

UC Davis

UC Davis Previously Published Works

Title

Meaning guides attention during scene viewing, even when it is irrelevant

Permalink

<https://escholarship.org/uc/item/264393sh>

Journal

Attention, Perception, & Psychophysics, 81(1)

ISSN

1943-3921

Authors

Peacock, Candace E

Hayes, Taylor R

Henderson, John M

Publication Date

2019

DOI

10.3758/s13414-018-1607-7

Peer reviewed



HHS Public Access

Author manuscript

Atten Percept Psychophys. Author manuscript; available in PMC 2020 January 01.

Published in final edited form as:

Atten Percept Psychophys. 2019 January ; 81(1): 20–34. doi:10.3758/s13414-018-1607-7.

Meaning Guides Attention During Scene Viewing Even When It Is Irrelevant

Candace E. Peacock^{1,2}, Taylor R. Hayes¹, and John M. Henderson^{1,2}

¹Center for Mind and Brain, University of California, Davis

²Department of Psychology, University of California, Davis

Abstract

During real-world scene viewing, humans must prioritize scene regions for attention. What are the roles of low-level image salience and high-level semantic meaning on attentional prioritization? A previous study suggested that when salience and meaning are directly contrasted in scene memorization and preference tasks, attentional priority is assigned by meaning (Henderson & Hayes, 2017). Here we examined the role of meaning on attentional guidance using two tasks in which meaning is irrelevant and saliency is relevant: a brightness rating task and a brightness search task. Meaning was represented by meaning maps that captured the spatial distribution of semantic features. Meaning was contrasted with image salience represented by saliency maps. Critically, both maps were represented similarly, allowing us to directly compare how meaning and salience influenced the spatial distribution of attention as measured by fixation density maps. Our findings suggest that even in tasks for which meaning is irrelevant and salience is relevant, meaningful scene regions are prioritized for attention over salient scene regions. These results support theories in which scene semantics play a dominant role in attentional guidance in scenes.

Keywords

attention; scene perception; eye movements; meaning; saliency

Because we can only attend to a small portion of the visual information available to us, we have to select some regions of the visual scene for preferential analysis at the expense of others via attention. It is therefore important to understand the mechanisms by which we guide our attention through real-world scenes. A good deal of work on attentional guidance in scenes has focused on the idea that attention is driven by bottom-up, low-level image features such as color, luminance, and edge orientation that are combined into saliency maps (Borji, Parks, & Itti, 2014; Borji, Sihite, & Itti, 2013; Harel, Koch, & Perona, 2006; Itti & Koch, 2001). Saliency maps are appealing because they are computationally tractable and neurobiologically plausible (Henderson, 2007, 2017).

Correspondence: Candace E. Peacock, Center for Mind and Brain, 267 Cousteau Place, University of California, Davis, CA 95618, cepeacock@ucdavis.edu.

The authors declare no competing financial interests.

Additional Information

Additionally, these data have not previously been published.

At the same time, there is strong evidence that visual attention is influenced by cognitive factors such as the semantic informativeness of objects and entities within a scene (Antes, 1974; Henderson, 2017; Henderson, Brockmole, Castelhana, & Mack, 2007; Mackworth & Morandi, 1967), along with the viewer's task and current goal (Buswell, 1935; Hayhoe & Ballard, 2005; Hayhoe, Shrivastava, Mruzek, & Pelz, 2003; Henderson, 2007, 2017; Henderson & Hollingworth, 1999; Navalpakkam & Itti, 2005, 2007; Rothkopf, Ballard, & Hayhoe, 2016; Tatler, Hayhoe, Land, & Ballard, 2011; Yarbus, 1967). Yet much of the research on attentional guidance has continued to focus solely on image salience. One reason for the popularity of image salience is that it is relatively straightforward to compute and represent. In contrast, it has been less clear how to generate and represent the spatial distribution of semantic features across a scene. To directly compare image salience to semantic informativeness, it is necessary to represent scene meaning in a format equivalent to that of image salience.

