMs, denoted by $P_i' = G_{i+1:L}(\mathbf{T}_i^{(l)})$

Optimization

The objective of DyVTE can be defined by $\argmin_{\theta_k} d(P, P_i')$

We compare the discrete outputs of the MLLM with and without DyVTE, e.g., the answer strings A.

$$y = \begin{cases} 1, & A_i' = A, \\ 0, & A_i' \neq A. \end{cases}$$

Where A_i' denotes the answer predicted with visual-token exit at the i -th layer.

Besides, to make this supervision more robust, we also consider the prediction uncertainty as a regularization term, thereby, making the MLLM behaviors closer to the default inference:

$$y = \begin{cases} 1, & A_i' = A \wedge p_c < \tau, \\ 0, & \text{otherwise.} \end{cases}$$

Here, p_c denotes the prediction uncertainty valued by cross-entropy and multiplied by a scaling factor.

With this supervision, the hyper-networks in DyVTE can be optimized by the cross-entropy loss:

$$\mathcal{L}_D = -(\mathbf{y} \cdot \log(\mathbf{p}_1) + (1 - \mathbf{y}) \cdot \log(\mathbf{p}_0)).$$

EXPERIMENT

Tab 1. Results of MLLMs with and without DyVTE on five MLLM benchmarks.

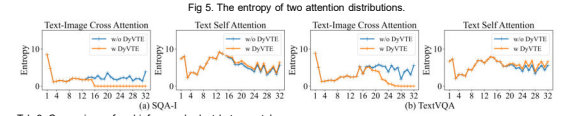
Method	SEED Accuracy	SEED FLOPs	Score	MMB Accuracy	MMB FLOPs	POPE Accuracy	POPE FLOPs	MM-Vet Accuracy	MM-Vet FLOPs
EAGLE-XS-7B	73.9 (4.0%)	43.0 (1.0%)	1581.7 (4.5%)	20.3 (2.7%)	68.8 (6.0%)	23.7 (1.9%)	88.4 (4.5%)	20.0 (2.7%)	37.8 (4.1%)
EAGLE-DyVTE-7B	61.7 (9.2%)	3.7 (0.1%)	1489.2 (8.9%)	36.3 (9.6%)	69.9 (9.8%)	9.8 (0.2%)	86.3 (7.7%)	36.7 (10.1%)	6.6 (2.4%)
VILA-7B	61.8 (4.2%)	5.3 (0.4%)	1503.1 (4.0%)	4.6 (0.3%)	69.8 (4.0%)	6.0 (0.3%)	83.6 (4.0%)	4.5 (0.4%)	36.7 (4.1%)
InternVL-7B	59.2 (16.0%)	1525.1 (15.5%)	15.5 (6.6%)	16.2 (6.6%)	86.4 (15.4%)	31.2 (15.4%)	29.5 (15.4%)	13.0 (15.0%)	8.3 (17.4%)
InternVL-DyVTE-7B	59.1 (4.2%)	11.9 (2.6%)	1474.1 (3.9%)	10.9 (2.6%)	64.4 (3.0%)	12.0 (2.6%)	81.3 (3.9%)	10.9 (2.6%)	29.5 (4.4%)
LLaVA-1.5-7B	58.6 (9.2%)	9.2 (0.5%)	1510.7 (8.9%)	8.9 (6.3%)	9.6 (8.5%)	8.8 (3.0%)	85.9 (8.8%)	30.5 (8.7%)	8.3 (7.7%)
LLaVA-DyVTE-7B	58.6 (0.0%)	5.0 (0.7%)	1491.4 (1.3%)	4.3 (0.7%)	64.7 (0.0%)	5.4 (0.6%)	81.6 (1.3%)	4.1 (0.6%)	5.7 (2.6%)
LLaVA-1.5-13B	61.6 (1.7%)	17.6 (1.5%)	1531.3 (1.6%)	67.7 (1.8%)	18.3 (5.9%)	36.1 (16.7%)	34.8 (3.8%)	10.6 (3.8%)	8.3 (7.7%)
LLaVA-DyVTE-13B	59.3 (3.7%)	7.1 (0.7%)	1540.4 (0.4%)	7.2 (0.7%)	66.0 (1.2%)	7.8 (1.4%)	84.8 (1.3%)	7.6 (1.4%)	10.6 (3.8%)

Tab 2. Results of MLLMs with and without DyVTE on four VL benchmarks.

