

Foveated Multi-Modal Large Language Model Maps to Predict Time to Understand Scenes

Motivation & Introduction

- The rapid improvement of Vision Language Models (VLM) makes the idea of image-computable model for predicting Response Time (RT) for understanding a scene possible
- ► VLM can generate scene descriptions for any image input
- Similarity metrics based on language embeddings allow quantitative evaluations of the similarity of a human description of a scene to a gold standard description.
- ▶ What cause the main bottleneck in human scene understanding?
- Hypothesis: The interaction between the foveated nature of the human visual system and the spatial distribution across the image of the visual information is critical to comprehending the scene.



A baseball player gets read to hit a ball as another gets ready to catch it

Mean RT: 2.01s



QQQ There is a dangerous looking man holding a knife and a man and a woman standing across the table from him trying to deescalate the situation.

Mean RT: 4.50s Mean #Saccades: 12.67 Difficult!

Mean #Saccades: 5.50 Easy! Figure 1: Comparison between low-effort and high-effort scene understanding.



Figure 2: Overview of human psychophysics

Psychophysics

► (1) Response Time Study:

- Participants (N = 17) will do the free viewing until they understand the scene, and then type in the scene descriptions. Response time and collected eye-tracking data, including the total number of saccades, are collected.
- ► (2) Saccade-Limited Presentation Study: Each scene was displayed for a restricted number of eye movements, either 2 or 4 saccades. Participants (N = 16) were then instructed to provide descriptions of the scene based on what they observed.

Ziqi Wen, Jonathan Skaza, Shravan Murlidaran, William Wang, Miguel Eckstein University of California, Santa Barbara



experiments). S will be further normalized to range from 0 to 1 and let a higher value means the scene is more difficult to understand.



Results

As shown in Fig.5: The F-SUM Score significantly outperforms all baseline metrics in predicting both response time and saccade

> count (p < .005). As shown in Fig.6: The F-SUM Score significantly (p < .0001) outperforms all baseline metrics in predicting participants' ability to extract and articulate the gist of a scene when access to visual information was limited by spatial viewing constraints.

Image complexity can be biased by the sheer number of visual elements—such as the densely packed carrots (top-right example in Fig.7) in the upper-left image—which may not

standard under limited saccades

> Language entropy captures the diversity of sampled descriptions but may miss the underlying effort required to generate them. In the bottom-left example in Fig.7, captions consistently mention a person reaching into a tree while others look on. Although the descriptions seem coherent and varied, they require multiple eye movements to integrate

Acknowledgements

Reference

[1] Tinglei Feng, Yingjie Zhai, Jufeng Yang, Jie Liang, Deng-Ping Fan, Jing Zhang, Ling Shao, and Dacheng Tao. Ic9600: A benchmark dataset for automatic image complexity assessment.

[2] Louis Mahon and Thomas Lukasiewicz. Minimum description length clustering to measure meaningful image complexity. Pattern Recognition, 145:109889, 2024.

[4] Jeffrey S Perry and Wilson S Geisler. Gaze-contingent real-time simulation of arbitrary visual fields. In Human vision and electronic imaging VII, volume 4662, pages 57-69. SPIE, 2002. [5] Brian D Ripley. The second-order analysis of stationary point processes. Journal of applied probability, 13(2):255–266, 1976.

[6] Ruth Rosenholtz, Yuanzhen Li, Jonathan Mansfield, and Zhenlan Jin. Feature congestion: a measure of display clutter. In Proceedings of the SIGCHI conference on Human factors in computing systems, pages 761–770, 2005.

[7] Ruth Rosenholtz, Yuanzhen Li, and Lisa Nakano. Measuring visual clutter. Journal of vision, 7(2):17-17, 2007.

