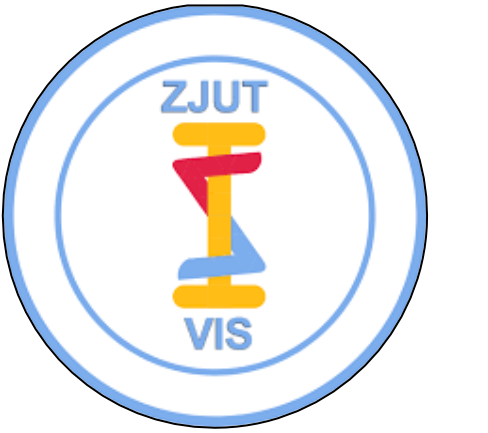# CompoVis: Probing the Compositional Understanding in VLMs with Visualization Representation and Insights

**Tong Li**
ZJUT Vis
litong@zjut.edu.cn

**Guodao Sun**
ZJUT Vis
guodao@zjut.edu.cn

**Xueqian Zheng**
ZJUT Vis
xqzheng@zjut.edu.cn

**Haixia Wang**
ZJUT Vis
hxwang@zjut.edu.cn

**Ronghua Liang**
ZJUT Vis
hxwang@zjut.edu.cn

PacificVis 2025

## 01 Introduction

In vision-language research, "compositional understanding" refers to the ability of a model to jointly handle images and texts.

The model should be able to recognize, comprehend and align each component in vision and textual modality (such as *"red"*, *"hydrant"*, *"man"*) and their combination (the scene semantics of *"a man in white T-shirt leans on a red hydrant"*).
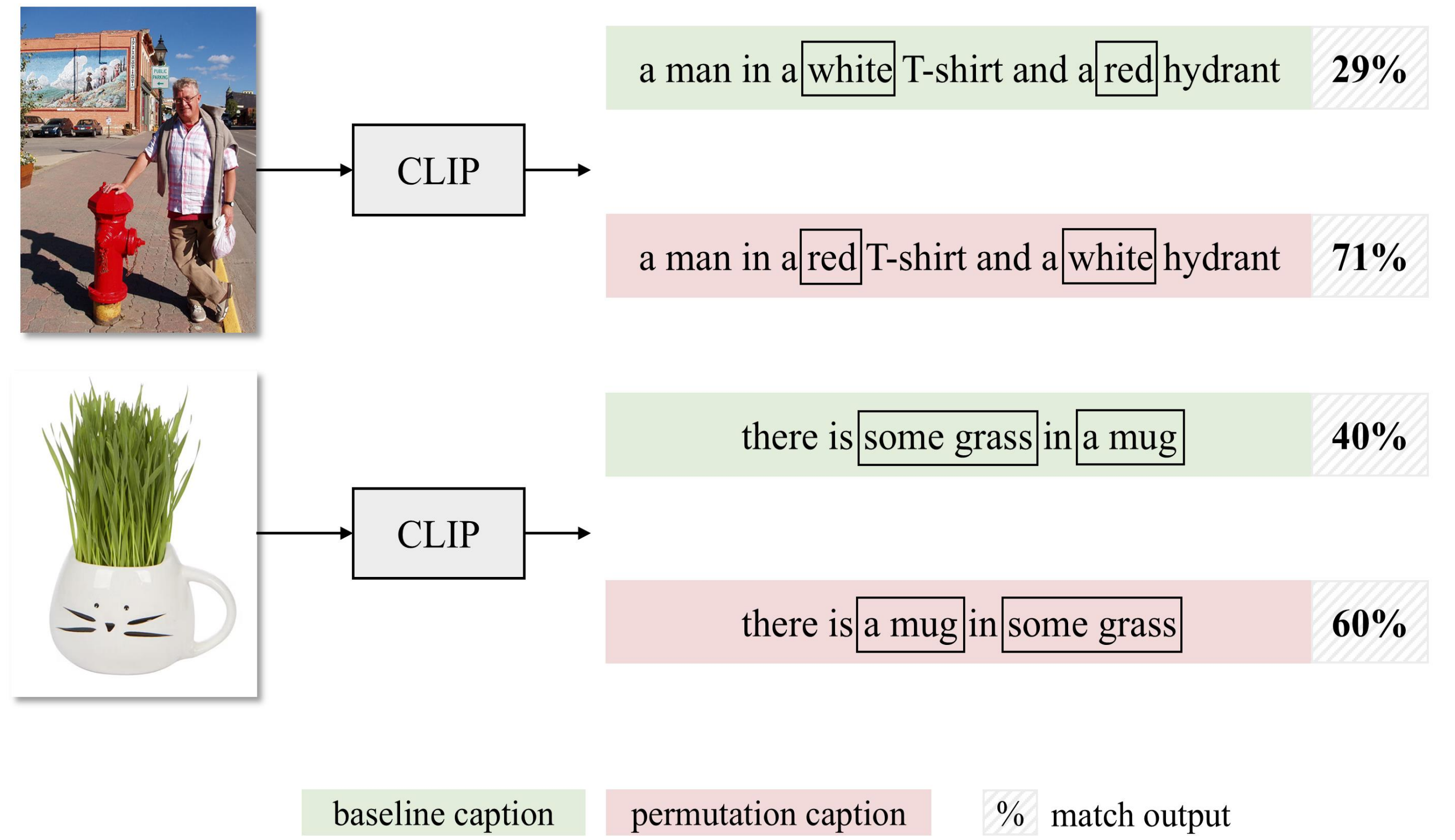
Recent studies reveal even advanced vision-language models (VLMs) struggle with "compositional understanding", particularly when dealing with fine-grained linguistic phenomena.