To address this issue, we have recently introduced the concept of meaning maps as a way to represent the spatial distribution of scene semantics (Henderson & Hayes, 2017). To generate meaning maps, Henderson and Hayes (2017) used crowd-sourced responses in which naïve participants rated the meaning of image patches from real-world scenes. Specifically, photographs of scenes were divided into a dense array of objectively defined circular overlapping patches at a coarse and a fine spatial scale. These patches were then shown to raters who rated how informative or recognizable each patch was (see also Antes, 1974; Mackworth & Morandi, 1967). Finally, meaning maps of each scene were created by interpolating the ratings at each spatial scale and averaging across the two scales.

Meaning maps provide a pixel-by-pixel prediction of semantic content across a scene just as saliency maps provide a pixel-by-pixel prediction of saliency across a scene. Since meaning maps are represented in the same format as saliency maps, their predictions for visual attention can be directly compared to saliency maps using the methods that have typically been used to compare the relationship between saliency maps and attention (Carmi & Itti, 2006; Itti, Koch, & Niebur, 1998; Parkhurst, Law, & Niebur, 2002; Torralba, Oliva, Castelhana, & Henderson, 2006). In this way, meaning maps and saliency maps together provide a way to compare how meaning and salience influence visual attention during real-world scene viewing.

Henderson and Hayes (2017) investigated the degree to which meaning maps and saliency maps predicted visual attention in real-world scenes during memorization and aesthetic judgment tasks. In that study, attention maps were created based on the locations of eye fixations. The results showed that meaning maps and saliency maps were highly correlated, and both were able to predict the spatial distribution of attention in scenes. Importantly, in both tasks meaning accounted for significantly more of the variance in attention than image salience. Further, when the variance due to salience was controlled, meaning accounted for significantly more of the remaining variance in attention, but when meaning was controlled, no additional variance in attention was accounted for by salience. These results held across the entire viewing time. Henderson and Hayes (2018) replicated this pattern of results using attention maps constructed from duration-weighted fixations, and Henderson, Hayes, Rehrig, and Ferreira (2018) showed that the results extended to scene description tasks. In

total, the findings showed that meaning (rather than image salience) was the main driver of visual attention.

Although the data favoring meaning over image salience have been clear, it could be argued that the viewing tasks used to compare meaning and image salience were biased toward meaning. That is, it might be that memorization, aesthetic preference, and scene description tasks require the viewer to focus on the semantic features of scenes. If this is true, then it may be that the advantage for meaning over salience is restricted to viewing tasks that specifically require analysis of meaning. To address this hypothesis, the current study investigated whether attention continues to be guided by meaning during scene viewing even when saliency is relevant, and meaning is irrelevant to the viewer's task.

Specifically, in the present study we used two tasks that were designed to emphasize salience and eliminate the need for meaning in attentional guidance: a brightness rating task in which participants rated scenes for overall brightness, and a brightness search task in which participants counted the number of bright patches within scenes (Figure 1). Critically, these tasks were designed to make meaning task-irrelevant and salience task-relevant. If the use of meaning to guide attention is task-based, then the relationship between meaning and attention found in our earlier studies should no longer be observed in these conditions. On the other hand, if the use of meaning to guide attention during scene viewing is a fundamental property of the attention system, then we should continue to observe a relationship between meaning and attention even when only salience is relevant to the task.

Method

Eye-tracking

Participants.—Thirty University of California, Davis undergraduate students with normal or corrected-to-normal vision participated in the experiment (25 females, average age = 20.84). All participants were naïve concerning the purpose of the experiment and provided verbal consent. The eye-movement data from each participant were filtered for excessive track losses due to blinks or loss of calibration. Following Henderson & Hayes (2017), we averaged the percent signal ($[\text{number of good samples} / \text{total number of samples}] \times 100$) for each trial and participant using custom MATLAB code. The percent signal for each trial was then averaged for each subject and compared to an *a priori* 75% criterion for signal. Overall, all participants had greater than 75% signal resulting in no removed subjects.

