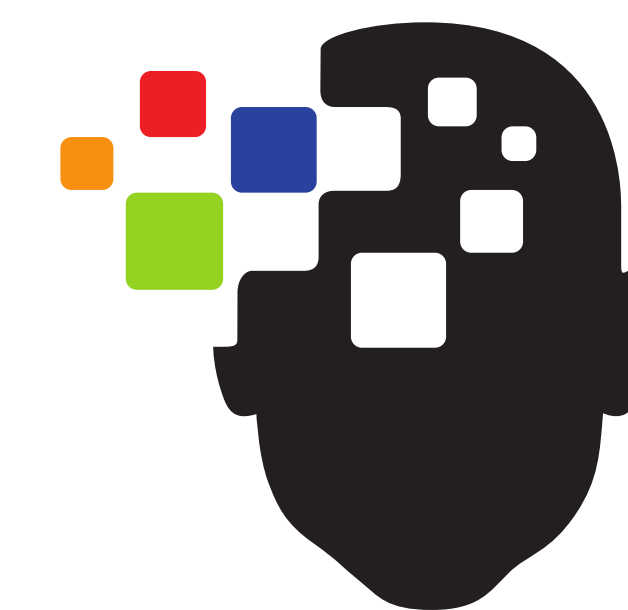# R2-T2: Re-Routing in Test-Time for Multimodal Mixture-of-Experts

Zhongyang Li[1], Ziyue Li[2], Tianyi Zhou[2]

Johns Hopkins University[1], University of Maryland, College Park[2]

## ◆ Introduction

R2-T2 improves expert selection in multimodal Mixture-of-Experts models by locally optimizing routing weights at test time — using nearby successful examples and without changing any model parameters. Our key contributions are:

➤ **Test-Time Re-Routing Framework**: We formalize adjusting routing outputs at inference via reference examples.

➤ **Three Optimization Methods**: We propose Neighborhood Gradient Descent, Kernel Regression, and Mode Finding for per-input weight optimization.

➤ **Significant Performance Gains**: We demonstrate consistent, significant gains across eight benchmarks, nearing oracle performance without any retraining.

## ◆ Reference Set & Experts

Our reference set spans three tasks—visual understanding, reasoning, and OCR—and uses six experts:

$I_{AUX}$ cross-attends visual features to structured CV outputs

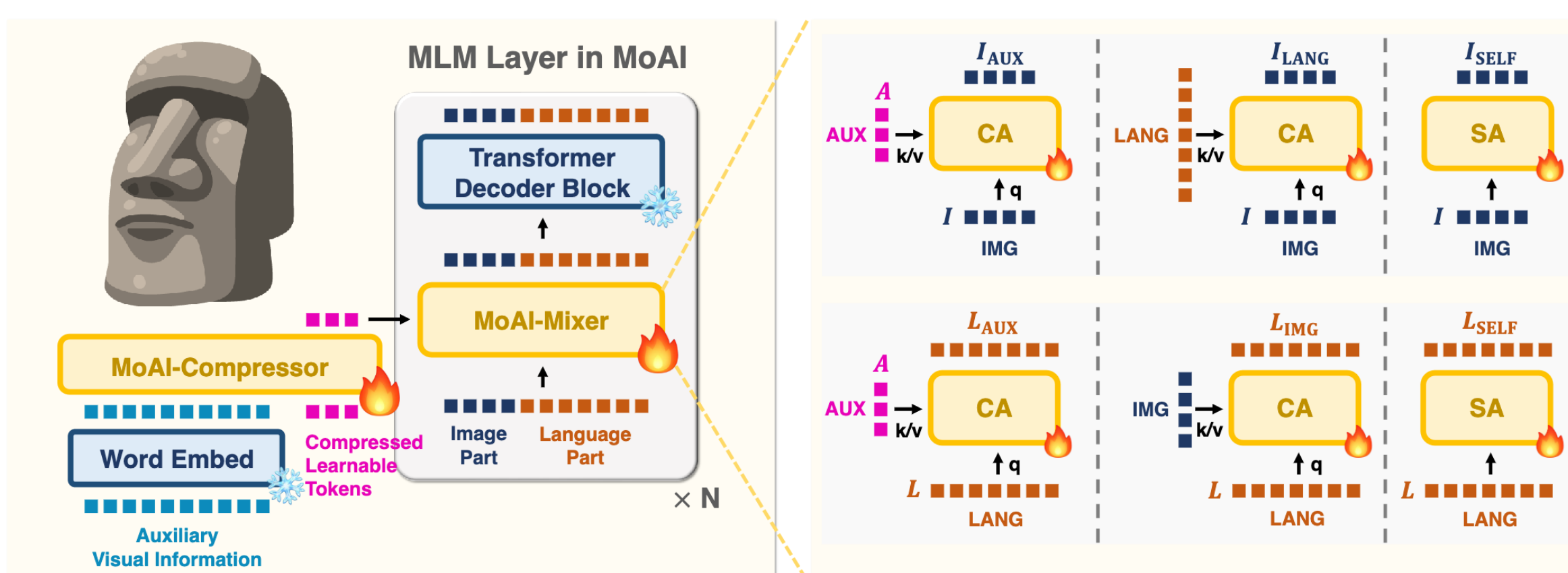$I_{LANG}$ aligns visual features with language semantics

