Conceptual priorities shape individual gaze patterns during 1 naturalistic visual attention 2 3 Amanda J. Haskins (ajhaskins@ucsd.edu), Katherine O. Packard 4 (kateopackard@gmail.com), Caroline E. Robertson (cerw@dartmouth.edu) Psychological & Brain Sciences, Dartmouth College, 3 Maynard Street Hanover, NH 03755 USA 6 7

Abstract

10 In daily life, we readily recognize people, places, and 11 objects (e.g., "soldier," "stadium," "flag"), as well as 12 the conceptual links between them (e.g., 13 "patriotism"). Here, we show that conceptual-level 14 information shapes individuals' gaze patterns in a 15 naturalistic eye tracking paradigm. Participants (N = 16 61) explored a large set of real-world photospheres 17 (N = 100) in headmounted VR while their gaze was 18 continuously monitored using in-headset eye 19 tracking. To assess the informational priorities 20 guiding individual differences in gaze, we leveraged 21 the embedding spaces of large vision and language 22 models. We found that individually-specific gaze 23 patterns across diverse real-world photospheres can 24 be captured by a large language model (LLM) that 25 encodes abstract relationships beyond the visual 26 image content. We demonstrate that the embedding 27 spaces of language and vision models explain 28 unique variance in gaze behavior, and that 29 LLM-based models capture individually specific 30 attentional priorities. These results highlight a new 31 dimension of human selective attention: namely, the 32 influence of individuals' unique conceptual-level 33 information seeking priorities.

34

37

5

8

9

Keywords: individual differences, naturalistic 35

visual attention, information seeking, concepts 36

Introduction

38 What guides individuals' selective attention when 39 viewing real-world scenes? Understanding gaze 40 behavior has been a central goal in psychology, with 41 focus on the roles of "bottom-up" (visual salience) 42 and "top-down" (semantic meaning) scene features 43 (Henderson & Haves, 2017). However, this 44 dichotomy overlooks the influence of contextual 45 factors on gaze behavior (Awh et al., 2012). While 46 external contexts like tasks and rewards are known 47 to shape attention (Borji & Itti, 2014; Tong et al., 48 2017), less is understood about how intrinsic factors, 49 like an individual's unique conceptual knowledge 50 and interests, guide their gaze. This gap is important 51 for current theories that forward gaze behavior as a 52 proactive, information-seeking process (Haskins et 53 al., 2020; Hayhoe, 2017). Individuals explore scenes 54 in the real world proactively: they use their 55 conceptual knowledge to form and test hypotheses 56 about object relationships. Here, we explored how 57 conceptual information beyond the visual domain 58 influences gaze patterns during scene viewing. We 59 hypothesized that gaze patterns reflect, in part, 60 individuals' stable, conceptual priorities. To test this, 61 gaze behavior was analyzed from adult participants 62 who explored real-world photospheres in virtual 63 reality (VR). We developed a computational 64 approach using a large language model (LLM; 65 (Devlin et al., 2019)) to approximate conceptual 66 information in scenes and its role in shaping gaze 67 patterns. In brief, we found that individual gaze 68 patterns are strongly predicted by the LLM's 69 conceptual feature space, which accounted for 70 significant unique variance beyond well-established 71 predictors of gaze (i.e., visual, motor features).

Methods

72

73 Participants. 66 adults (n = 36 female; mean age 74 20.67 years) participated in this study. Following 75 exclusion, data were analyzed from 61 participants. 76 Stimuli. We used 100 diverse, information-rich 77 photospheres, selected based on pilot data showing 78 consistent attention to both social (e.g., people) and 79 nonsocial elements. On each 16-second trial. 80 participants were instructed to "look around each 81 scene naturally, like you would in daily life."



82 83

Figure 1: A) LLM embeddings capture unique variance in gaze patterns, relative to control feature spaces. B) LLM gaze models can be used to generate individually-specific gaze predictions on left-out scenes.

84 85

Gaze data. Duration-weighted fixation density maps 86 were individually plotted for each subject and scene. 87 88 Modeling conceptual information using an LLM. 89 We used the embedding space of an LLM to 90 characterize the conceptual-level information in each 91 photosphere. First, we divided scenes into densely 92 sampled, overlapping "tiles" that were captioned by 93 independent human raters; then, we transformed 94 tiles into language model embeddings. As a control 95 model to the LLM, we used a visual transformer 96 model (ViT; (Dosovitskiy et al., 2020)) trained on 97 image classification, to model the visual content 98 depicted at each image tile. As a second control 99 model to the LLM, we used the tile's equirectangular 100 spatial coordinates (X,Y). These control feature 101 spaces were used to test whether the LLM 102 embedding space captured unique variance in gaze 103 patterns, beyond control spaces.