Apparatus.—Eye movements were recorded using an EyeLink 1000+ tower mount eyetracker (spatial resolution 0.01° rms) sampling at 1000 Hz (SR Research, 2010b). Participants sat 85 cm away from a 21" monitor, so that scenes subtended approximately 26.5° x 20° of visual angle at 1024 × 768 pixels. Head movements were minimized by using a chin and forehead rest. Although viewing was binocular, eye movements were recorded from the right eye. The experiment was controlled with SR Research Experiment Builder software (SR Research, 2010a).

Stimuli.—Stimuli consisted of the 40 digitized photographs (1024 × 768 pixels) of indoor and outdoor real-world scenes. Scenes were luminance matched across the scene set by

converting the RGB image of the scene to LAB space and scaling the luminance channel of all scenes from 0 to 1. All instruction, calibration, and response screens were luminance matched to the average luminance ($M = 0.45$) of the scenes.

Procedure.—Each participant completed two scene-viewing conditions in a within-subject design: a brightness rating task and a brightness search task (Figure 1). During the brightness rating task, participants were instructed to rate the overall brightness of the scene on a scale from 1 to 6 (1 = very dark; 2 = dark; 3 = somewhat dark; 4 = somewhat bright; 5 = bright; 6 = very bright). During the brightness search task, participants were instructed to count the number of bright patches within the scene. Because the goal of this study was to assess whether we could eliminate the relationship between meaning and attention in tasks that did not require the use of meaning, we emphasized speed and accuracy during both tasks. Participants were given a maximum scene-viewing time of 12 s (as done in Henderson & Hayes, 2017), but had the option to terminate the scene and continue to the response screen earlier by pressing a key on a button box (RESPONSEPixx; VPixx Technologies, Saint-Bruno, CA). We included the early termination option so that could focus on task-relevant eye movement behavior. Following their button press or the maximum 12 s of scene presentation, participants were shown a response screen in which the number 0 was enclosed in a square (Figure 1). Then, participants used left and right buttons on the button box to respectively increase or decrease the value of the number until it matched their rating or patch count for that scene. They then pressed the center key to continue to the next scene.

Before starting the experiment, participants completed two practice trials in which they were familiarized with each condition and the button-box. After the practice trials, a 13-point calibration procedure was performed to map eye position to screen coordinates. Successful calibration required an average error of less than 0.49° and a maximum error of less than 0.99° . Presentation of each scene was preceded by a drift correction procedure, and the eye-tracker was recalibrated when the calibration was not accurate. The calibration was also repeated between task blocks.

The 40 scene stimuli were randomly divided into two scene sets (set A and set B), each composed of 20 scenes, and for each subject each set was assigned to one task. Task order and scene set assignment was fully counterbalanced across all participants. Additionally, scenes within each set were presented in a randomized order for each participant in each condition.

Analysis

All analyses were chosen *a priori* and based on our previous work (Henderson & Hayes, 2017, 2018; Henderson et al., 2018).

Data Segmentation and Outliers

Fixations and saccades were segmented with EyeLink's standard algorithm using velocity and acceleration thresholds ($30^\circ/s$ and $9500^\circ/s^2$; SR Research, 2010b). Eye movement data were imported offline into Matlab using the EDFConverter tool. The first fixation on each scene, always located at the center of the display as a result of the pretrial fixation marker,

was eliminated from analysis. Additionally, any fixations that were shorter than 50ms and longer than 1500ms were eliminated as outliers. This outlier removal process resulted in loss of 3.94% of the data.

Attention Maps

Attention maps were generated as described in Henderson and Hayes (2017). Briefly, a fixation frequency matrix based on the locations (x,y coordinates) of all fixations was generated across participants for each scene. A Gaussian low-pass filter with a circular boundary and a cutoff frequency of -6dB was applied to each matrix to account for foveal acuity and eyetracker error (Figure 2). The spatial extent of the low pass filter was 236 pixels in diameter.

Meaning Maps

Meaning maps were generated as per Henderson and Hayes (2017). Because the nature of our tasks resulted in peripheral fixations, we used both unbiased and center-biased meaning maps (Figure 2). Overall, the unbiased maps provided better predictive power than the center-biased maps. However, we included analyses from both because center-biased maps are standard in the literature and thus provide a basis for comparison with previous studies. The center-biased meaning maps were generated by applying a multiplicative center bias operation to the meaning maps using the same center bias present in the saliency maps.