$I_{SELF}$ preserves spatial detail via self-attention

$L_{AUX}$ integrates CV outputs into language understanding
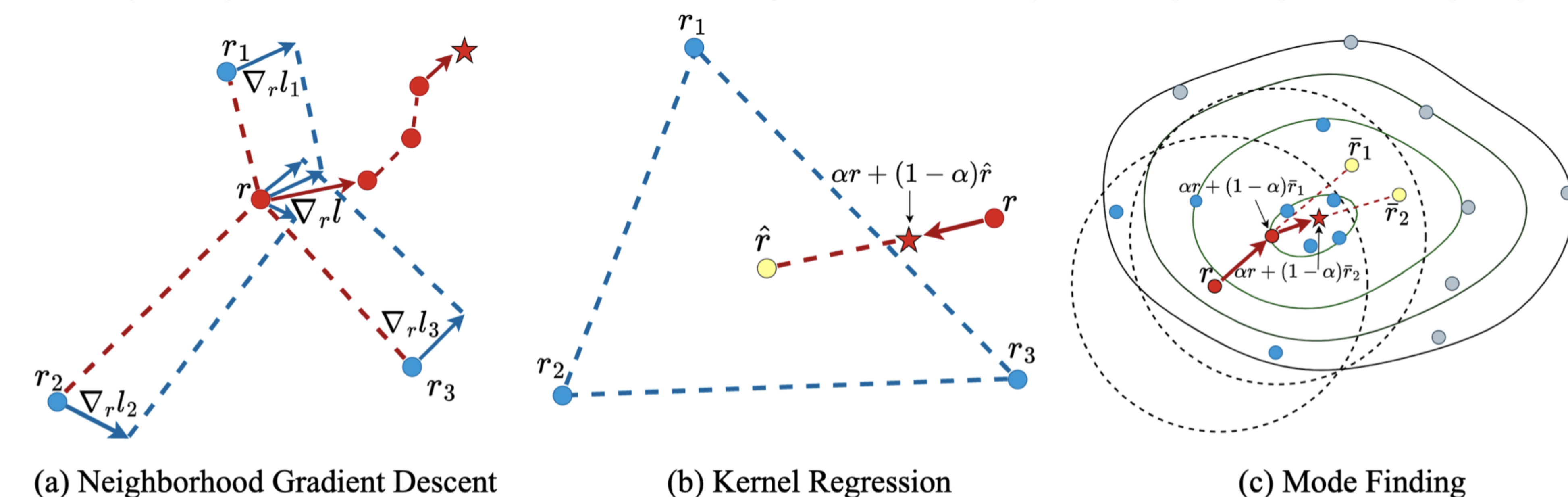
$L_{IMG}$ grounds language in visual context

$L_{SELF}$ ensures coherent text generation.