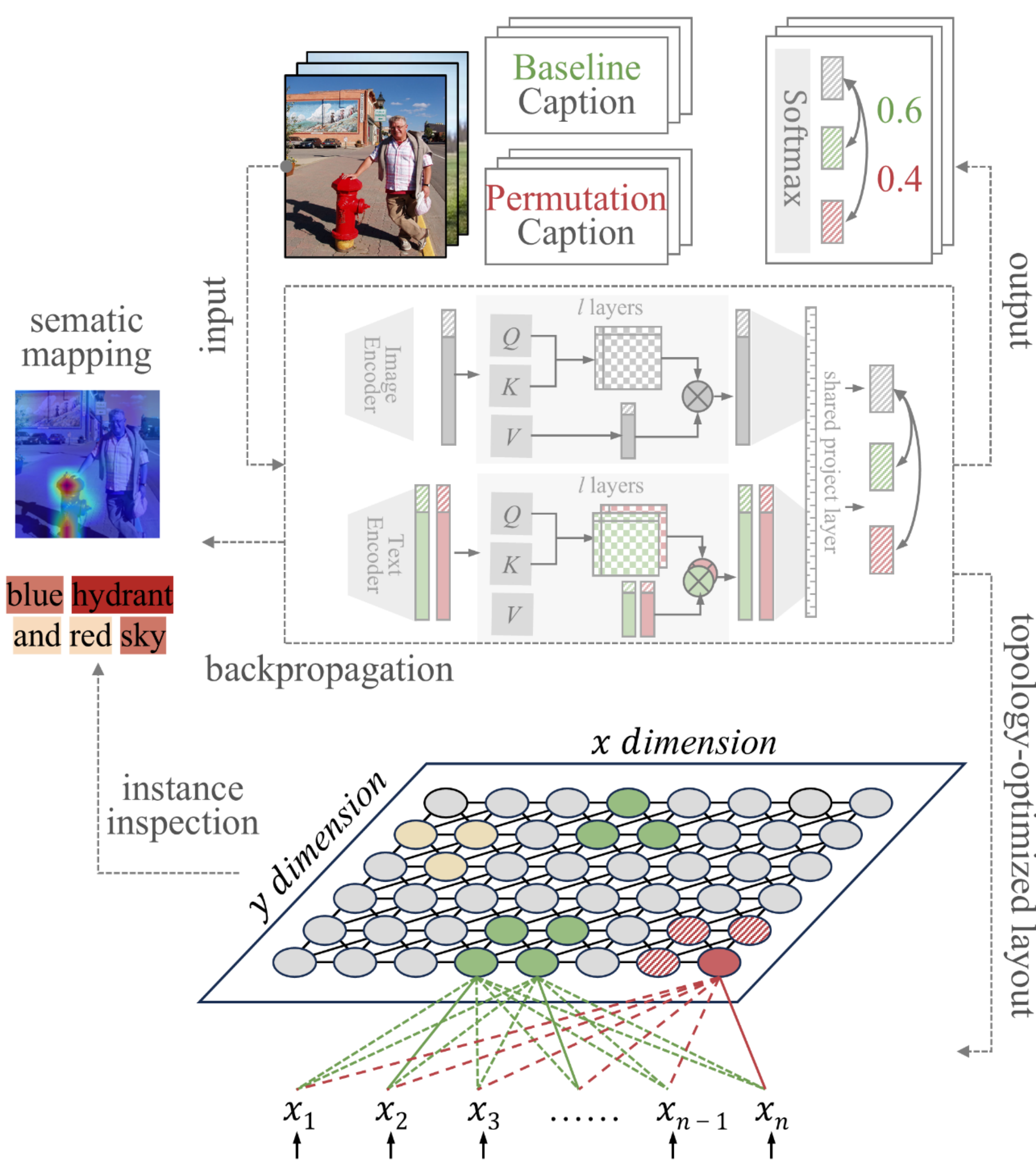**vision-linguistic compositionality**

As shown in the figure, CLIP assigns higher scores to the permutation caption. Computer vision methods focus on quantitative metrics and model architectures, constrained by closed datasets and predefined parameters.