Subjects.—Scene patches were rated by 165 subjects on Amazon Mechanical Turk. Participants were recruited from the United States, had a HIT (human intelligence task) approval rate of 99% and 500 HITs approved, and were only allowed to participate in the study once. Participants were paid \$0.50 cents per assignment, and all participants provided informed consent.

Stimuli.—Stimuli consisted of the 40 digitized photographs used in the current experiment. Each scene was decomposed into a series of partially overlapping and tiled circular patches at coarse and fine spatial scales. The full patch stimulus set consisted of 12,000 unique fine patches and 4,320 unique coarse patches for a total of 16,320 scene patches.

Procedure.—Each participant rated 300 random scene patches extracted from the scenes. Participants were instructed to assess the meaningfulness of each patch based on how informative or recognizable they thought it was. During the instruction period, participants were provided with examples of two low-meaning and two high-meaning scene patches to make sure they understood the task. They then rated the meaningfulness of test patches on a six-point Likert scale ('very low', 'low', 'somewhat low', 'somewhat high', 'high', 'very high'). Patches were presented in random order and without scene context, so ratings were based on context-independent judgments. Each unique patch was rated three times by three independent raters for a total of 48,960 ratings. However, owing to the high degree of overlap across patches, each fine patch contained rating information from 27 independent raters, and each coarse patch from 63 independent raters.

Meaning maps were generated from the ratings by averaging, smoothing and combining the fine and coarse maps from the corresponding patch ratings. The ratings for each pixel at each scale in each scene were averaged, producing an average fine and coarse rating map for each scene. The average fine and coarse rating maps were then smoothed using thin-plate spline interpolation based on the center of each patch (MATLAB ‘fit’ using the ‘thinplateinterp’ method). Finally, the smoothed fine and coarse maps were averaged to produce the meaning map for each scene.

Saliency Maps

Saliency maps for each scene were computed using the Graph-Based Visual Saliency (GBVS) toolbox with default settings (Harel et al., 2006). GBVS is a prominent saliency model that combines maps of low-level image features to create image-based saliency maps (Figure 2).

Center bias is a natural feature of the GBVS saliency maps. To compare to unbiased meaning maps, we also generated GBVS maps without center bias (Figure 2). These maps were created using a whitening method (Rahman & Bruce, 2015), a 2-step normalization approach in which each saliency map is normalized to have 0 mean and unit variance. After this, a second pixel-wise normalization is performed so each pixel location across all the saliency maps has 0 mean and unit variance.

Histogram matching.—Following Henderson and Hayes (2017), meaning and saliency maps were normalized to a common scale using image histogram matching with the fixation density map for each scene serving as the reference image for the corresponding meaning and saliency maps. This was accomplished by using the Matlab function ‘imhistmatch’ from the Image Processing Toolbox.

Results

Task Comparisons

Scene viewing.—Because we gave participants the option to terminate each presentation trial early, we began by comparing the average scene-viewing (scene onset to response) time for each scene during each condition as well as the number of fixations per scene in each task (Figure 3). The average scene viewing time for the brightness rating task was 5262.55ms ($SD = 3141.39$) with 15.56 fixations ($SD = 9.84$), and for the brightness search task was 10726.52ms ($SD = 2420.55$) with 32.28 fixations ($SD = 8.20$). Because the distributions were not normal (Figure 3), Wilcoxon rank sum tests were conducted and showed that the scene viewing times and number of fixations were significantly different between the rating and search tasks: Z 's > 5.50 ; p 's < 0.001 . These results showed that participants tended to view scenes during the rating task for shorter durations than the search task, with participants much more likely to use the entire 12s in the search compared to the rating task. The finding that the rating task produced significantly shorter viewing durations than the search task suggests that participants only viewed the scenes for the amount of time necessary to complete each task. Given that the viewing times and number of fixations were

very different between the tasks, we treated the two tasks separately in the following analyses.

