# Department of Computer Science and Engineering, CUHK

# Empowering Mobile Devices with Multimodal Large Language Models in 6G Era

Chengyao Wang, from Prof. Jiaya Jia's Group

## Motivation

The rapid advancements in multimodal large language models offer significant potential for 6G wireless communication systems, where AI-driven applications will require real-time data processing, visual understanding, and reasoning capabilities. However, the computational demands of these models, especially for cloud-based inference, remain a challenge. To address this, we propose Mini-Gemini, an initiative focused on improving open-source models and optimizing them for efficient local inference on mobile devices and personal computers, reducing reliance on closed-source systems. This approach aligns with 6G's vision of low-latency, distributed AI, enabling seamless, accessible, and efficient AI-powered applications across a wide range of devices in the 6G era.

## Challenge

1. There is a huge performance gap between close source and open source VLMs.
2. Large language model is hard to deployed on edge devices, due to huge inference cost.

## Contribution

1. With data iteration and better utilization of the image data resolution, we propose token simplification to reduce computation.
2. Our approach attains leading performance in various settings and even surpasses many private models.
3. With instruction driven visual information mining approach, our model can efficiently deployed on laptop and mobile phone.
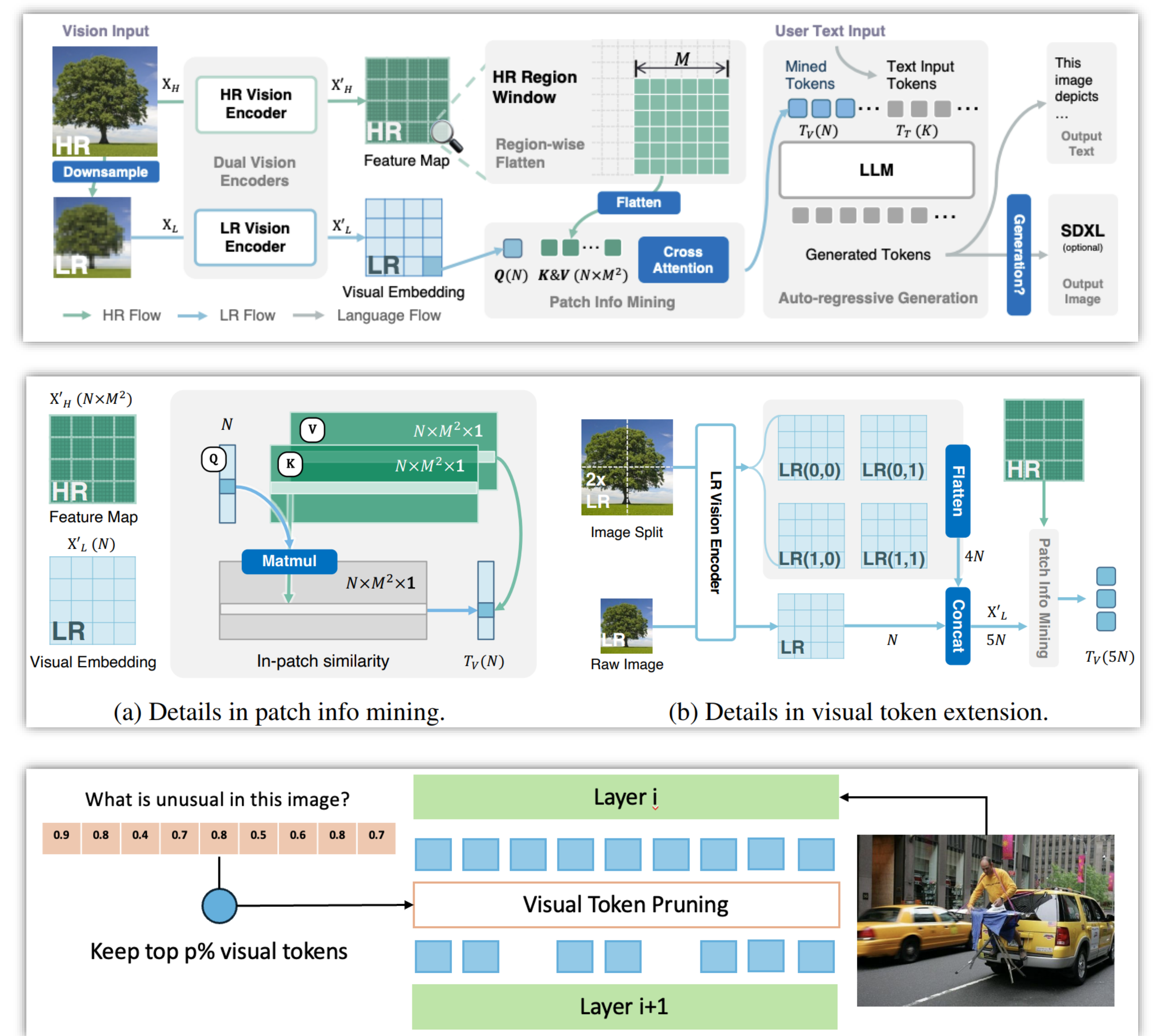


(a) Details in patch info mining.  (b) Details in visual token extension.



## Image Understanding and Generation



High-Resolution Understanding

Benchmark Performance

Generation with Reasoning

## Multimodal Agent



(a) The detailed results on Android with observation, thoughts, and next action.