104

Results

105 LLM embeddings explain unique gaze variance. 106 We first asked whether LLM embeddings capture 107 information in gaze patterns that is distinct from 108 established predictors: individual spatial and 109 visual-level feature biases (e.g., center bias, object 110 biases). We performed a variance partitioning 111 analysis on an "omni gaze model" that included LLM 112 features alongside the two control feature spaces. 113 This analysis revealed that LLM features explained 114 significant variance beyond control feature sets 115 (adjusted R², all-features model: M = 0.14, SD =116 0.02; adjusted R², non-language features model: M 117 = 0.11, SD = 0.01; t(60) = 49.6, p < 0.001; Fig 1A).

118 LLM embeddings predict unique gaze patterns. 119 Next, we built an "LLM gaze model" for each 120 participant by using L2-regularized linear regression 121 to relate their gaze distribution to only the LLM 122 feature space, iteratively training on N-1 scenes, and 123 generating a predicted gaze map for each left-out 124 scene. Then, to test whether LLM gaze models were 125 individually specific. we used the same 126 leave-one-out approach to iteratively assess the 127 accuracy of an individual's own LLM gaze model 128 (correlation, predicted vs. actual map), as compared 129 with the LLM gaze models of all other individuals. 130 Specifically, we computed an "own-other difference 131 score" for each participant - subtracting the average 132 accuracy gain (i.e., improvement over the accuracy 133 of a generic baseline model) of others' models from 134 the accuracy gain of the participant's own model. 135 This difference was significant (own: M = 0.07, SD =136 0.03; other: M = 0.03, SD = 0.02, t(60) = 11.22, p < 0.02137 0.001; Fig. 1B), demonstrating clear accuracy gains 138 for one's own LLM gaze model prediction vs. other 139 participants' models in almost every participant.

140 Discussion

141 Our results show that gaze patterns offer insight into
142 the conceptual priorities different individuals bring to
143 bear while exploring their visual environment.
144 Specifically, we find that the embedding space of an
145 LLM can be used to capture variance both within
146 individual gaze patterns (variance partitioning) and
147 across individuals (own-other difference). Overall,
148 we find that each individual's gaze reveals their
149 personal conceptual priorities during scene viewing.

Acknowledgments

151 This work was supported by grants from the Nancy
152 Lurie Marks Family Foundation and the Neukom
153 Institute for Computational Science (C.E.R). A.J.H.
154 was supported by a grant from the National Institute
155 of Mental Health (1F99NS135812).

156

150

157 References

- 158 Awh, E., Belopolsky, A. V., & Theeuwes, J. (2012).
- 159 Top-down versus bottom-up attentional control:
- a failed theoretical dichotomy. *Trends in*
- 161 *Cognitive Sciences*, *16*(8), 437–443.
- 162 Borji, A., & Itti, L. (2014). Defending Yarbus: Eye

movements reveal observers' task. *Journal ofVision*, 14(3), 29–29.

- 165 Devlin, J., Chang, M. W., Lee, K., & Toutanova, K.
- 166 (2019). BERT: Pre-training of deep bidirectional
- transformers for language understanding.
- 168 NAACL HLT 2019 2019 Conference of the
- 169 North American Chapter of the Association for
- 170 Computational Linguistics: Human Language
- 171 Technologies Proceedings of the Conference,
- **172** *1*, 4171–4186.
- 173 Dosovitskiy, A., Beyer, L., Kolesnikov, A.,
- 174 Weissenborn, D., Zhai, X., Unterthiner, T.,
- 175 Dehghani, M., Minderer, M., Heigold, G., Gelly,
- 176 S., Uszkoreit, J., & Houlsby, N. (2020). An
- image is worth 16x16 words: Transformers for
- image recognition at scale. In *arXiv* [cs.CV].
- 179 Haskins, A. J., Mentch, J., Botch, T. L., & Robertson,
- 180 C. E. (2020). Active vision in immersive, 360°
- 181 real-world environments. *Scientific Reports*,
- **182** *10*(1), 14304.
- 183 Hayhoe, M. M. (2017). Vision and action. Annual
- 184 *Review of Vision Science*, *3*(1), 389–413.
- 185 Henderson, J. M., & Hayes, T. R. (2017).
- 186 Meaning-based guidance of attention in scenes
- as revealed by meaning maps. *Nature Human*
- 188 Behaviour, 1(10), 743–747.
- 189 Tong, M. H., Zohar, O., & Hayhoe, M. M. (2017).
- 190 Control of gaze while walking: Task structure,
- 191 reward, and uncertainty. *Journal of Vision*,
- **192** *17*(1), 28.