Response agreement.—To verify that subjects were staying on task and attending to brightness during the study, we examined response agreement in the rating and search tasks. If subjects were on-task, then their responses should vary as a function of scene and be consistent within scenes. That is, subjects should generally agree in their judgements of brightness in the rating task and the number of bright regions in the search task. On the other hand, if subjects were simply attending to scene content rather than following instructions, then responses should be unsystematic across scenes and subjects. As can be seen in Figure 4, the former was true, suggesting that subjects were indeed following instructions.

Overall Scene Analyses

Following Henderson and Hayes (2017), we used squared linear and semi-partial correlations to quantify the degree to which meaning maps and saliency maps accounted for shared and unique variance in the attention maps. Specifically, we conducted two-tailed, two-sample t-tests for the correlations across scenes to statistically compare the relative ability of meaning and salience to predict attentional guidance.

For comparison to the literature, we tested how well traditional center-biased meaning and saliency maps could account for attention. In addition, because the center bias was substantially reduced in the brightness search task compared to the brightness rating task (Figure 5), we also conducted analyses using unbiased meaning and saliency maps that excluded a center bias.

Brightness rating task.—Using the center-biased maps, for squared linear correlations on average across all 40 scenes, meaning accounted for 55% of the variance in fixation density ($M = 0.55$, $SD = 0.12$) and salience accounted for 33% of the variance in fixation density ($M = 0.33$, $SD = 0.14$) (Figure 6). This difference between meaning and saliency maps was significant: $t(78) = 7.31$, $p < 0.001$, 95% CI = [0.16, 0.28]. Similarly, for squared semi-partial correlations, meaning accounted for 24% of the variance in fixation density ($M = 0.24$, $SD = 0.13$) controlling for salience, but salience accounted for only 3% of the variance in fixation density controlling for meaning ($M = 0.03$, $SD = 0.03$) (Figure 6). This difference was again significant: $t(78) = 10.57$, $p < 0.001$, 95% CI = [0.17, 0.25]. This pattern of results did not change when using the unbiased meaning and saliency maps (linear: $t(78) = 8.79$, $p < 0.001$, 95% CI = [0.16, 0.25]; semi-partial: $t(78) = 9.62$, $p < 0.001$, 95% CI = [0.16, 0.25]) (Figure 6). These findings suggest that meaning played a dominant role in the guidance of attention even though meaning was irrelevant and salience was central to the brightness rating task.

Brightness search task.—Using the center-biased maps, meaning accounted for 22% of the variance in fixation density ($M = 0.22$, $SD = 0.13$) and salience accounted for 24% of the variance in fixation density ($M = 0.24$, $SD = 0.12$) (Figure 7). This difference was not significant: $t(78) = -0.33$, $p = 0.74$, 95% CI = [-0.07, 0.05]. Similarly, for the semi-partial correlations, meaning accounted for 5% of the variance in fixation density controlling for

salience ($M = 0.05$, $SD = 0.07$) and salience accounted for 6% of the variance in fixation density controlling for meaning ($M = 0.06$, $SD = 0.07$) (Figure 7). Again, this difference was not significant: $t(78) = -0.59$, $p = 0.56$, 95% CI = $[-0.04, 0.02]$. Importantly, however, this pattern of results changed when using the unbiased meaning and saliency maps (Figure 7). Using the unbiased maps, meaning accounted for 22% of the overall variance in attention ($M = 0.22$, $SD = 0.11$) whereas salience explained only 4% of the variance ($M = 0.04$, $SD = 0.05$) for the linear correlations, $t(78) = 6.42$, $p < 0.001$, 95% CI = $[0.10, 0.18]$. Similarly, for the semi-partial correlations, meaning accounted for 18% of the total variance in attention ($M = 0.18$, $SD = 0.11$) whereas salience explained only 1% of the variance ($M = 0.04$, $SD = 0.04$), $t(78) = 7.42$, $p < 0.001$, 95% CI = $[0.10, 0.18]$. These findings suggest that when the more distributed nature of attention away from scene centers and to scene peripheries in the brightness search task was taken into account, meaning influenced attentional guidance more than salience even though meaning was irrelevant and saliency was central to the task.

