



# AdaLLaVA: Learning to Inference Adaptively for Multimodal Large Language Models

Zhuoyan Xu $^{*1}$ , Khoi Duc Nguyen $^{*1}$ , Preeti Mukherjee $^2$ , Saurabh Bagchi $^2$ , Somali Chaterji $^2$ , Yingyu Liang $^{1,3}$ , Yin Li $^1$ <sup>1</sup>University of Wisconsin-Madison, <sup>2</sup>Purdue University, <sup>3</sup>The University of Hong Kong

## Multimodal Large Language Models Visual --- Connecter -----Modality **→** Encoder Key Challenge Answer 100% Compute Simple request 100% Latency Conventional MLLMAnswer 100% Compute Hard request 100% Latency Fixed model under different input 100% compute resources Answer 100% Compute 100% Latency Hard request Conventional 50% compute resources Answer 100% Compute 200% Latency Same content with different compute budget Can we make MLLM adaptive to varying compute resources and input contents? Goal

60% compute budget

85% compute budget

00% compute budget

What is the image showing?

The image is showing a painting or

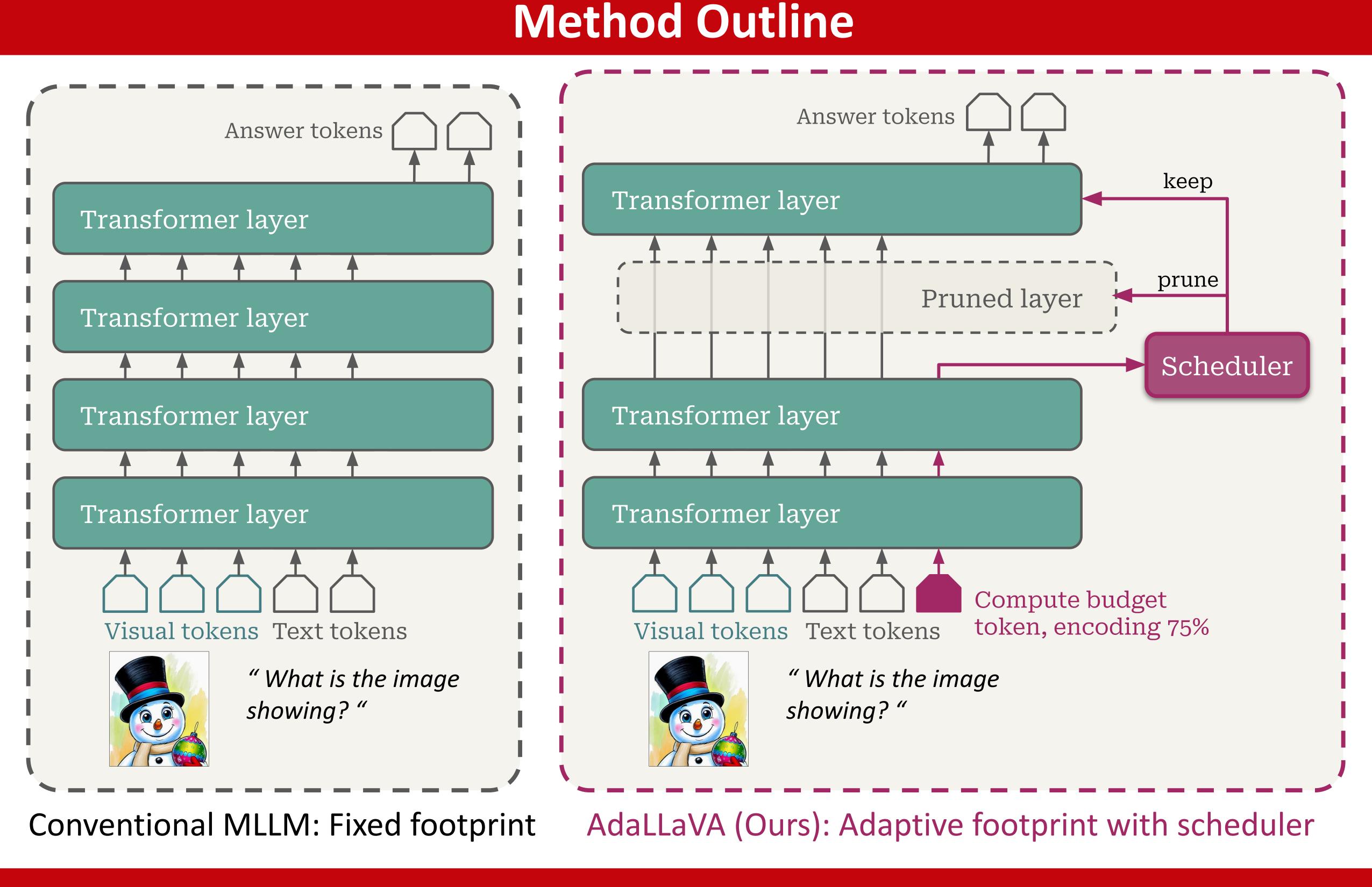
drawing of a snowy winter scene.

