

Chart-G1: Visually Grounded Chart Reasoning by Rewarding Multimodal Large Language Models

Under Review

Xingchen Zeng, Yuanbang Liu, Jianing Hao, Wei Zeng

Quick look of Chart-G1

GRPO-based Visual Grounding

- Pipeline Illustration

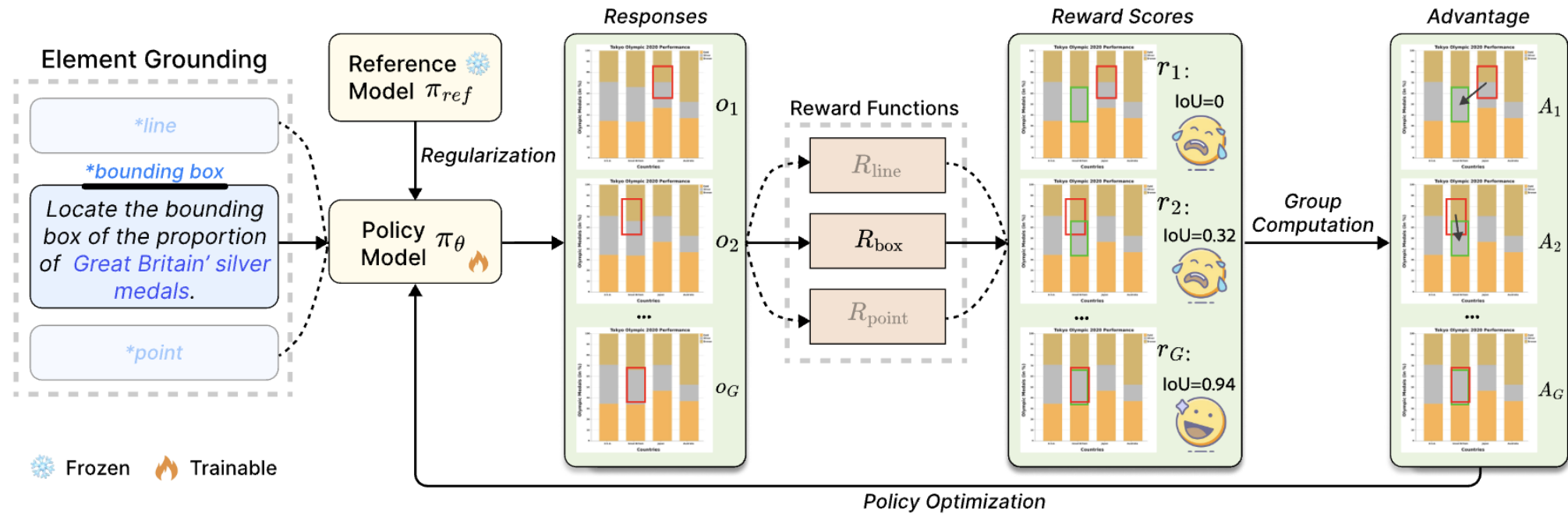
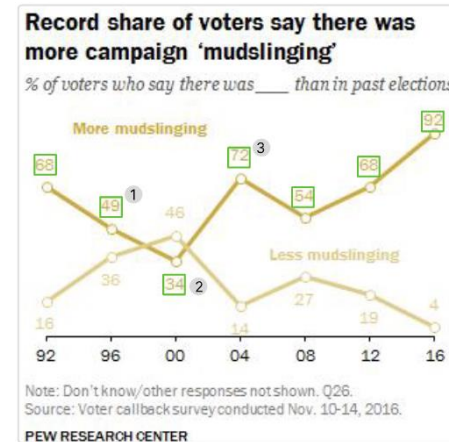


Fig. 4: Illustration of the GRPO framework, instantiated with an element grounding task. Given an instruction (e.g., "Locate the bounding box..."), the policy model generates a group of candidate responses. A task-specific reward function (R_{box}) then quantifies the quality of each response. These reward scores are used to compute advantage values for the group. Finally, the advantages guide the policy optimization step, updating the model's parameters. Throughout this process, a reference model regularizes the updates to stabilize training and maintain the model's generalizability.

Quick look of Chart-G1

Dual-phase GRPO Training

- Dual-phase
 - **Element grounding** rewards the MLLM to learn to perceive accurately, ensuring the faithful capture of visual information;
 - **Grounded reasoning** rewards the MLLM to integrate coordinates as vision anchors into their reasoning chain to perform diverse chart reasoning tasks, thus improving the reasoning reliability and interpretability.



What is the average of "more mudslinging"?