Fixation by Fixation Analyses

Previously, it has been posited that attention during scene viewing might initially be guided by salience, but that as time progresses, meaning begins to play an increasing role (Anderson, Donk, & Meeter, 2016; Anderson, Ort, Kruijne, Meeter, & Donk, 2015; Henderson Ferreira, 2004; Henderson & Hollingworth, 1999; Parkhurst et al., 2002). On the other hand, in two studies investigating the roles of meaning and salience in memorization and scene description tasks, we did not observe this change from guidance by salience to guidance by meaning (Henderson & Hayes, 2017; Henderson et al., 2018). Instead, meaning was found to guide attention from the first saccade. Because the current tasks were designed to make meaning irrelevant and salience central, they provide another opportunity to test this hypothesis.

We conducted a temporal time-step analysis in which a series of attention maps were generated from each sequential fixation (1st fixation, 2nd fixation, 3rd fixation, etc.) for each scene in each task. We then correlated each attention map for each fixation and scene using both the center-biased and unbiased meaning and saliency maps to calculate the squared linear and semi-partial correlations. Then the correlations for each scene and fixation were averaged across scenes to assess how meaning and image salience predicted attention on a fixation by fixation basis. The prediction of the salience first hypothesis is that the correlation between saliency and attention maps should be greater for earlier than later fixations, with salience dominating meaning in the earliest fixations.

Brightness rating task.—Using the center-biased maps, meaning accounted for 34%, 23%, and 17% of the variance in the first 3 fixations whereas salience accounted for 8%, 12%, and 11% of the variance in the first 3 fixations, respectively, for the linear correlations (Figure 8). Two-sample, two-tailed t-tests compared meaning and salience for all 8 initial fixations using p-values corrected for multiple comparisons using a false discovery rate (FDR) correction (Benjamini & Hochberg, 1995). Overall, this confirmed the advantage for meaning over salience for all 8 fixations (all FDR p s < 0.05). Similarly, for the semi-partial correlations, meaning accounted for 28%, 14%, and 9% of the variance in the first 3 fixations and salience accounted for 2%, 3%, and 3% of the variance in the first 3 fixations

(Figure 8). Again, meaning predicted attention significantly better than salience for all 8 initial fixations (all FDR p s < 0.001). Using the unbiased maps, this overall pattern of results did not change (linear and semi-partial correlations: all FDR p s < 0.001) (Figure 8). These results do not support the hypothesis that the influence of meaning on attentional guidance was delayed to later fixations.

Brightness search task.—Using the center-biased maps, meaning accounted for 30%, 14%, and 7% of the variance in the first 3 fixations and salience accounted for 11%, 16%, and 14% in the first 3 fixations, respectively, for the linear correlations (Figure 9). Here, meaning produced an advantage over salience on the first fixation (FDR p < 0.001) but not fixations 2 through 8 (FDR p > 0.05). For the semi-partial correlations, meaning explained 22%, 8%, and 3% of the variance in the first 3 fixations and salience accounted for 3%, 10%, and 10% in the first 3 fixations (Figure 9). A significant advantage for meaning was observed on the first fixation (FDR p < 0.001) and for salience on the third fixation (FDR p < 0.05), with no other comparisons reaching significance (FDR p s > 0.05).

Using the unbiased meaning and saliency maps, the pattern of results changed. For the linear correlations, meaning accounted for 5%, 6%, and 4% of the variance and salience accounted for 1%, 4%, and 5% of the variance in attention in the first 3 fixations. Turning to the semi-partial correlations, meaning accounted for 5%, 6%, and 3% of the variance and salience accounted for 0.1%, 3%, and 4% of the variance in attention in the first 3 fixations. Meaning still produced an advantage over salience for the first fixation (linear and semi-partial FDR p < 0.05) with all other fixations nonsignificant (linear and semi-partial FDR p > 0.05). The advantage for saliency over meaning for the third fixation seen in the center-biased maps was not observed with the unbiased maps.