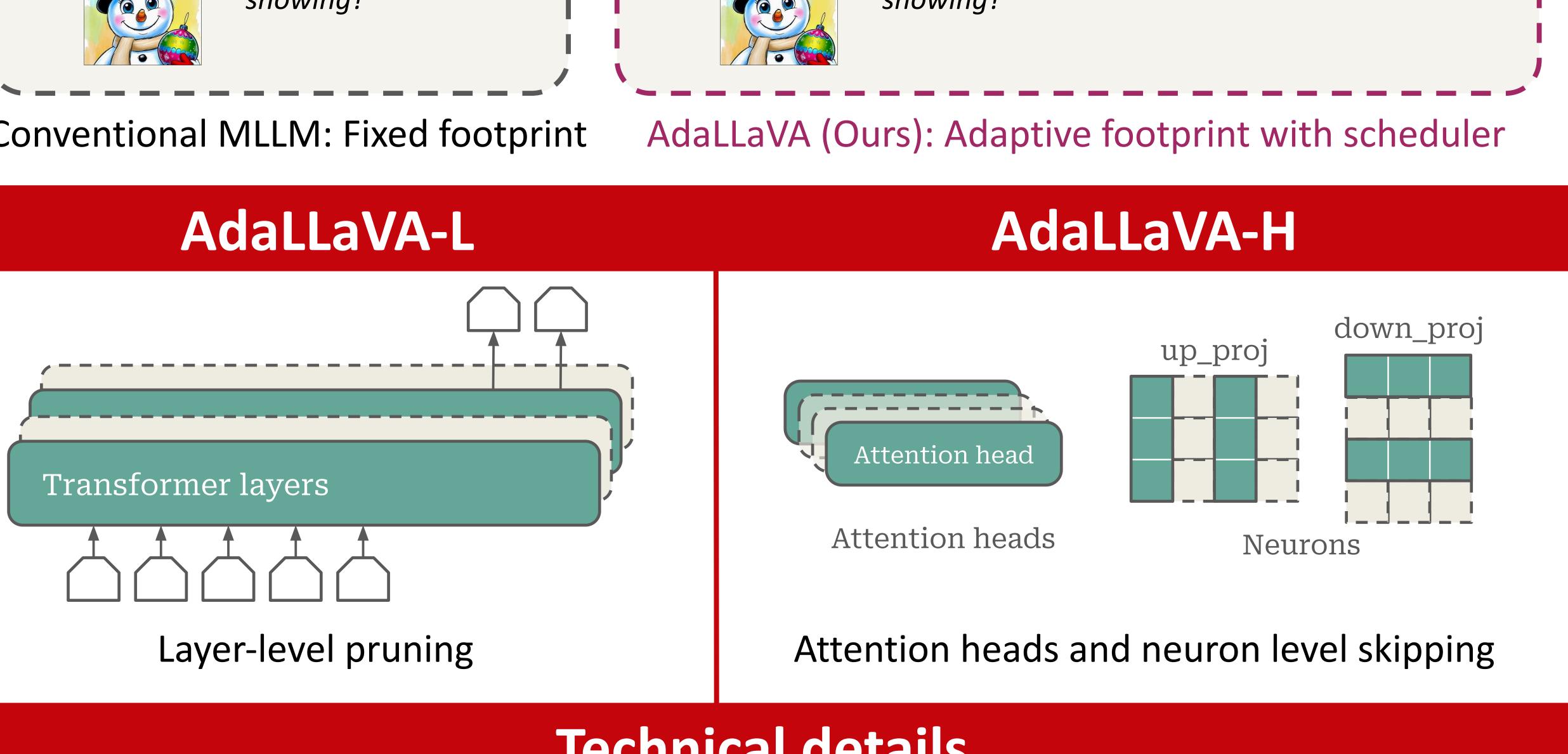
The image shows a snowman holding