To our knowledge, we are the first to elucidate the "bag-of-objects" behavior of VLMs from a visualization perspective. In summary, the contributions of our paper are as follows:
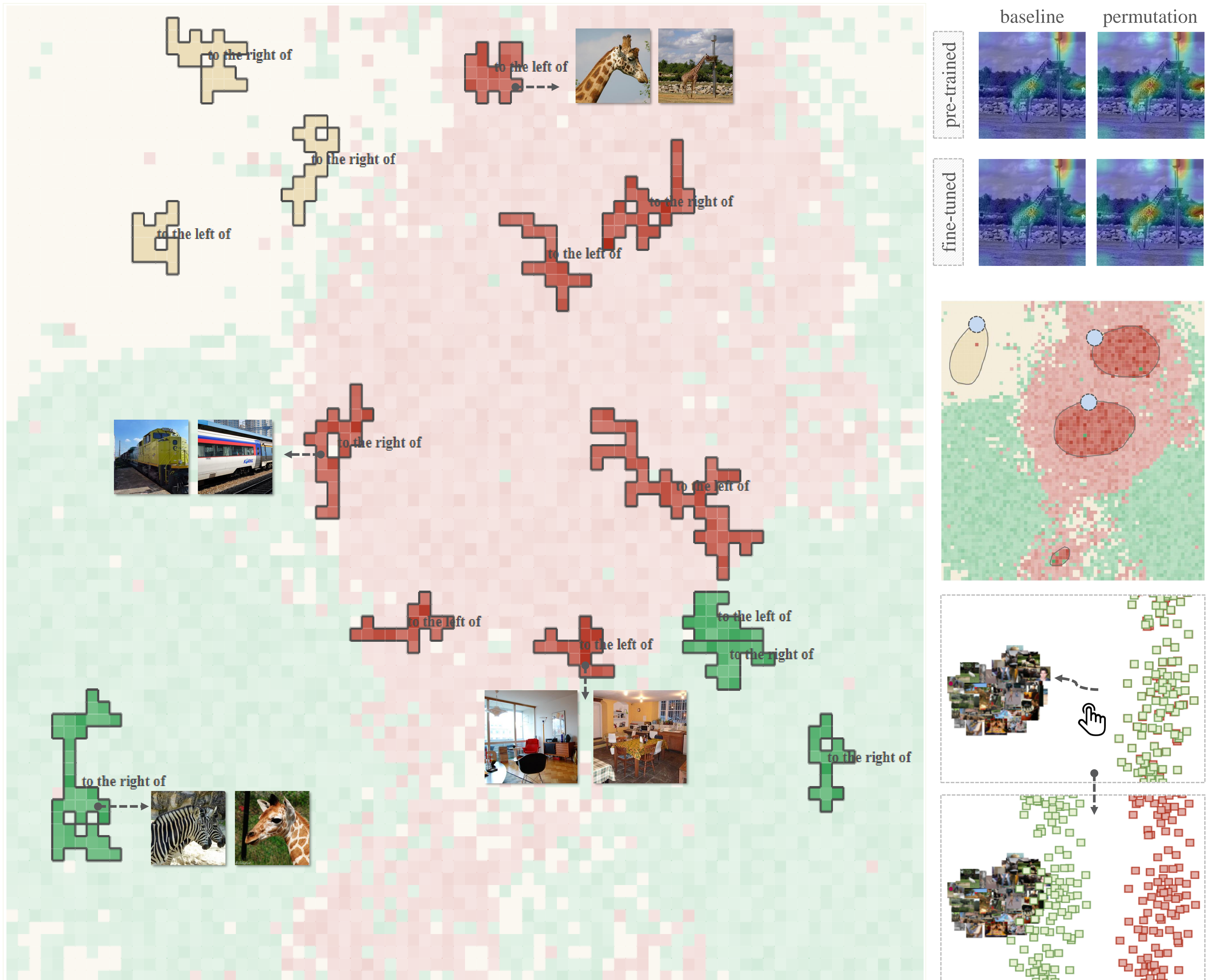
➤ We introduce a visual analysis tool that supports non-expert users in exploring comprehension defects in VLMs, promoting user-guided fine-tuning evaluation.

