

# Global Compression Commander: Plug-and-Play Inference Acceleration for High-Resolution Large Vision-Language Models

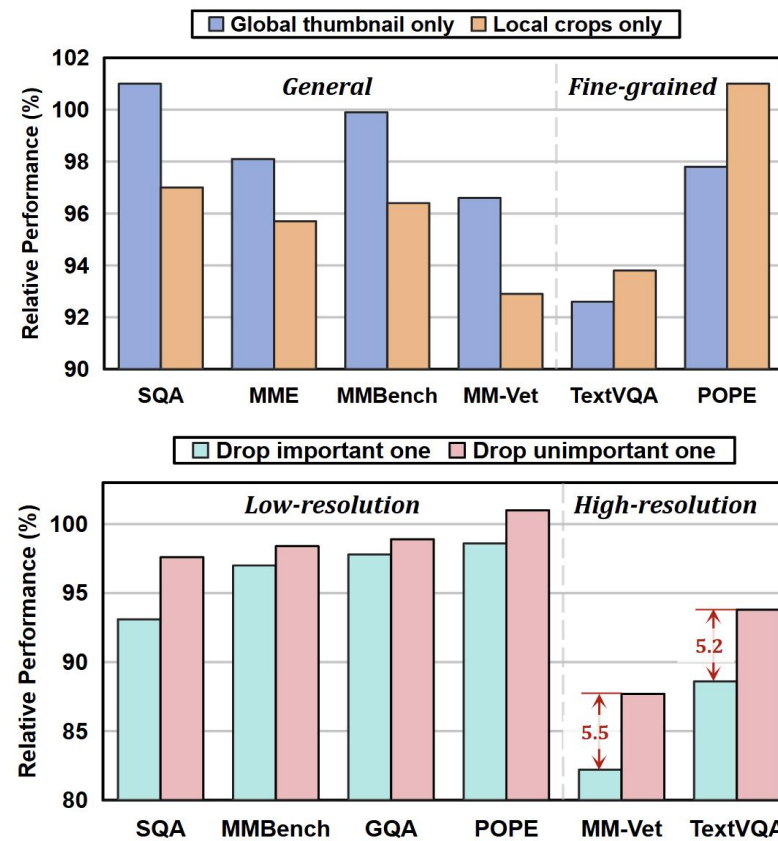
Xuyang Liu<sup>1\*</sup>, Ziming Wang<sup>2</sup>, Junjie Chen<sup>1</sup>, Yuhang Han<sup>3</sup>, Yingyao Wang<sup>2</sup>, Jiale Yuan<sup>2</sup>, Jun Song<sup>2</sup>✉, Siteng Huang<sup>4</sup>, Honggang Chen<sup>1</sup>✉

<sup>1</sup>Sichuan University, <sup>2</sup>Taobao & Tmall Group of Alibaba, <sup>3</sup>Westlake University, <sup>4</sup>Zhejiang University

\*This work was done during an internship at Alibaba. ✉ Corresponding author: honggang\_chen@scu.edu.cn