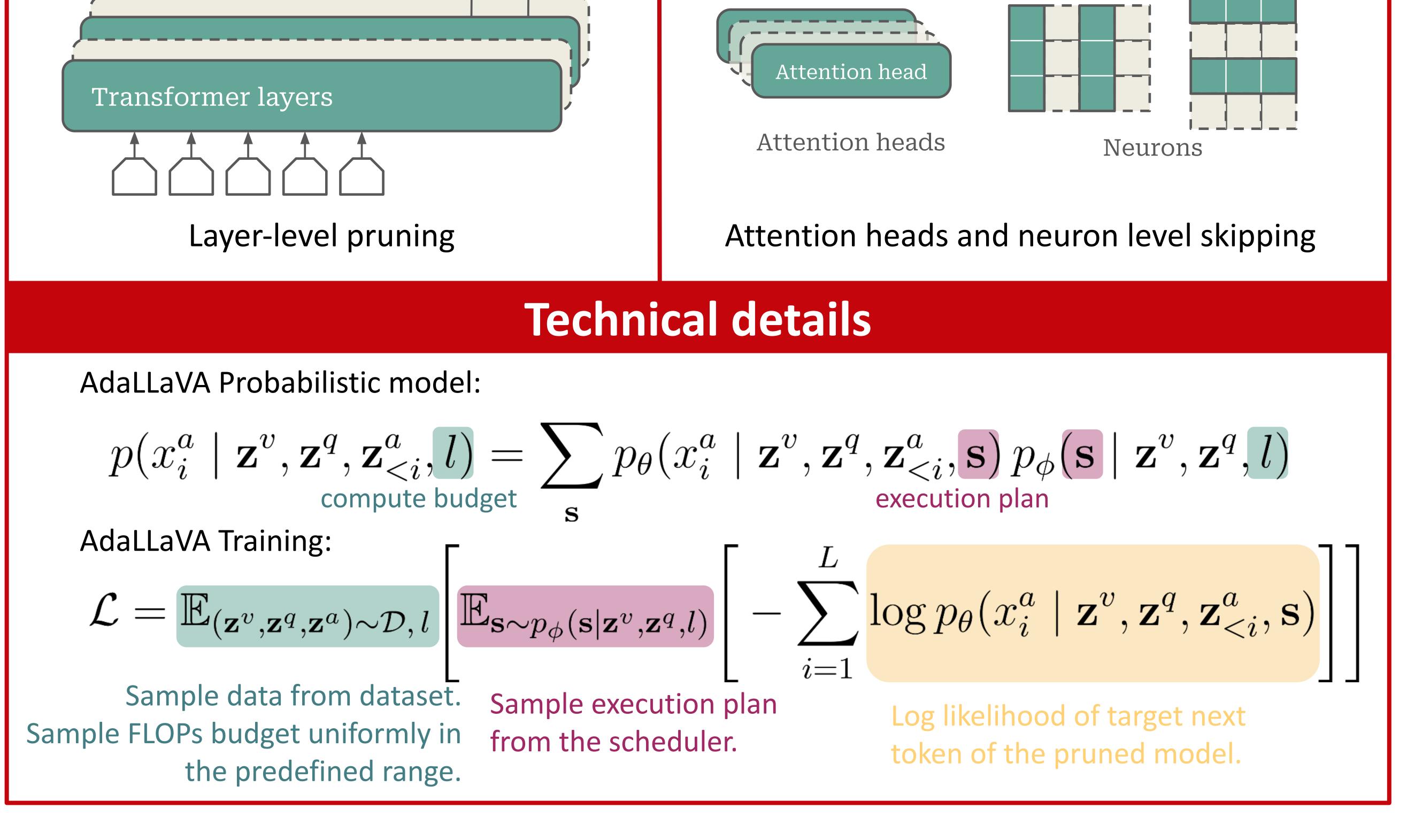
The image shows a snowman holding