The fixation by fixation analyses were not consistent with the salience first hypothesis. In the analyses using both the center-biased and unbiased maps, meaning was more important than salience at the first fixation. Using the center-biased maps, salience was stronger in the third fixation. This result, however, was not true using the unbiased maps, suggesting that the advantage for saliency in the center-biased maps was driven by the center bias rather than saliency itself. Overall, the results are not consistent with the hypothesis that attentional guidance transitions from salience to meaning over time.

Saccade Amplitude Analyses

In the analyses thus far, fixations following both shorter and longer saccades were included. It could be that meaning guides attention within local scene regions, whereas salience guides attention as it moves from one scene region to another. To test this hypothesis, we analyzed the role of meaning on attentional guidance as a function of saccade amplitude. If meaning plays a greater role for local (e.g., within-object) shifts of attention, then meaning should be more related to attentional selection following shorter saccades versus longer saccades. Such a pattern might be more likely in the case of the current study because meaning was not relevant to the tasks. To investigate this hypothesis, we assessed how meaning and salience related to attention following saccades of shorter to longer amplitudes (Figure 10). Specifically, saccade amplitudes were binned by decile, and fixation density maps were

created for each saccade amplitude decile. Meaning and salience maps were then correlated with the fixation density maps for each decile. We conducted these analyses using both the center-biased and unbiased meaning and saliency maps. The saccade amplitude average for the rating task was 5.37° ($SD = 3.41$) and for the search task was 4.61° ($SD = 3.51$).

Brightness rating task.—For the brightness rating task, using the center-biased maps, meaning produced an advantage over salience for saccade amplitude deciles 1 through 7 and 9 (FDR $p < 0.05$) but not deciles 8 and 10 (FDR $p > 0.05$). For the semi-partial correlations, meaning explained significantly more of the variance in fixation density than salience for all 10 saccade amplitude deciles (all FDR $ps < 0.05$). When using the unbiased meaning and saliency maps, this pattern of results became stronger as meaning produced an advantage over saliency across all deciles in both the linear and semi-partial correlations (FDR $p < 0.05$).

Brightness search task.—For the brightness search task, using the center-biased maps, there were no significant differences between meaning and salience for any saccade amplitude deciles in either the linear or semi-partial correlations (all FDR $ps > 0.05$). When using the unbiased maps, on the other hand, this pattern of results changed as meaning produced an advantage over saliency for saccade amplitude deciles 1 through 9 (FDR $p < 0.05$) but not 10 (FDR $p > 0.05$).

Overall, it appears that meaning was used to guide attention for both short and long shifts of attention, though there was some evidence that this influence was reduced when the scene peripheries were removed from the analyses (i.e., with the center-biased maps) and for the longest shifts of attention.

General Discussion

Past research has emphasized image salience as a key basis for attentional selection during real-world scene viewing (Borji et al., 2014; Borji et al., 2013; Harel et al., 2006; Itti & Koch, 2001; Koch & Ullman, 1985; Parkhurst et al., 2002). Although this previous work has provided an important framework for understanding attentional guidance in scenes, it downplays the fact that attention is strongly guided by cognitive factors related to semantic features that are relevant to understanding the scene in the context of the task (Buswell, 1935; Hayhoe & Ballard, 2005; Hayhoe et al., 2003; Henderson, Brockmole, Castelano, & Mack, 2007; Henderson, Malcolm, & Schandl, 2009; Land & Hayhoe, 2001; Rothkopf et al., 2016; Yarbus, 1967). With the development of meaning maps, which capture the spatial distribution of semantic content in scenes in the same format that saliency maps capture the spatial distribution of image salience, it has become possible to directly compare the influence of meaning and image salience on attention in scenes (Henderson & Hayes, 2017).

In prior studies comparing meaning and image salience during scene viewing, meaning has better explained the spatial and temporal patterns of attention (Henderson & Hayes, 2017, 2018; Henderson et al., 2018). However, those studies used memorization, aesthetic judgement, and scene description viewing tasks, and it could be argued that those tasks were biased towards attentional guidance by meaning. The current study sought to determine

whether the influence of meaning on attention would be eliminated in tasks that do not require any semantic analysis of the scenes. To test this hypothesis, we used two viewing tasks designed to eliminate the need for attending to meaning: a brightness rating task in which participants rated the overall brightness of scenes, and a brightness search task in which participants counted the number of bright areas in scenes.

