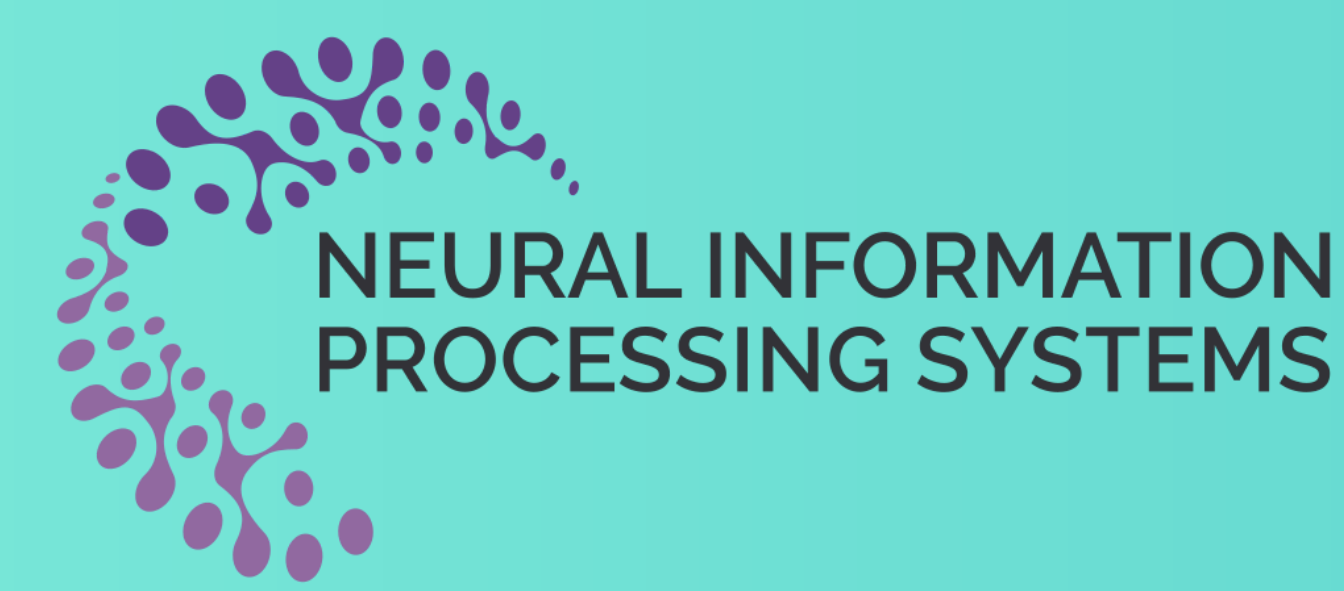# VLMs underpreform on visual data compared to text.

# Why?

Same Task, Different Circuits: Disentangling Modality-Specific Mechanisms in VLMs

Yaniv Nikankin, Dana Arad, Yossi Gandelsman, Yonatan Belinkov
yaniv.n@campus.technion.ac.il

TECHNION Israel Institute of Technology

BAIR BERKELEY ARTIFICIAL INTELLIGENCE RESEARCH

NEURAL INFORMATION PROCESSING SYSTEMS

**Q:** *How many "banana" are in the sequence/image?*

"car flower *banana* chair ... *banana* *banana* chair flower *banana*"

| | |
|---|---|
| **Data** | car flower |
| **Query** | How many |
| **Generation** | ? |

→ 4 ✓

| | |
|---|---|
| **Data** | <img> <img> |
| **Query** | How many |
| **Generation** | ? |

→ 3 ✗

We compare **circuits** — task-specific computational sub-graphs — used to solve textual and visual variants of the same task.

## (1) VLMs use **different** computation paths for vision and language variants of the same task.

We find low overlap in **data** and **query** positions, and some overlap in the **generation** position, calculated by IoU of circuit components.



Qwen2-VL-7B — Normalized IoU (Counting, Arithmetic, Color Ordering, Factual Recall, Sentiment Analysis)

## (2) **Query** and **Generation** components implement **similar** functionality across modalities; **Data** components are modality-specific.

We swap components at matched positions and re-evaluate faithfulness, quantifying how much behavior is preserved.



Qwen2-VL-7B — Normalized Interchange Faithfulness (Counting, Arithmetic, Color Ordering, Factual Recall, Sentiment Analysis)

## (3) **Visual** data tokens align with **text representations** only in late layers, too late to influence subsequent positions. Back-patching closes a third of the performance gap!



Qwen2-VL-7B — Cosine Similarity vs Layer

| Model | Counting | Arithmetic | Sentiment Analysis |
|---|---|---|---|
| **Qwen2-7B-VL** | +4.2% | +9.4% | +2.1% |
| **Pixtral-12B** | +1.5% | +7.6% | +17.8% |
| **Gemma-3-12B** | +2.9% | +6.7% | +5.5% |