

## WHY DUET-VLM?

Vision tokens dominate VLM compute. DUET-VLM first merges redundant visual tokens, then lets text saliency drop the rest inside the LLM — **keeping accuracy while cutting compute.**

**>99%**

**accuracy retained**  
67% fewer visual tokens

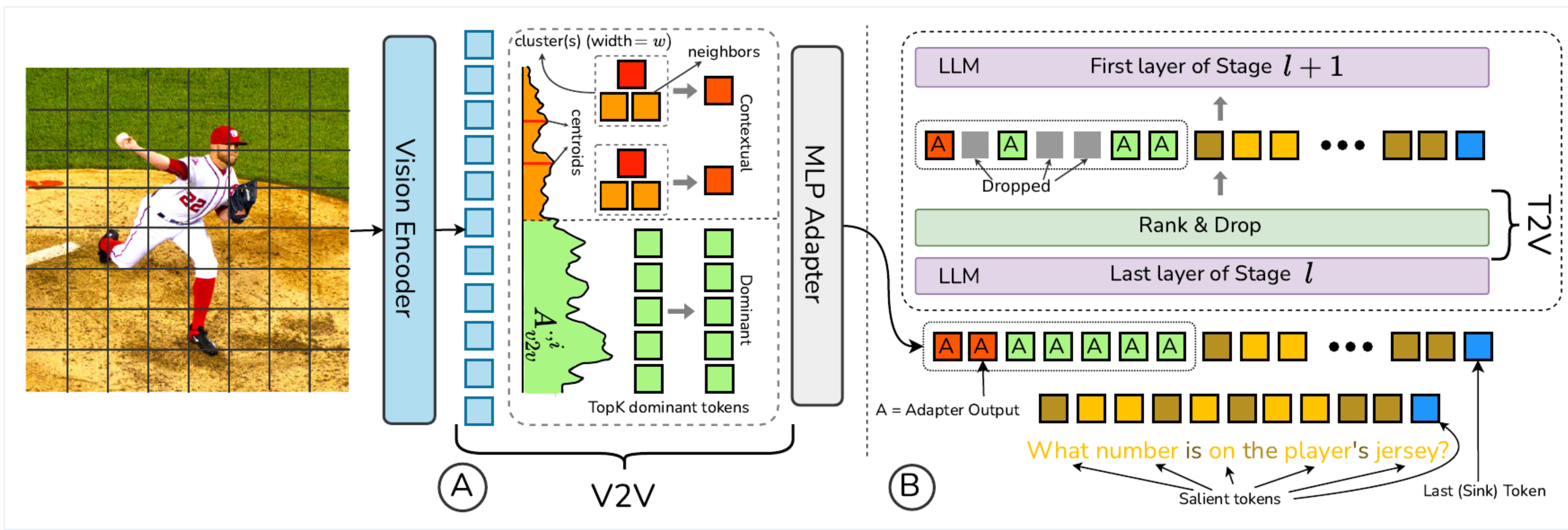
**>97%**

**accuracy retained**  
89% fewer visual tokens

**31%**

**training time saved**  
LLaVA-1.5-7B backbone

## METHOD: TWO-STAGE TOKEN COMPRESSION



### (A) V2V merge

keep dominant tokens, locally merge context

### (B) T2V prune

drop visual tokens guided by text saliency

Training uses the same compressed-token path as inference, so the model adapts to small visual budgets instead of treating pruning as a post-hoc trick.