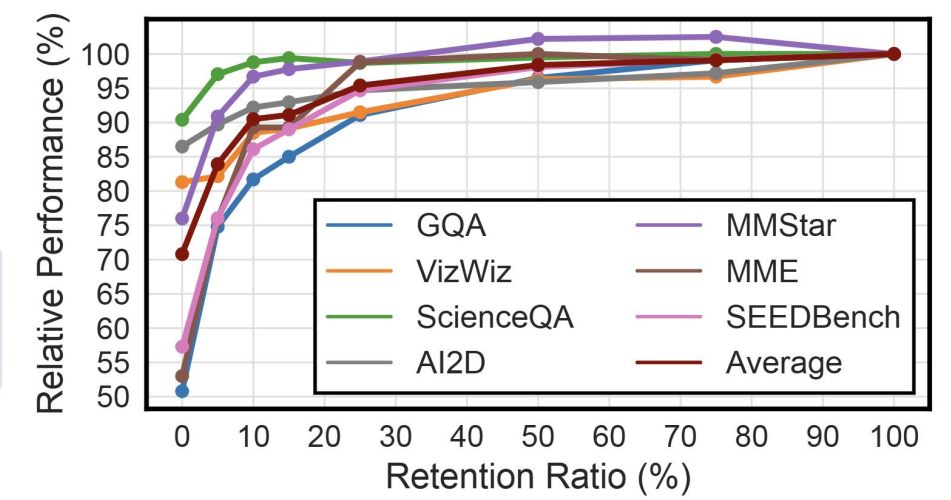
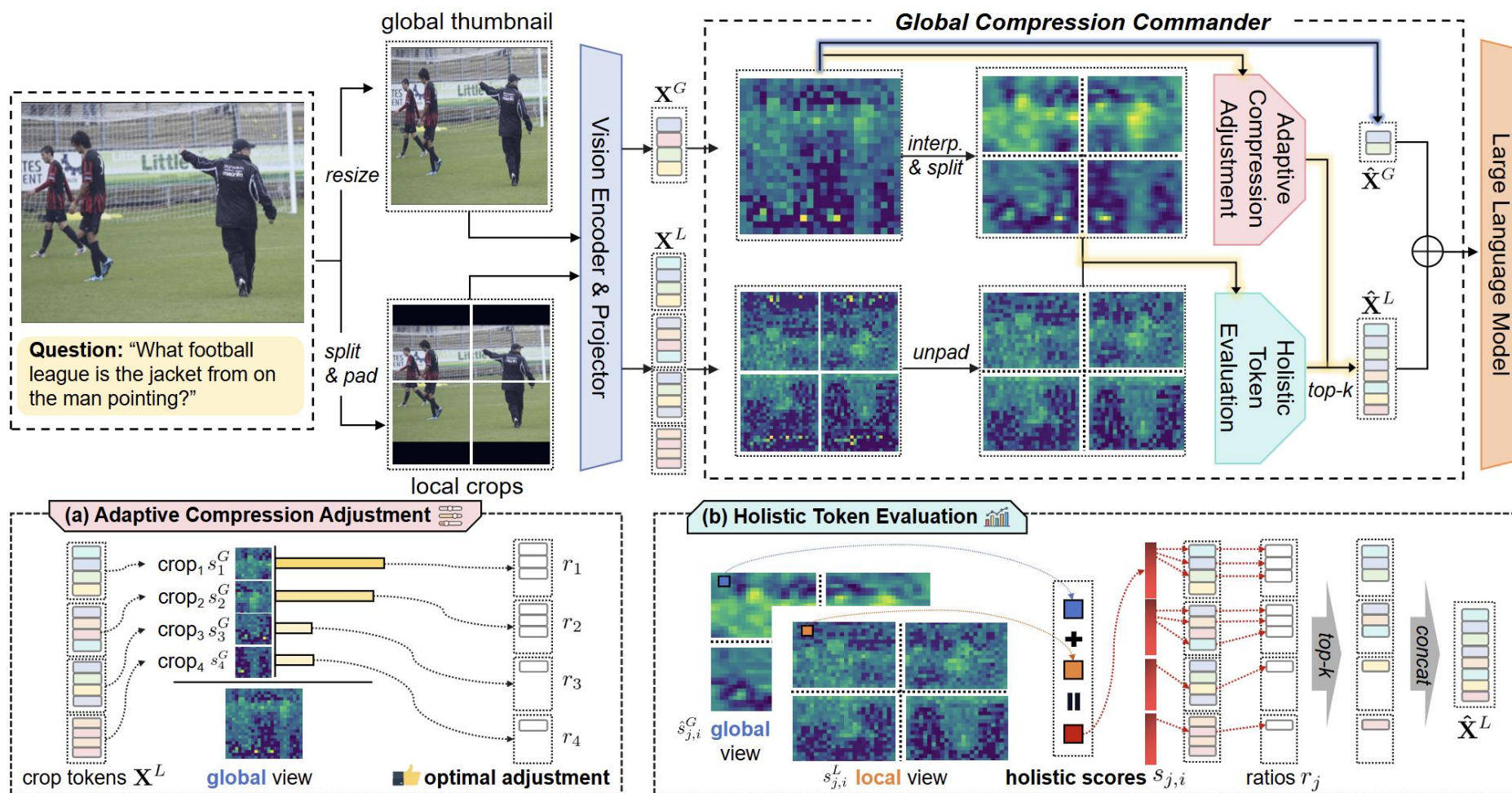
## Motivation and Research Status: The Characteristics of HR-LVLMs and the Limitations of Existing Methods



**Takeaways:** We derive *two observations* for LVLMs with dynamic cropping: **(i)** Thumbnails and crops serve complementary roles in HR-LVLMs with dynamic cropping. **(ii)** Crops exhibit varying information richness, leading to different contributions.