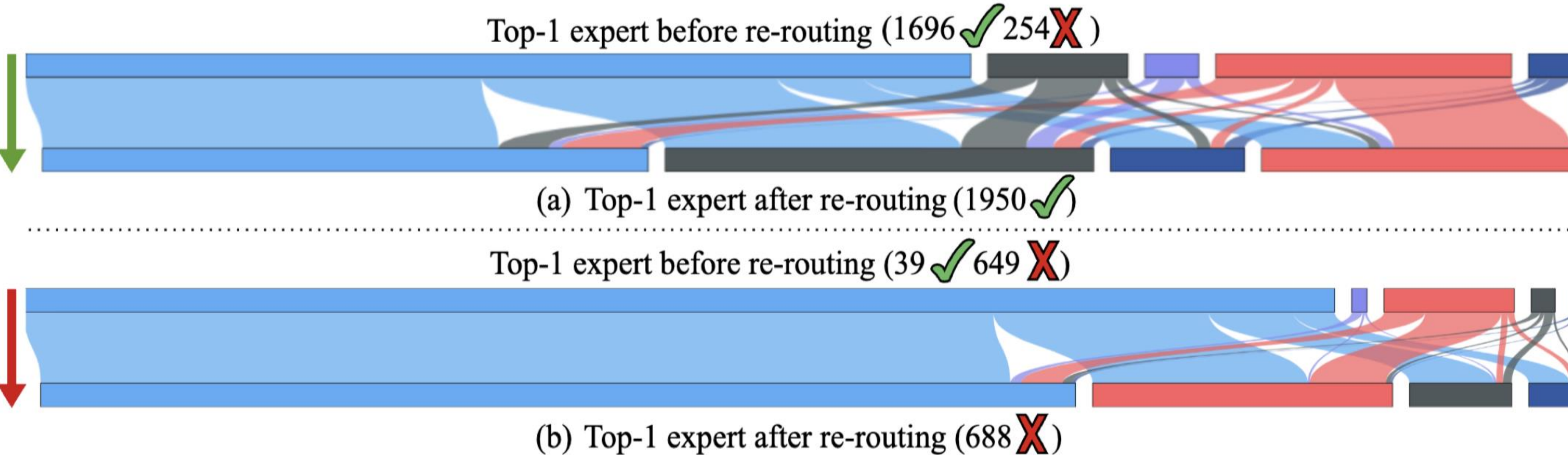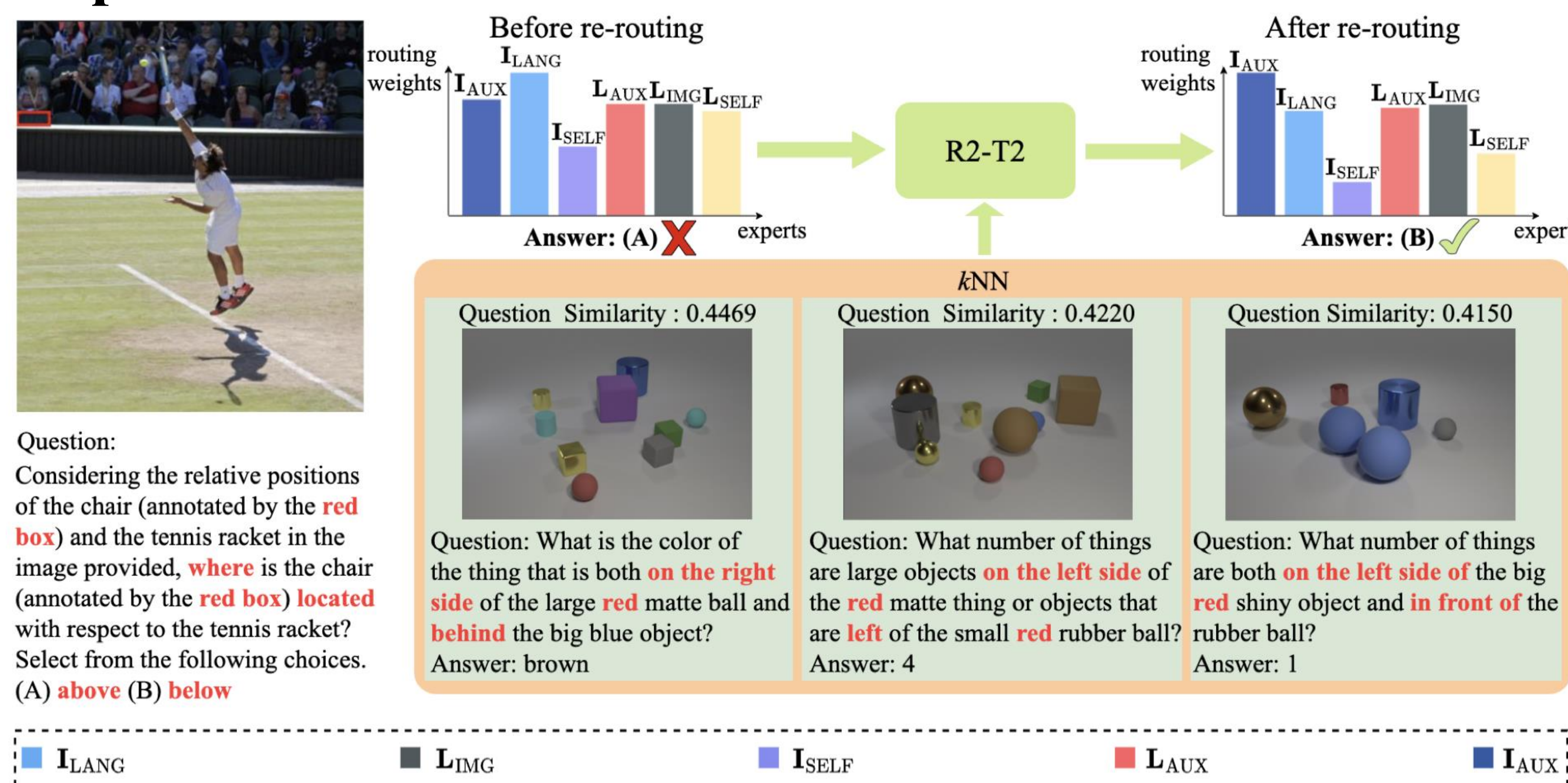| Task Type | Reference | Size | Evaluation | Size |
|---|---|---|---|---|
| General Visual Understanding | VQA-V2 | 5,000 | MMBench | 2,374 |
| | Visual7W | 5,000 | MME-P | 2,114 |
| | COCO-QA | 5,000 | CVBench$^{2D/3D}$ | 2,638 |
| | CLEVR | 5,000 | GQA | 1,590 |
| Knowledge-Based Reasoning | A-OKVQA | 5,000 | SQA-IMG | 2,017 |
| | TQA | 5,000 | AI2D | 3,087 |
| | MathVista | 5,000 | | |
| Optical Character Recognition | ST-VQA | 5,000 | TextVQA | 5,734 |
| | DocVQA | 5,000 | | |