## RESULTS: ACCURACY SURVIVES COMPRESSION

### Accuracy stays high as tokens disappear

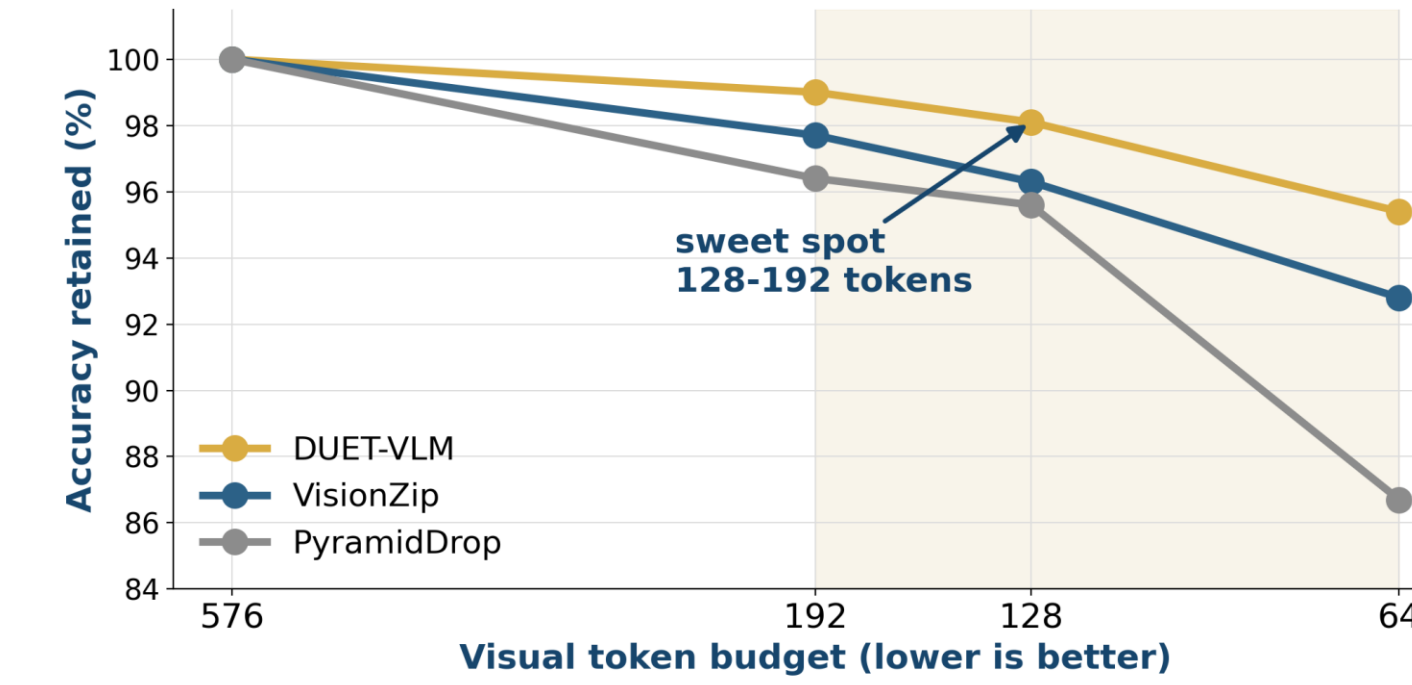


Fig. 1 Inference accuracy vs visual-token budget on LLaVA-1.5-7B (avg of 5 benchmarks).