Method	GQA Accuracy	GQA FLOPs	VQA Score	VQA FLOPs	TextVQA Accuracy	TextVQA FLOPs	SQA-1 Accuracy	SQA-1 FLOPs	Average Accuracy	Average FLOPs
EAGLE-XS-7B	64.9 (2.7%)	27.8 (1.1%)	83.4 (2.7%)	27.8 (1.1%)	71.2 (2.9%)	29.5 (1.1%)	69.8 (2.9%)	29.2 (1.1%)	72.3 (2.8%)	28.6 (1.1%)
EAGLE-DyVTE-7B	62.4 (1.9%)	21.7 (2.1%)	82.6 (1.9%)	21.6 (2.1%)	70.2 (1.4%)	24.5 (1.6%)	71.7 (2.7%)	23.3 (2.3%)	71.7 (1.9%)	22.8 (2.0%)
VILA-7B	63.1 (8.8%)	80.3 (8.8%)	62.2 (8.8%)	80.3 (8.8%)	69.5 (8.8%)	9.5 (9.5%)	69.5 (8.8%)	9.5 (9.5%)	68.8 (9.2%)	9.2 (9.2%)
VILA-DyVTE-7B	61.9 (1.9%)	5.5 (0.7%)	79.2 (1.4%)	5.4 (0.6%)	61.2 (2.2%)	7.2 (2.2%)	69.5 (0.9%)	6.1 (3.7%)	67.9 (1.3%)	6.0 (3.4%)
InternVL-7B	62.9 (15.4%)	79.3 (15.4%)	57.0 (16.1%)	79.3 (15.4%)	66.2 (16.4%)	16.4 (16.4%)	66.2 (16.4%)	16.4 (16.4%)	66.4 (15.8%)	16.4 (15.8%)
InternVL-DyVTE-7B	61.3 (1.2%)	11.8 (1.3%)	77.6 (1.1%)	11.7 (1.3%)	55.8 (1.3%)	13.5 (1.6%)	66.2 (0.9%)	12.1 (1.3%)	65.2 (1.8%)	12.3 (1.2%)
LLaVA-1.5-7B	62.0 (8.8%)	78.5 (8.8%)	58.2 (9.5%)	78.5 (8.8%)	69.1 (9.8%)	9.8 (9.8%)	67.0 (9.2%)	9.2 (9.2%)	67.0 (9.2%)	9.2 (9.2%)
LLaVA-DyVTE-7B	60.0 (1.2%)	5.3 (0.8%)	76.6 (1.2%)	5.1 (0.8%)	56.6 (2.7%)	6.7 (2.9%)	69.6 (0.6%)	5.5 (4.0%)	65.7 (1.9%)	5.6 (0.9%)
LLaVA-1.5-13B	63.3 (16.8%)	80.0 (16.8%)	61.3 (18.1%)	80.0 (16.8%)	72.3 (18.1%)	18.1 (18.1%)	72.3 (18.1%)	18.1 (18.1%)	69.4 (17.6%)	17.6 (17.6%)
LLaVA-DyVTE-13B	62.3 (0.4%)	9.0 (0.4%)	78.8 (1.1%)	8.9 (0.4%)	58.9 (1.3%)	10.8 (0.9%)	72.3 (0.8%)	8.2 (3.5%)	68.1 (1.9%)	9.2 (4.7%)

Tab 3. Ablation study of different token statuses for DyVTE.