## ◆ Method: Test-Time Re-Routing



● Routing weights of neighbors    ● Routing weights of the test sample in re-routing    ★ Routing weights after re-routing
→ Neighbors' gradient descent direction    → Re-routing direction    ○ Weighted average of neighbors' routing weights

(a) Neighborhood Gradient Descent    (b) Kernel Regression    (c) Mode Finding
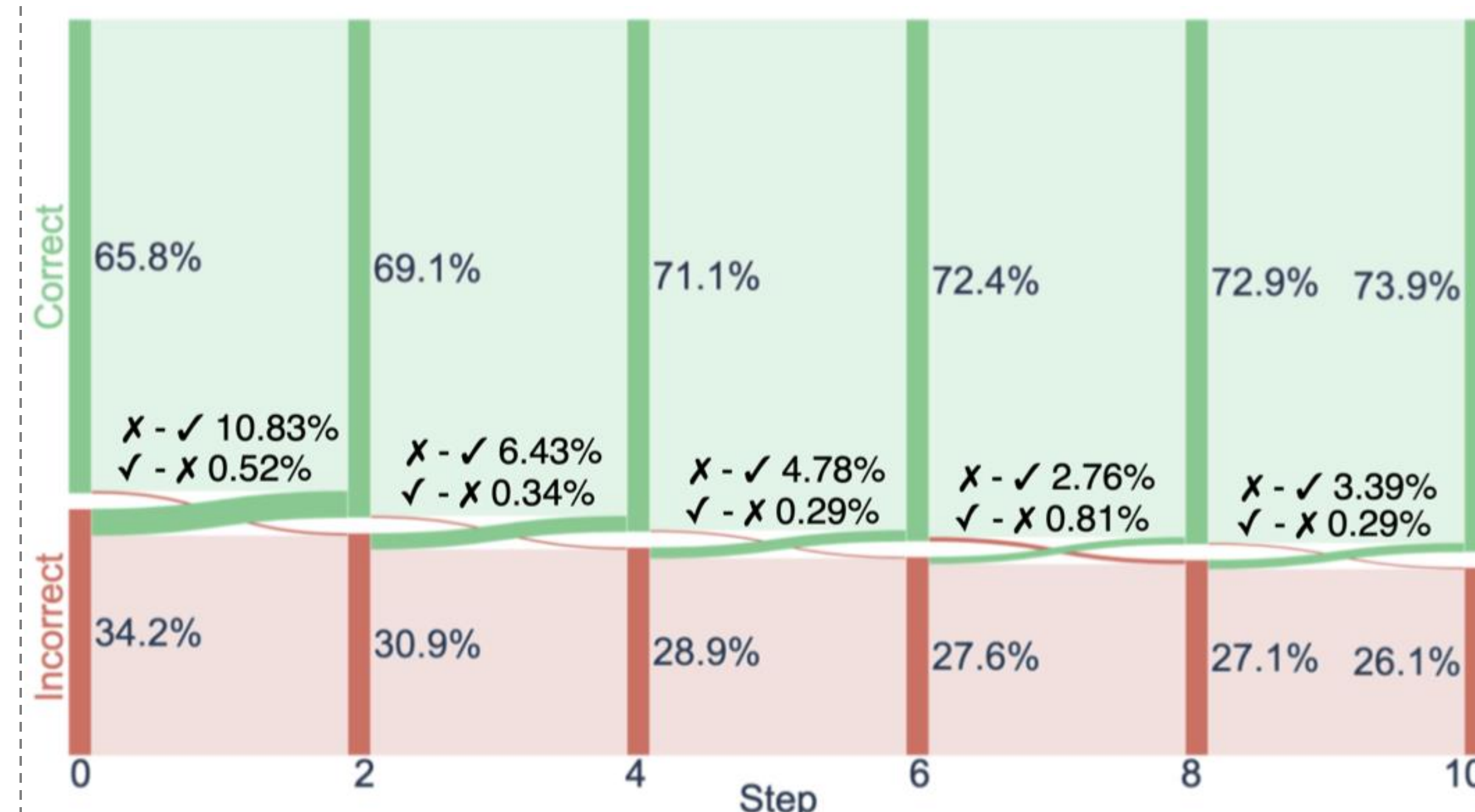
➤ **Neighborhood Gradient Descent** estimates the gradient of $r$ using the loss function of the nearest neighbors in reference set and take gradient steps on $r$ to minimize this loss.