Chart-G1: <think>\nTo calculate the average of 'more mudslinging', I need to sum up all the values for 'more mudslinging' and then divide by the number of data points.\n\nThe values for 'more mudslinging' are:\n<bbox>[25,189,35,201]</bbox>: 68\n<bbox>[87,232,93,243]</bbox>: 49\n<bbox>[153,255,161,265]</bbox>: 46\n<bbox>[217,180,226,192]</bbox>: 72\n<bbox>[285,216,293,228]</bbox>: 54\n<bbox>[352,189,360,200]</bbox>: 68\n<bbox>[419,135,426,146]</bbox>: 92\n\nSumming these values: 68 + 49 + 34 + 72 + 54 + 68 + 92 = 439\n\nThere are 7 data points, so the average is 432 / 7 = 62.42\n\n</think>\n<boxed>{62.42}</boxed>

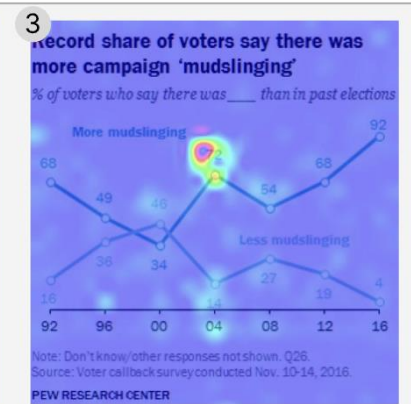
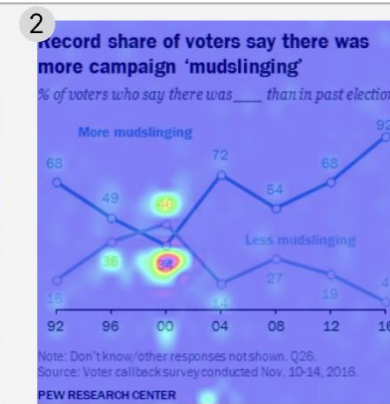
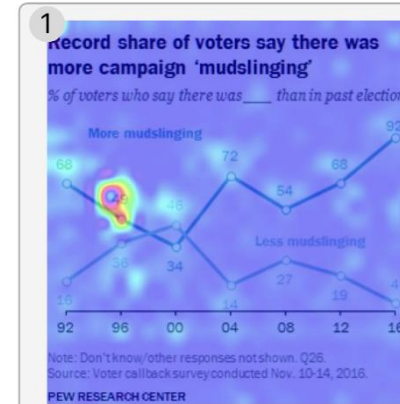


Fig. 6: Illustration of the visual grounded chart reasoning pattern that emerges during the grounded reasoning rewarding phase.



Quick look of Chart-G1

GRPO-based Visual Grounding

- Phase1 *Element Grounding*
 - Three different grounding formats, including bboxes, points, and lines.

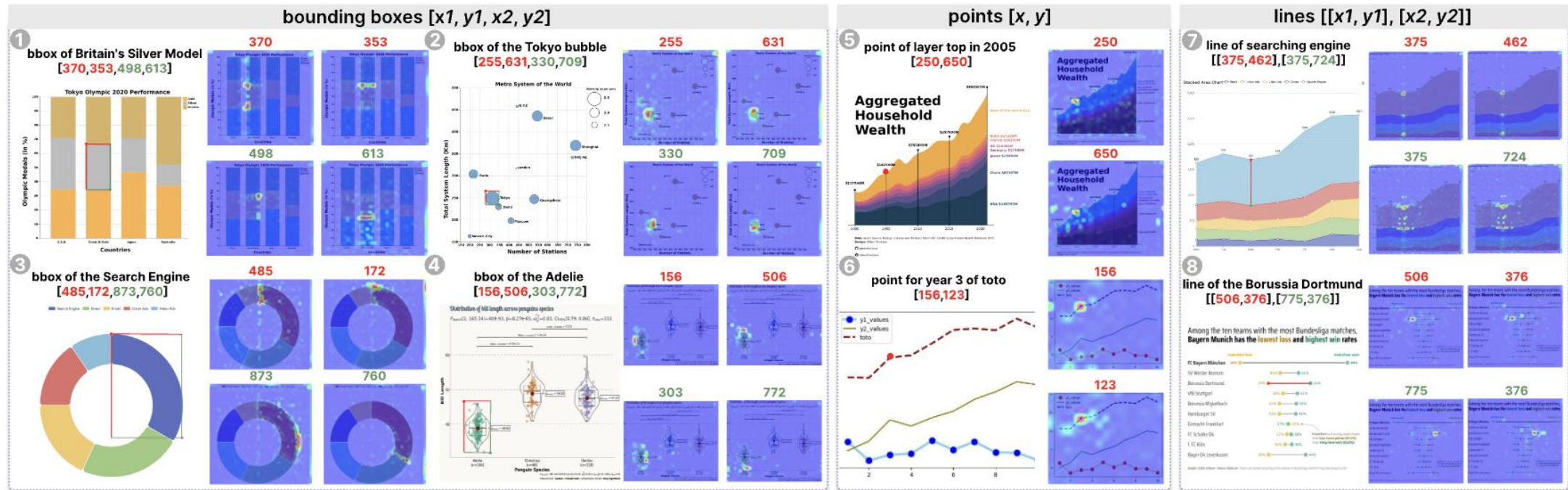


Fig. 5: Visualizations of Chart-G1's visual attention distributions when grounding chart elements with three different formats: bboxes, points, and lines (from left to right). Chart-G1 exhibits highly precise visual attention across diverse geometric representations of chart elements.



Quick look of Chart-G1

Data Engine for Coordinates Annotation

- SVG-as-the-intermediary Data Engine