Method	State	Mean	Last	Mean	Last	GQA Acc.	Exit Layer	TextVQA Acc.	TextVQA FLOPs	MM-Vet Acc.	Exit Layer	Average Acc.	Average FLOPs
✓	✓	✓	✓	✓	✓	61.4	21.6	57.3	21.0	34.5	21.4	69.7	21.4
✓	✓	✓	✓	✓	✓	61.2	21.3	57.1	21.1	30.0	21.0	69.7	21.4
✓	✓	✓	✓	✓	✓	60.1	21.4	57.6	22.1	31.1	21.1	69.7	20.4
✓	✓	✓	✓	✓	✓	59.8	16.3	56.3	19.1	29.9	21.0	69.8	13.8
✓	✓	✓	✓	✓	✓	61.1	21.2	57.3	21.3	31.6	21.3	69.7	21.0
✓	✓	✓	✓	✓	✓	58.7	16.5	57.6	22.2	32.7	22.5	69.7	14.1
✓	✓	✓	✓	✓	✓	59.4	16.3	56.5	19.9	30.9	21.8	69.7	14.4
✓	✓	✓	✓	✓	✓	59.4	16.7	57.4	21.8	31.3	22.0	69.6	20.3
✓	✓	✓	✓	✓	✓	59.3	16.4	56.3	19.9	31.2	21.6	69.6	14.4
✓	✓	✓	✓	✓	✓	60.0	18.4	56.6	19.7	31.9	21.1	69.6	13.4



Tab 3. Comparison of real inference budget between token pruning methods and vanilla decoding.

Method	Latency	Acc.	Latency	Acc.
LLaVA-13B	0.237s	72.9	0.236s	67.7
FastV	0.171s (-27.7%)	73.1 (+0.3%)	0.174s (-26.2%)	68.6 (+1.3%)
ToMe	0.175s (-26.2%)	73.2 (+0.4%)	0.179s (-24.2%)	67.4 (+0.4%)
DyVTE	0.161s (-32.1%)	72.3 (+0.8%)	0.163s (-30.9%)	66.0 (-2.5%)
DyVTE+FastV	0.154s (-34.9%)	73.0 (+0.6%)	0.155s (-34.4%)	68.5 (+1.2%)

Fig 6. Examples on LLaVA-1.5 with DyVTE.



Tab 5. Comparison between DyVTE and token pruning methods.

Method	SQA-1 Accuracy	SQA-1 FLOPs	MM-Vet Accuracy	MM-Vet FLOPs	SEED Accuracy	SEED FLOPs	MMB Accuracy	MMB FLOPs	Average Accuracy	Average FLOPs
LLaVA-7B	69.4	9.8	30.5	8.7	58.6	9.2	64.3	9.6	55.7	9.3
ToMe	69.6 (+0.3%)	5.9 (-39.8%)	30.6 (+0.3%)	4.9 (-43.7%)	57.8 (-1.4%)	5.5 (-40.2%)	63.7 (+0.9%)	5.7 (-40.6%)	55.4 (+0.5%)	5.5 (-40.9%)
FastV [7]	69.0 (+0.6%)	6.2 (-36.7%)	31.3 (+2.6%)	5.2 (-35.8%)	58.2 (+0.4%)	5.8 (-37.4%)	64.4 (+0.2%)	6.0 (-37.5%)	55.7 (+0.0%)	5.8 (-37.6%)
DyVTE	69.6 (+0.3%)	5.5 (-43.9%)	31.9 (+0.6%)	6.3 (-37.6%)	58.6 (0.0%)	5.0 (-43.4%)	64.7 (+0.0%)	5.2 (-43.4%)	56.2 (+0.9%)	5.5 (-40.9%)
DyVTE+FastV	68.9 (-0.7%)	4.8 (-51.0%)	29.8 (-2.3%)	4.0 (-54.0%)	58.2 (-0.4%)	4.6 (-50.0%)	64.3 (+0.2%)	4.8 (-51.2%)	55.3 (-0.7%)	4.5 (-51.5%)

Tab 6. Comparison between DyVTE and token pruning methods.

MLLM	Trained	GQA Acc.	Exit Layer	VQA-2 Acc.	Exit Layer	SEED Acc.	Exit Layer	MMB Acc.	Exit Layer	MME Acc.	Exit Layer
VILA-7B	VILA-7B	63.1	-	80.3	-	61.7	-	69.9	-	1489.2	-
InternVL-7B	InternVL-7B	61.3	17.4	79.2	16.7	61.8	15.8	69.8	16.2	1503.1	13.1
LLaVA-7B	LLaVA-7B	62.9	79.3	77.6	15.8	59.1	13.9	64.4	13.6	1474.1	12.1
LLaVA-7B	LLaVA-7B	61.3	17.0	77.7	16.8	59.2	18.1	64.9	17.8	1516.9	14.1