**Contributions:** **(i)** analyze dynamic-cropping HR-LVLMs, revealing global-context neglect, crop informativeness disparity, and content-agnostic positional bias; **(ii)** propose GlobalCom<sup>2</sup>, a training-free plug-and-play global-to-local compressor; **(iii)** retain >90% performance while pruning 90% visual tokens.

## Our Solution: Global Compression Commander (GlobalCom<sup>2</sup>) - "Global-to-Local" Guided Compression Philosophy



**Performance:** GlobalCom<sup>2</sup> maintains over 90% performance while compressing 90% visual tokens across multiple vision-language understanding benchmarks.

**Efficiency:** GlobalCom<sup>2</sup> cuts FLOPs to 9.1% and peak GPU memory to 60%, delivering a 1.8× throughput gain.

## Experimental Results: Optimal Performance-Efficiency Tradeoffs

Method	GQA	VizWiz	SQA	MMB	POPE	VQA <sup>T</sup>	MME	MM-Vet	Average
<i>Upper Bound, 2880 Tokens</i>									
LLaVA-NeXT-7B	64.2	57.6	70.1	67.4	86.5	64.9	1519.0	43.9	100.0%
<i>Ratio=50%, Retain up to 1440 Tokens</i>									
FastV (ECCV24)	61.8	54.9	69.0	67.4	85.5	59.6	1490.3	37.6	95.5%
PDrop (CVPR25)	63.7	<b>57.9</b>	<b>69.2</b>	<b>67.7</b>	87.9	61.6	1499.6	37.5	97.4%
SparseVLM (ICML25)	63.7	57.2	68.3	67.6	87.9	60.5	1507.2	36.8	96.8%
FasterVLM (2024.12)	63.4	56.4	69.1	67.4	87.7	58.9	1533.3	39.6	97.3%
<b>GlobalCom<sup>2</sup></b>	<b>63.9</b>	56.5	68.5	67.6	<b>88.1</b>	<b>62.3</b>	<b>1552.9</b>	<b>40.4</b>	<b>98.5%</b>
<i>Ratio=25%, Retain up to 720 Tokens</i>									
FastV (ECCV24)	60.4	54.2	<b>68.8</b>	65.6	83.1	58.4	1477.3	35.4	93.4%
PDrop (CVPR25)	60.3	<b>56.8</b>	68.5	65.6	85.5	59.8	1473.7	31.1	93.3%
SparseVLM (ICML25)	59.9	56.0	67.5	65.6	85.0	58.3	1465.9	38.5	94.6%
FasterVLM (2024.12)	61.3	55.4	67.1	<b>66.0</b>	87.2	58.8	1454.6	37.8	94.8%
<b>GlobalCom<sup>2</sup></b>	<b>61.5</b>	55.7	68.1	65.9	<b>87.6</b>	<b>60.9</b>	<b>1493.5</b>	<b>40.7</b>	<b>96.7%</b>
<i>Ratio=10%, Retain up to 288 Tokens</i>									
FastV (ECCV24)	55.9	53.1	68.1	61.6	71.7	55.7	1282.9	27.2	85.4%
PDrop (CVPR25)	54.5	54.4	67.7	59.0	77.6	54.4	1262.1	24.0	84.3%
SparseVLM (ICML25)	56.3	52.1	68.5	60.0	80.1	53.9	1334.2	26.5	86.1%
PruMerge (ICCV25)	53.6	54.0	66.4	61.3	60.8	50.6	1149.3	25.5	80.6%
FasterVLM (2024.12)	56.9	52.6	66.5	61.6	83.6	56.5	1359.2	35.0	89.9%
<b>GlobalCom<sup>2</sup></b>	<b>57.1</b>	<b>54.6</b>	<b>68.7</b>	<b>61.8</b>	<b>83.8</b>	<b>58.4</b>	<b>1365.5</b>	<b>36.4</b>	<b>91.6%</b>

## Ablation Studies

Method	SQA POPE VQA <sup>T</sup> MME MM-Vet					Avg.
<i>Upper Bound, 2880 Tokens</i>						
Vanilla	70.1	86.5	64.9	1519.0	43.9	100.0%
<i>Ratio=25%, Retain up to 720 Tokens</i>						
Uniform	67.1	87.2	60.1	1454.6	37.8	94.2%
$\Pi_{\text{top-}k}$	67.4	87.3	59.8	1471.5	35.7	94.5%
Softmax (max)	67.3	87.2	60.3	1462.6	38.4	94.7%
<b>Softmax (sum)</b>	<b>67.6</b>	<b>87.4</b>	<b>60.6</b>	<b>1473.3</b>	<b>39.6</b>	<b>95.6%</b>