### Training-time savings on LLaVA-1.5-7B

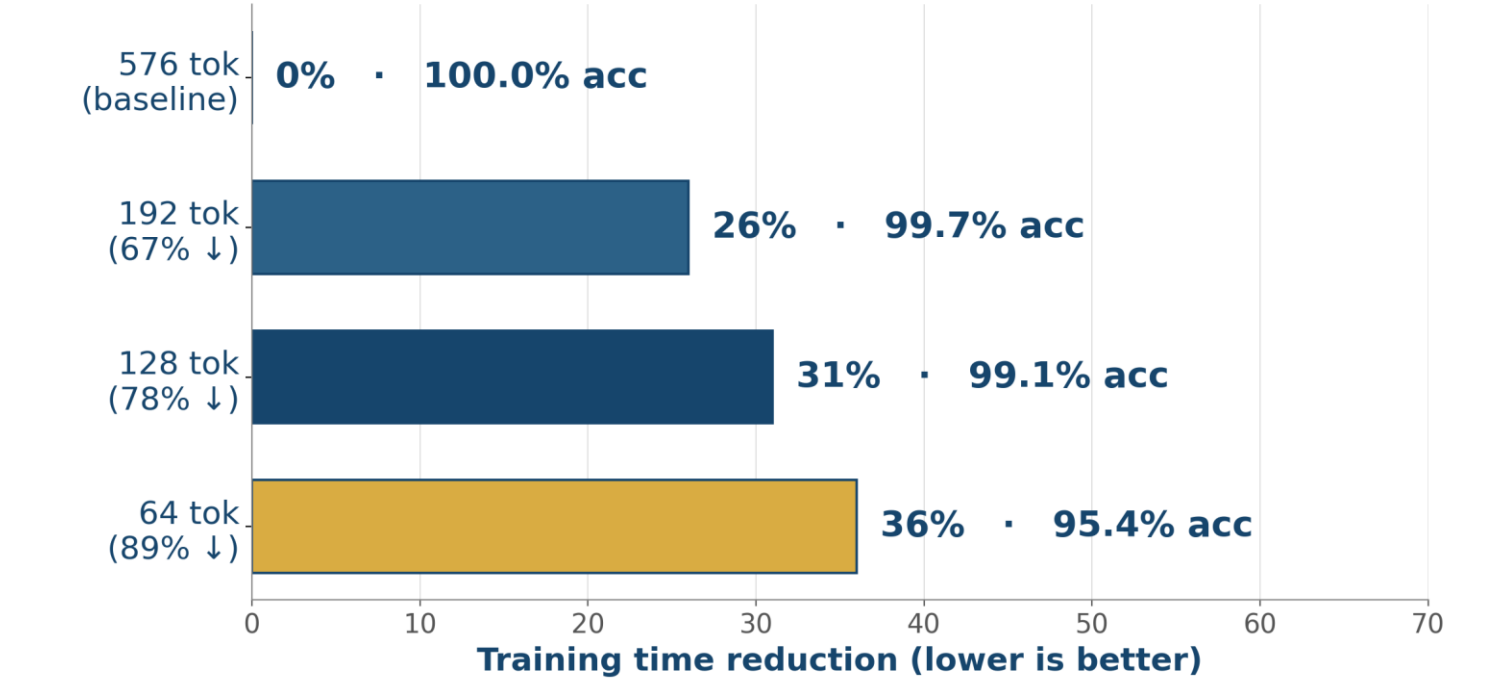


Fig. 2 End-to-end training time saved while preserving baseline accuracy.