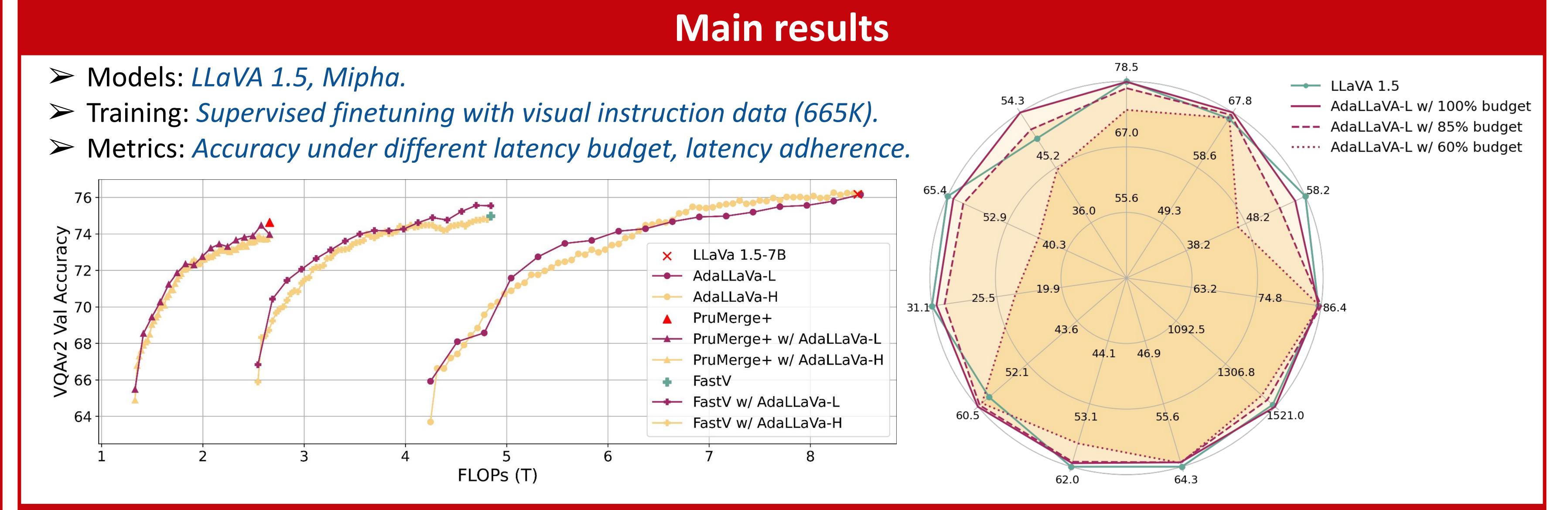
a colorful ball.

a colorful egg.