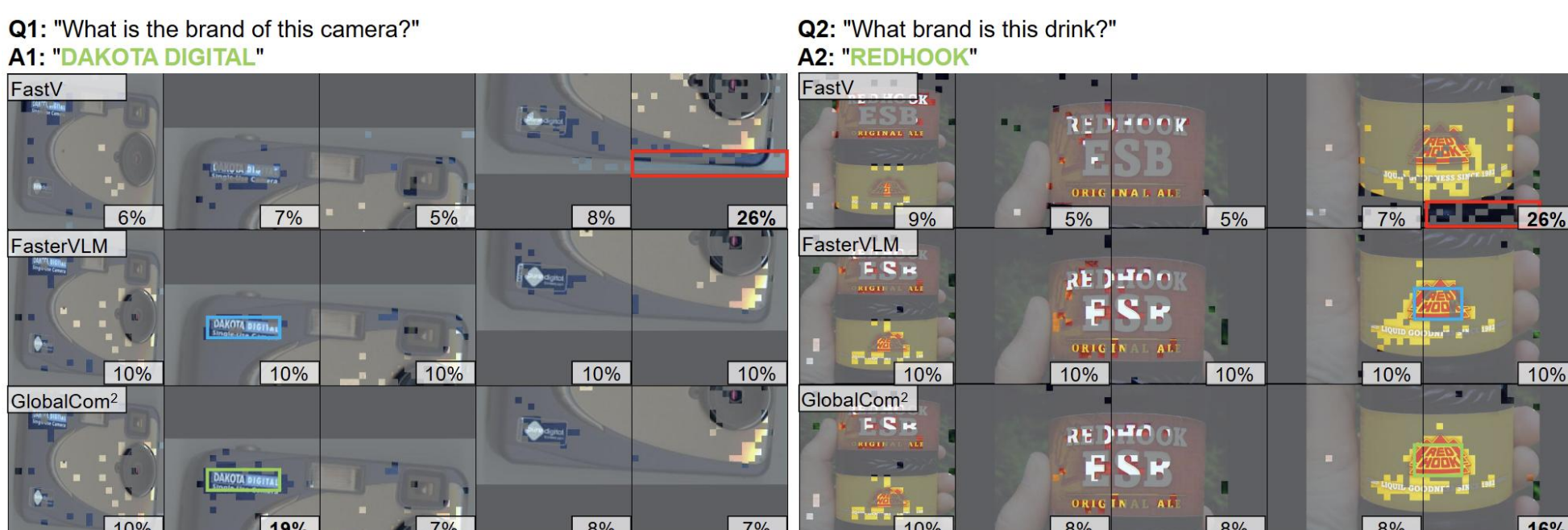
  

Method	SQA POPE VQA <sup>T</sup> MME MM-Vet					Avg.
<i>Upper Bound, 2880 Tokens</i>						
Vanilla	70.1	86.5	64.9	1519.0	43.9	100.0%
<i>Ratio=25%, Retain up to 720 Tokens</i>						
Local only	67.6	87.4	60.6	1473.3	39.6	95.6%
Global only	67.9	86.4	60.2	1488.5	37.8	94.7%
<b>Global and Local</b>	<b>68.1</b>	<b>87.6</b>	<b>60.9</b>	<b>1493.5</b>	<b>40.7</b>	<b>96.7%</b>

Method	TFLOPs↓	Memory↓	Throughput↑	Performance↑
<i>Upper Bound, 2880 Tokens</i>				
Vanilla	41.7	23.0	3.8	100%
<i>Ratio=10%, Retain up to 288 Tokens</i>				
SparseVLM	5.4 (↓87.1%)	24.2 (↑5.2%)	5.9 (1.6×)	85.7%
FasterVLM	<b>3.8</b> (↓90.9%)	<b>13.6</b> (↓40.1%)	<b>6.7</b> (1.8×)	89.5%
<b>GlobalCom<sup>2</sup></b>	<b>3.8</b> (↓90.9%)	<b>13.9</b> (↓40.0%)	<b>6.7</b> (1.8×)	<b>90.8%</b>

## Token Compression Visualizations: Better Information Preservation



## Broader Applicability of GlobalCom<sup>2</sup>

**Compatible with other methods:** Plug-in boost for SparseVLM and FastV.