(b) End-to-end episode demonstration for completing task on a Windows10 environment.

## Mobile Deployment



## Quantitative Results

| Method | LLM | Res. | VQA$^T$ | MMB | MME | MM-Vet | MMMU$_v$ | MMMU$_t$ | MathVista |
|---|---|---|---|---|---|---|---|---|---|
| *Normal resolution setting* | | | | | | | | | |
| MobileVLM[63] | MLLaMA 2.7B | 336 | 47.5 | 59.6 | 1289 | – | – | – | – |
| InstructBLIP[42] | Vicuna-7B | 224 | 50.1 | 36.0 | – | 26.2 | – | – | – |
| InstructBLIP[42] | Vicuna-13B | 224 | 50.7 | – | 1213 | 25.6 | – | – | 25.3 |
| Qwen-VL[23] | Qwen-7B | 448 | 63.8* | 38.2 | – | – | – | – | – |
| Qwen-VL-Chat[23] | Qwen-7B | 448 | 61.5* | 60.6 | 1488 | – | 35.9 | 32.9 | – |
| Shikra[64] | Vicuna-13B | 224 | – | 58.8 | – | – | – | – | – |
| IDEFICS-80B[65] | LLaMA-65B | 224 | 30.9 | 54.5 | – | – | – | – | – |
| LLaMA-VID[10] | Vicuna-7B | 336 | – | 65.1 | 1521 | – | – | – | – |
| LLaMA-VID[10] | Vicuna-13B | 336 | – | 66.6 | 1542 | – | – | – | – |
| LLaVA-1.5[43] | Vicuna-7B | 336 | 58.2 | 65.2 | 1511 | 31.1 | – | – | – |
| LLaVA-1.5[43] | Vicuna-13B | 336 | 61.3 | 69.2 | 1531/295 | 36.1 | 36.4 | 33.6 | 27.6 |
| Mini-Gemini | Gemma-2B | 336 | 56.2 | 59.8 | 1341/312 | 31.1 | 31.7 | 29.1 | 29.4 |
| Mini-Gemini | Vicuna-7B | 336 | 65.2 | 69.3 | 1523/316 | 40.8 | 36.1 | 32.8 | 31.4 |
| Mini-Gemini | Vicuna-13B | 336 | 65.9 | 68.5 | 1565/322 | 46.0 | 38.1 | 33.5 | 37.0 |
| Mini-Gemini | Mixtral-8x7B | 336 | 69.2 | 75.6 | 1639/379 | 45.8 | 41.8 | 37.1 | 41.8 |
| Mini-Gemini | Hermes-2-Yi-34B | 336 | 70.1 | 79.6 | 1666/439 | 53.0 | 48.7 | 43.6 | 38.9 |
| *High resolution setting* | | | | | | | | | |
| OtterHD[12] | Fuyu-8B | 1024 | – | 53.6 | 1314 | – | – | – | – |
| CogVLM-Chat[66] | Vicuna-7B | 490 | 70.4* | 63.7 | – | 51.1 | 41.1 | – | 34.5 |
| LLaVA-NeXT[11] | Vicuna-7B | 672 | 64.9 | 68.1 | 1519/332 | 43.9 | 35.8 | – | 34.6 |
| LLaVA-NeXT[11] | Vicuna-13B | 672 | 67.1 | 70.7 | 1575/326 | 48.4 | 36.2 | – | 35.3 |
| LLaVA-NeXT[11] | Hermes-2-Yi-34B | 672 | 69.5 | 79.6 | 1631/397 | 57.4 | 51.1 | 44.7 | 46.5 |
| Mini-Gemini-HD | Vicuna-7B | 672 | 68.4 | 65.8 | 1546/319 | 41.3 | 36.8 | 32.9 | 32.2 |
| Mini-Gemini-HD | Vicuna-13B | 672 | 70.2 | 68.6 | 1597/320 | 50.5 | 37.3 | 35.1 | 37.0 |
| Mini-Gemini-HD | Mixtral-8x7B | 672 | 71.9 | 74.7 | 1633/356 | 53.5 | 40.0 | 37.0 | 43.1 |
| Mini-Gemini-HD | Hermes-2-Yi-34B | 672 | 74.1 | 80.6 | 1659/482 | 59.3 | 48.0 | 44.9 | 43.3 |
| *Private models* | | | | | | | | | |
| Gemini Pro[5] | Private | – | 74.6 | 75.2 | – | 64.3 | 47.9 | – | 45.2 |
| Qwen-VL-Plus[23] | Private | – | 78.9 | 66.2 | – | 45.2 | 40.8 | – | 43.3 |
| GPT-4V[4] | Private | – | 78.0 | 75.1 | – | 67.6 | 56.8 | 55.7 | 49.9 |

## Our Related Publications

1. Li, Yanwei*, Chengyao Wang*, and Jiaya Jia. "Llama-vid: An image is worth 2 tokens in large language models." In *European Conference on Computer Vision*, pp. 323-340. Springer, Cham, 2025.

2. Li, Yanwei*, Yuechen Zhang*, Chengyao Wang*, Zhisheng Zhong, Yixin Chen, Ruihang Chu, Shaoteng Liu, and Jiaya Jia. "Mini-gemini: Mining the potential of multi-modality vision language models." *arXiv preprint arXiv:2403.18814* (2024).

## Acknowledgment