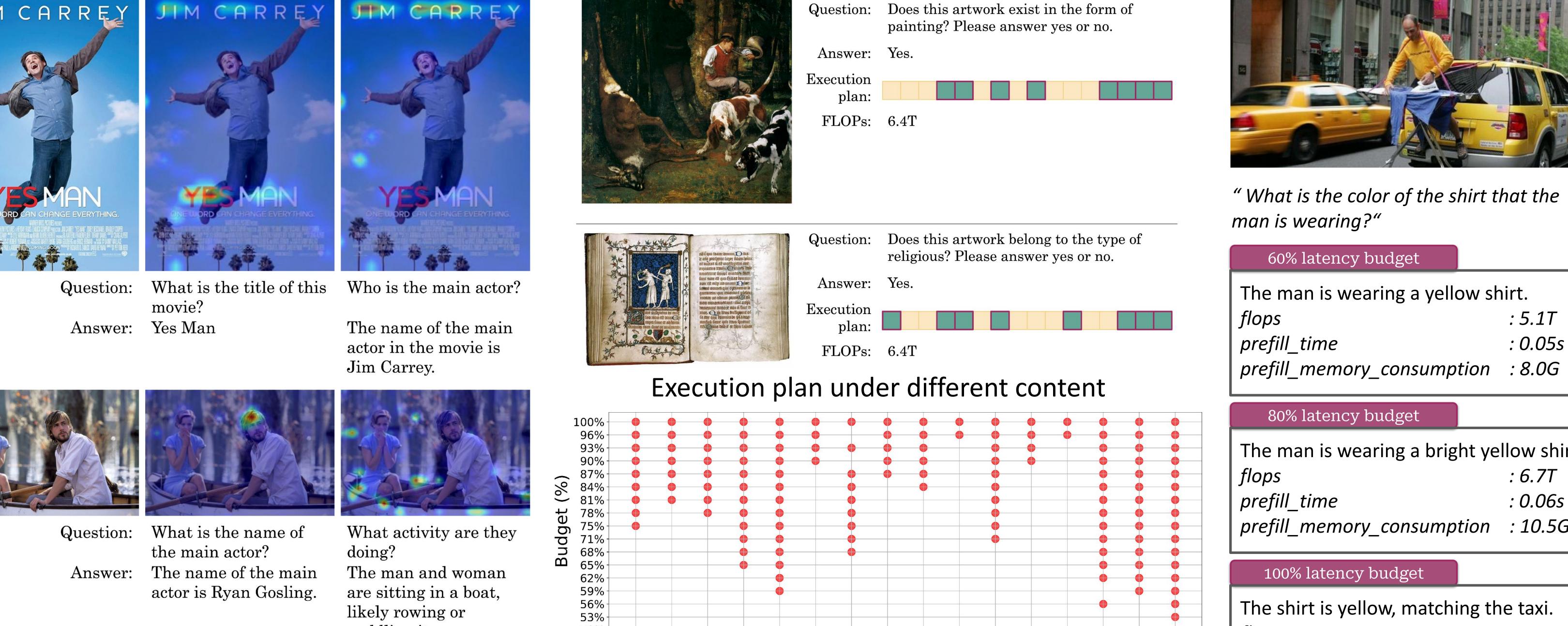


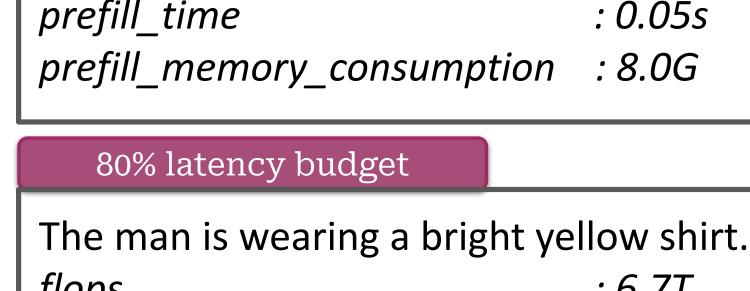




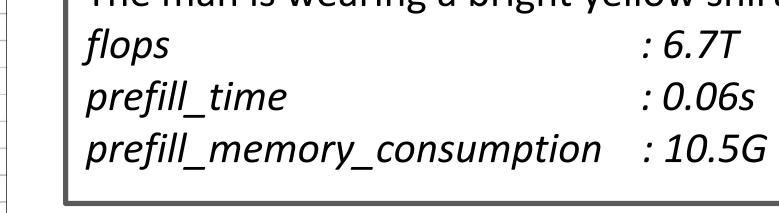


## Qualitative results

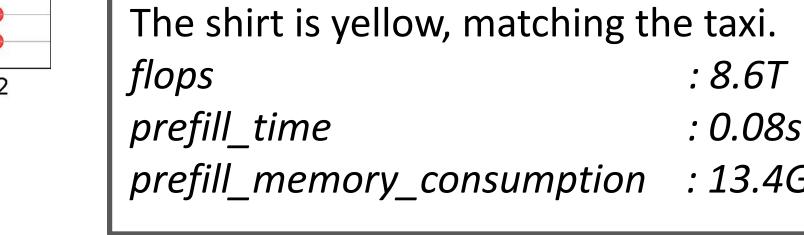




60% latency budget



### 100% latency budget



### Conclusions

- Novel adaptive inference MLLM framework adapt to compute budget and input content.
- Latency-aware scheduler and probabilistic modeling approach.

Latency token attention mapping

• Validated across benchmarks and models, integrates with token selection techniques.

This work was supported by National Science Foundation (CNS-2333487/2333491, CNS-2146449, IIS-2442739), Army Research Lab (W911NF-2020221), Google, and AWS.

Execution plan under different budget