➤ We optimize the topological layout to enhance the visual proximity between semantically consistent samples in large-scale data, providing a reference for grid layout in the visualization community.

➤ We introduce a permutation dataset, supporting the visual probing of cross-modal alignment ability, and facilitating human-in-the-loop alignment evaluation.



a man in a white T-shirt and a red hydrant — 29%
a man in a red T-shirt and a white hydrant — 71%

there is some grass in a mug — 40%
there is a mug in some grass — 60%

baseline caption    permutation caption    % match output

## 02 Pipeline



Baseline Caption
Permutation Caption
Softmax 0.6 / 0.4

semantic mapping
blue hydrant and red sky
instance inspection

input / output / backpropagation / topology-optimized layout

$x_1$ $x_2$ $x_3$ ...... $x_{n-1}$ $x_n$
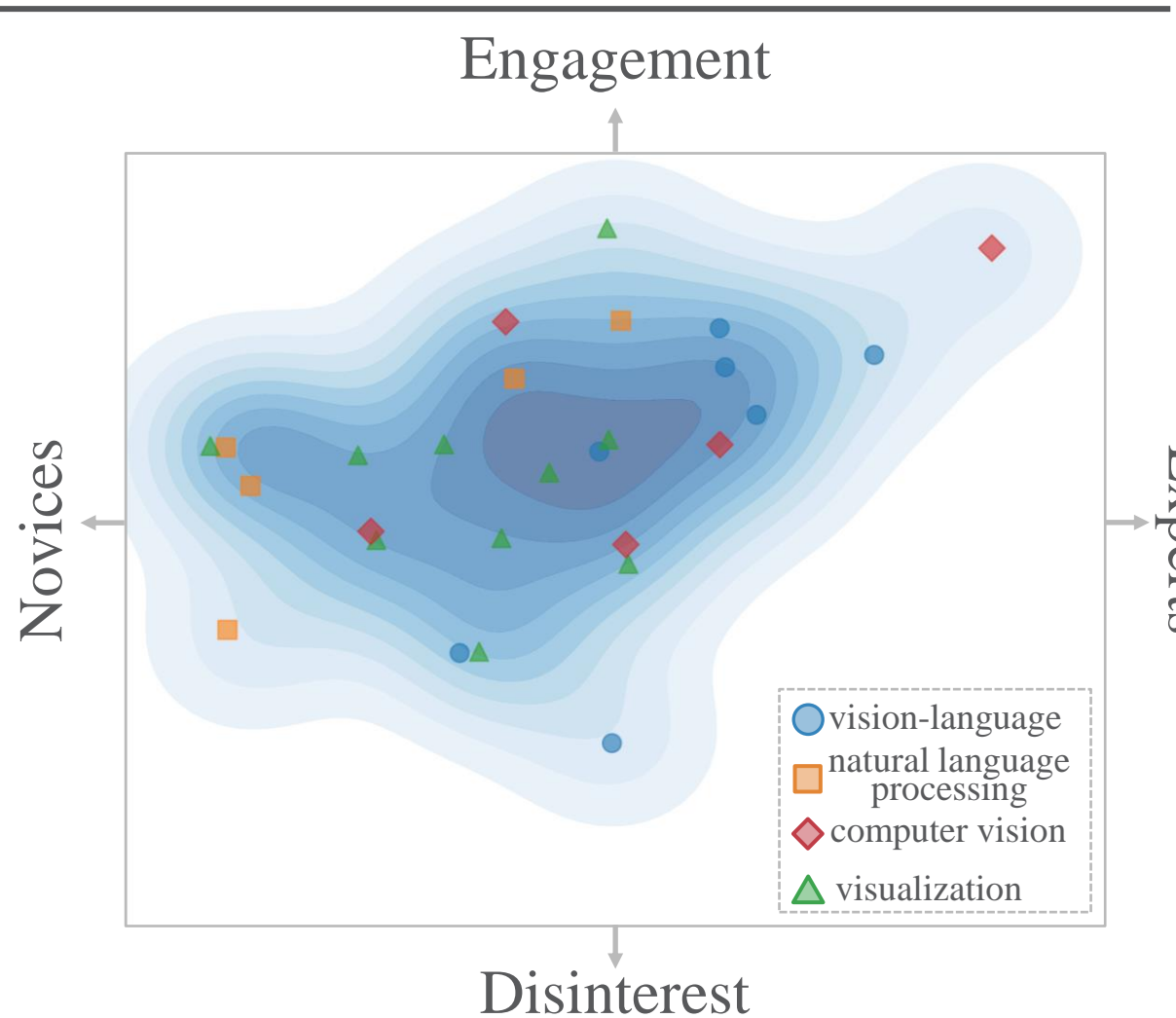
*x dimension* / *y dimension*

## 04 Evaluation

The table compares self-organizing map (SOM), resource-controlled self-organizing map (RC-SOM), and HERC-SOM (ours) in terms of quantization error (QE), topological error (TE), boundary entropy (BE), and neuron activation ratio (NA).

| Methods | Indicators | | | |
| --- | --- | --- | --- | --- |
| | QE. | TE. | BE. | NA. |
| SOM | 3.899 | 0.032 | 1.113 | 0.310 |
| RC-SOM | 3.969 | 0.040 | 1.256 | 1.0 |
| HERC-SOM(ours) | 4.027 | 0.059 | 0.918 | 1.0 |

Our work sparks discussions among participants from diverse backgrounds. Approximately 90% of participants (24 and 3) stated that, compared to methods relying solely on statistics metrics and closed dataset, CompoVis offers a more innovative and effective approach for investigating modality gaps in VLMs.



Engagement / Disinterest / Novices / Experts

vision-language
natural language processing
computer vision
visualization

## 03 Visualization Insights



to the right of / to the left of / to the right of

baseline / permutation
pre-trained / fine-tuned



data size: 0.9k    data size: 2.5k    data size: 3.6k

TongLi97.github.io    litong@zjut.edu.cn    ZJUT Vis, Zhejiang University of Technology

浙江工业大学 ZHEJIANG UNIVERSITY OF TECHNOLOGY