➤ **Kernel Regression** computes a kernel-weighted average $\hat{r}$ of neighbors' routing vectors. Then interpolate between the original $r$ and $\hat{r}$, using the binary search to find $\alpha$ that maximizes model confidence.

➤ **Mode Finding** identifies the high-density "mode" of neighbors' routings via a mean-shift update in the routing-weight space. Iteratively move $r$ toward this dense region.

## ◆ Expert Shift Patterns



Question:
Considering the relative positions of the chair (annotated by the **red box**) and the tennis racket in the image provided, **where** is the chair (annotated by the **red box**) **located** with respect to the tennis racket? Select from the following choices. (A) **above** (B) **below**

Question Similarity : 0.4469
Question: What is the color of the thing that is both **on the right side** of the large **red** matte ball and **behind** the big blue object?
Answer: brown

Question Similarity: 0.4220
Question: What number of things are large objects **on the left side** of the **red** matte thing or objects that are **left** of the small **red** rubber ball?
Answer: 4

Question Similarity: 0.4150
Question: What number of things are both **on the left side** of the big **red** shiny object and **in front of** the rubber ball?
Answer: 1

$I_{LANG}$    $L_{IMG}$    $I_{SELF}$    $L_{AUX}$    $I_{AUX}$

Top-1 expert before re-routing (1696 ✓ 254 ✗)
(a) Top-1 expert after re-routing (1950 ✓)

Top-1 expert before re-routing (39 ✓ 649 ✗)
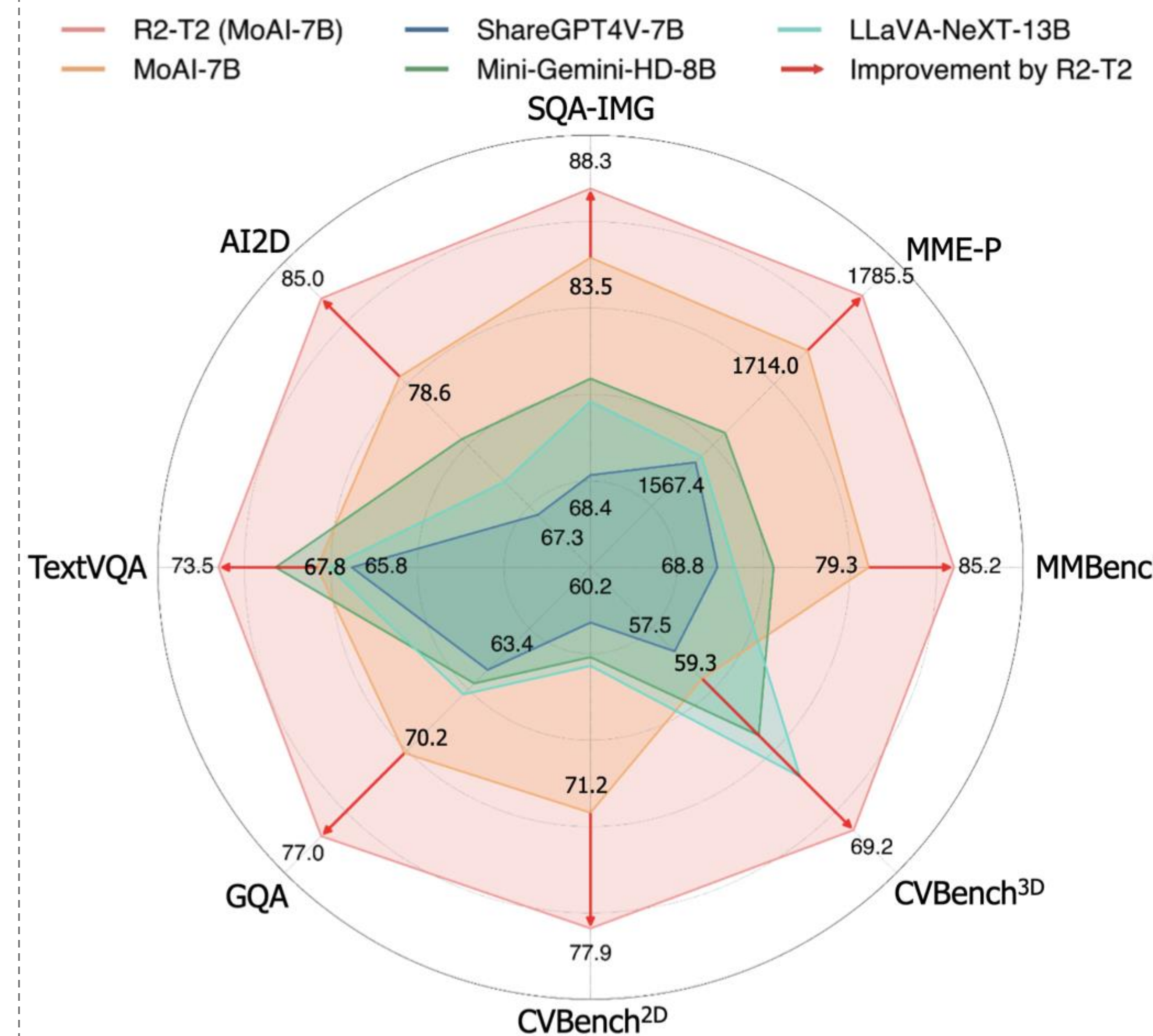(b) Top-1 expert after re-routing (688 ✗)

