

Scene Understanding Maps: Predicting Most Frequently Fixated Object during Scene Description with Multi-Modal Large Language Models <u>Shravan Murlidaran (smurlidaran@ucsb.edu)</u>, <u>Ziqi Wen², Jonathan Skaza³, and Miguel P. Eckstein¹</u>

Introduction

- Over the last decade, image computable models such as saliency (GBVS^[1]) or Deepgaze^[2] have been developed to predict locations people might fixate when viewing naturalistic scenes.
- A recent technique (Scene Understanding Map (SUM)) developed by Murlidaran et al.^[3] showed that people fixate on objects critical to understanding the scene. They digitally removed objects from a scene and quantified the impact of the removal by comparing the similarity of descriptions from humans with and without the object's presence.
- Here, we build upon their method by automating various components of their procedure. To this end, we use the Winograd images developed by Murlidaran et al. to compare the effect of automating various components on the agreement of the most critical object with the original method. We also compare the ability of all the methods to predict the most fixated object on a scene description task.

References

- Harel, J., Koch, C. and Perona, P., 2006. Graph-based visual saliency. Advances in neural information processing systems, 19.
- Kümmerer, M., Wallis, T.S. and Bethge, M., 2016. DeepGaze II: Reading fixations from deep features trained on object recognition. arXiv preprint arXiv:1610.01563.
- Murlidaran, S., & Eckstein, M. P. (2024). Eye Movements during Free Viewing Maximize Scene Understanding. Journal of Vision, 24(10), 1189-1189.
- Hurst, A., Lerer, A., Goucher, A. P., Perelman, A., Ramesh, A., Clark, A., ... & Kivlichan, I. (2024). Gpt-4o system card. arXiv preprint arXiv:2410.21276.
- Liu, S., Zeng, Z., Ren, T., Li, F., Zhang, H., Yang, J., ... & Zhang, L. (2023). Grounding dino: marrying dino with grounded pre-training for open-set object detection. arXiv abs/2303.05499 (2023).
- Ren, T., Liu, S., Zeng, A., Lin, J., Li, K., Cao, H., ... & Zhang, L. (2024). Grounded sam: Assembling open-world models for diverse visual tasks. *arXiv preprint arXiv:2401.14159*.
- Zhuang, Junhao, et al. "A task is worth one word: Learning with task prompts for highquality versatile image inpainting." European Conference on Computer Vision. Cham: Springer Nature Switzerland, 2024.













LLM - Large Language Model

¹Dept. of Psychological and Brain Sciences, ²Dept. of Computer Science, ³Dept. of Dynamical Neuroscience University of California, Santa Barbara

MLLM - Multi-Modal LLM

Heatmap Correlations across Winograd Pairs





Limitations for generalization

Deleting very large objects can induce artifacts that can impact the model's understanding of the scene, though it may not be a critical object



The person appears to be taking a break after exercising with dumbbells in a living room setting.

The model's understanding of the scene may not be similar to how humans understand the scene



Human Description: A man is bending to look for a tennis ball under sofa.

MLLM Description: A person is bending over, possibly looking at or interacting with something on a table in a living room.

- description task.

VIU vision o image understanding

Results









A person is examining a small diorama on the floor in a room that also contains a keyboard and guitar.

4. Conclusion

Scene Understanding Maps perform well at predicting the most fixated object during a scene

MLLMs can be used to automate the generation of Scene Understanding Maps. Although they have certain limitations, their performance is still at par with the original SUM maps