### Works beyond one backbone

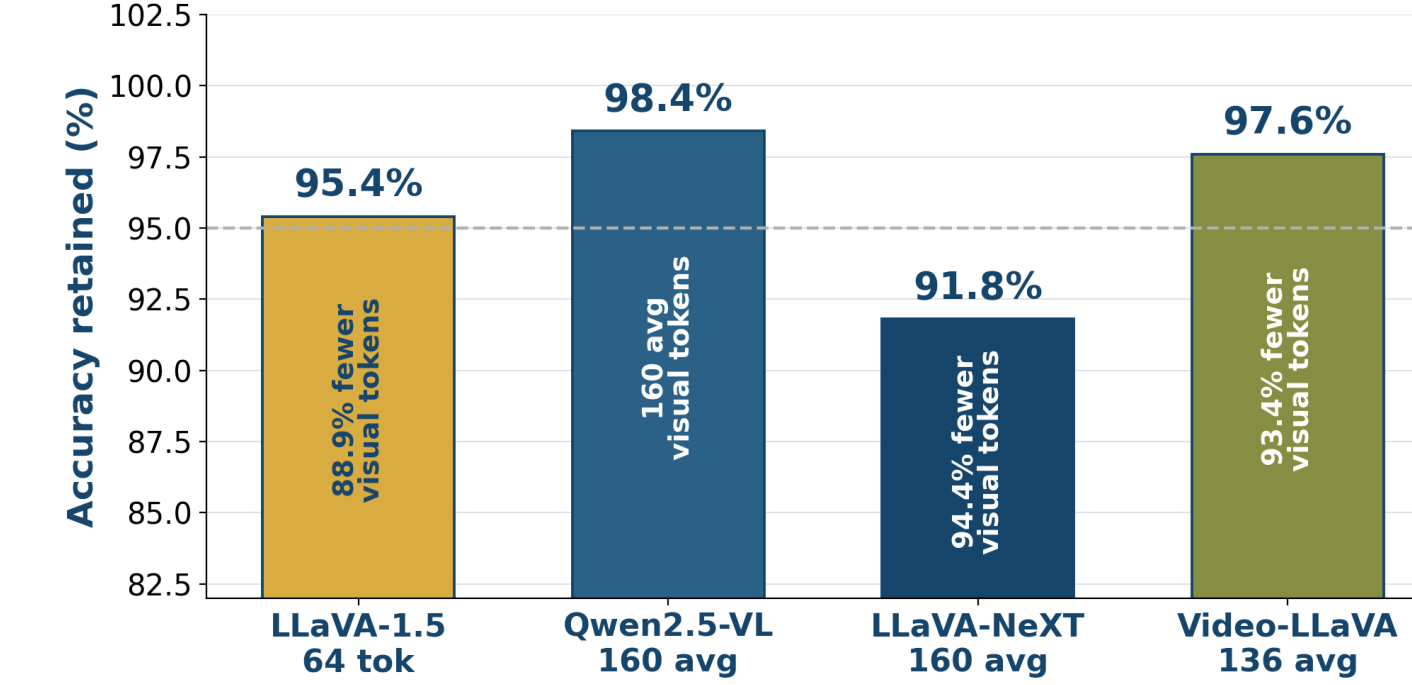


Fig. 3 DUET-VLM generalizes across LLaVA-1.5, LLaVA-NeXT, Qwen2.5-VL and Video-LLaVA backbones.