## ◆ Accuracy Transition Analysis



This figure illustrates the transition of predictions as NGD progresses over ten steps. During Step 0 to Step 10, a total of 28.12% of incorrect predictions have been converted to correct ones.

## ◆ Results

R2-T2 applied to MoAI-7B compared against 7/8/13B VLMs on 8 benchmarks, surpassing a recent 13B VLM.



R2-T2 (MoAI-7B)    ShareGPT4V-7B    LLaVA-NeXT-13B
MoAI-7B    Mini-Gemini-HD-8B    Improvement by R2-T2

| Method | MMBench | MME-P | SQA-IMG | AI2D | TextVQA | GQA | CVBench$^{2D}$ | CVBench$^{3D}$ |
|---|---|---|---|---|---|---|---|---|
| MoVA (base model) | 74.3 | 1579.2 | 74.4 | 74.9 | 76.4 | 64.8 | 61.6 | 62.3 |
| Mode Finding | 75.2 | 1587.1 | 74.9 | 75.8 | 77.3 | 65.7 | 62.5 | 63.2 |
| Kernel Regression | 77.9 | 1610.6 | 76.4 | 78.5 | 79.9 | 68.3 | 65.2 | 65.9 |
| NGD | 81.2 | 1645.3 | 79.1 | 82.4 | 81.8 | 83.2 | 71.5 | 68.9 |
| Oracle (upper bound) | 87.6 | 1735.4 | 87.3 | 88.4 | 89.5 | 76.2 | 72.5 | 73.2 |
| MoAI (base model) | 79.3 | 1714.0 | 83.5 | 78.6 | 67.8 | 70.2 | 71.2 | 59.3 |
| Mode Finding | 80.8 | 1725.2 | 84.1 | 79.8 | 66.5 | 71.4 | 70.0 | 60.1 |
| Kernel Regression | 83.7 | 1756.7 | 86.2 | 82.6 | 71.2 | 74.5 | 74.6 | 64.5 |
| NGD | 85.2 | 1785.5 | 88.3 | 85.0 | 73.5 | 77.0 | 77.9 | 69.2 |
| Oracle (upper bound) | 92.1 | 1860.2 | 93.8 | 91.2 | 79.6 | 84.0 | 84.0 | 76.8 |