### Video: accuracy survives 93.4% token drop

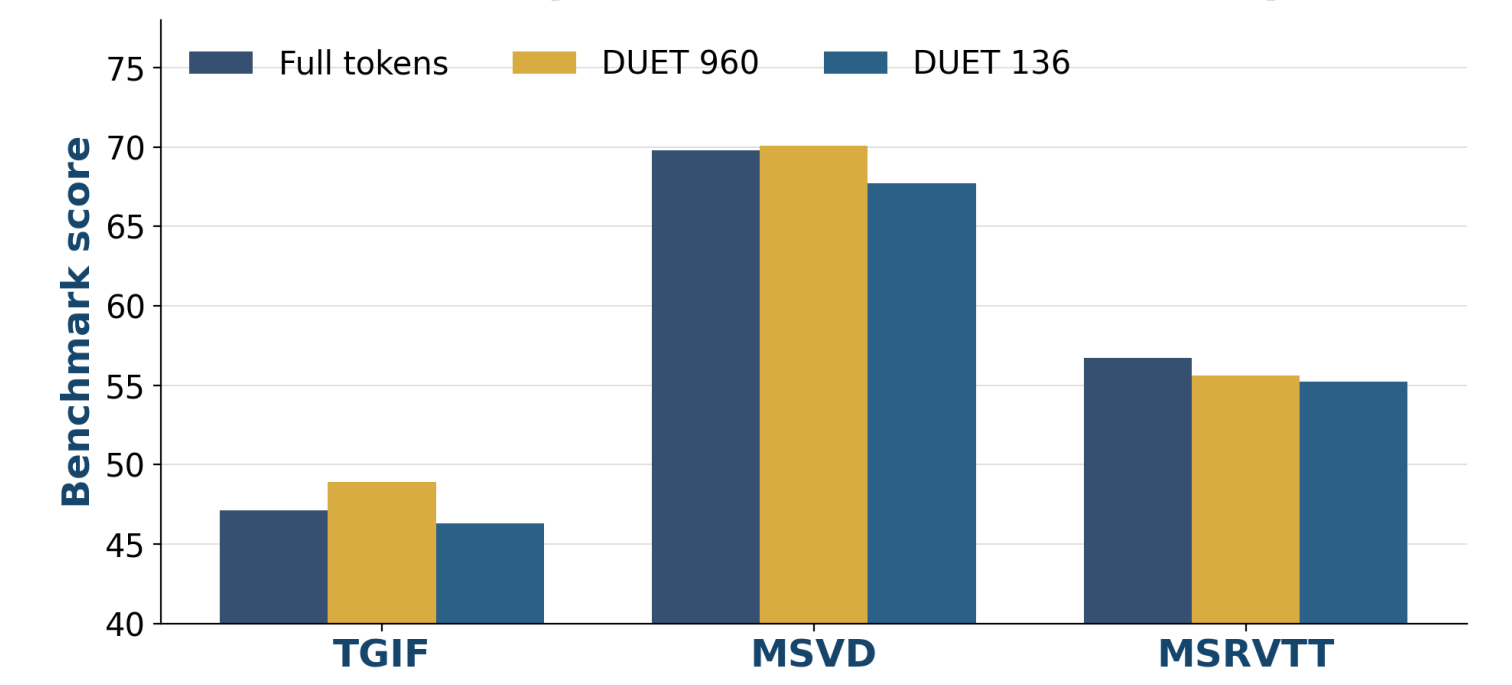
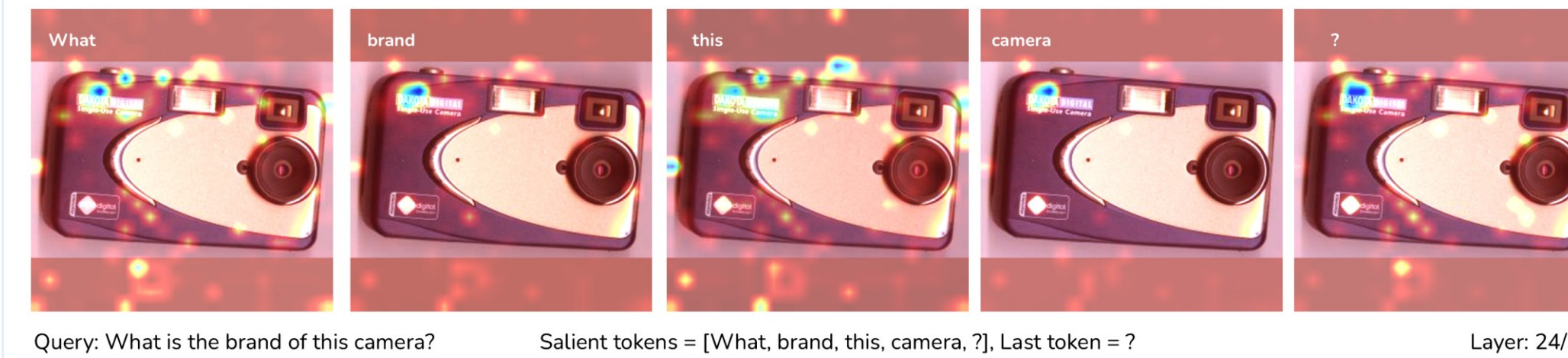


Fig. 4 Video-LLaVA-7B: 93.4% of visual tokens dropped, only ~2% accuracy lost.

## QUALITATIVE EVIDENCE

### Text-guided attention keeps the patches needed for the query.

The token budget shrinks, but semantically useful visual regions remain visible.



## TAKEAWAY + LINKS

### What to remember

DUET-VLM is a plug-in token budget controller. It keeps accuracy high across image and video VLMs.



Paper



Code



Project page