For the brightness rating task, we found that meaning better explained the spatial distribution of attention than image salience. This result was observed both overall and when the correlation between meaning and image salience was statistically controlled, and held for early scene viewing, for short and long saccades, and using center-biased and unbiased meaning and saliency maps. For the brightness search task, using center-biased meaning and saliency maps, there were no differences between meaning and salience overall or when controlling for their correlation. However, the center-biased maps did not capture the fact that during the search task, the center bias in attention was greatly attenuated because attention was distributed much more uniformly over the scenes. Meaning and saliency maps with center bias over-weight scene centers and ignore scene peripheries, opposite to the attention maps actually observed. When the attention maps were analyzed using meaning and saliency maps that did not include center bias, the results were similar to those of the brightness rating task: meaning explained the variance in attention better than salience both overall and after statistically controlling for the correlation between meaning and salience. This pattern held for short and long saccades, and for the first saccade.

Overall, the results provide strong evidence that the meaning of a scene plays an important role in guiding attention through real-world scenes even when meaning is irrelevant and image salience is relevant to the task. Converging evidence across two viewing tasks that focused on an image property related to image salience showed that meaning accounted for more variance in attentional guidance than salience, and critically, that when the correlation between meaning and salience was controlled, only meaning accounted for significant unique variance. These results indicate that the guidance of attention by meaning is not restricted to viewing tasks that focus on encoding the meaning of the scene, strongly suggesting a fundamental role of meaning in attentional guidance in scenes.

Although the main pattern of results was clear and generally consistent across the two tasks, a few points are worth additional comment. First, our results suggest that tasks can differ in the degree to which center bias is present. Here, center bias was much greater when judging overall scene brightness than when searching for bright scene regions. These differences in center bias for the rating and search tasks likely occurred due to differences in the requirements of the tasks. The rating task simply required participants to rate the overall brightness of scenes, so there was no particular reason for viewers to direct attention away from the centers and to the peripheries of the scenes. In comparison, the search task required participants to count individual bright regions, many of which appeared away from the scene centers and in the peripheries. This resulted in fewer central fixations and more peripheral fixations in the brightness search task than the brightness rating task. Because there were more peripheral fixations in the search task, the center-biased meaning and saliency maps did not have the same predictive power to capture the relationship between meaning, salience, and attention as they did for the brightness rating task. Indeed, for this reason

neither meaning nor saliency maps did a particularly good job of predicting attention when center-bias was included in the maps. However, when the center bias was removed from the two prediction maps, meaning maps were significantly better than saliency maps in accounting for attention.

The difference between the center-biased and unbiased maps was also evident in the analysis focusing on the earliest eye movements. According to the “salience first” hypothesis, we should have seen an initial bias of attention toward salient regions followed by a shift to meaningful regions. In our prior studies, we instead observed that meaning guided attention from the very first eye movement (Henderson & Hayes, 2018, Henderson et al., 2018). In the present study, when center-bias was included in the meaning and saliency maps in the brightness search task, meaning initially guided attention in the first eye movement, but there was a tendency for salience to take over for a few saccades before meaning again dominated. This pattern might offer some small support for salience first. However, as noted, viewers were much less likely to attend to scene centers and more likely to move their eyes to the edges of the scenes in the brightness search task. When the unbiased maps were used in the search task analysis, the trend from meaning to salience over the first few fixations was not observed. At best, then, there is a hint that when the viewer’s task is explicitly to find and count salient scene regions, they may be slightly more biased early on to attend to regions that are more salient. However, this result is weak at best given it appeared only in the third fixation and disappeared in the unbiased map analysis. Overall, even in a task that explicitly focused on salience and in which meaning was completely irrelevant, meaning played a stronger role in attentional guidance from the very beginning of viewing.

The type of meaning studied in the current work is what we refer to as context-free meaning, in that it is based on ratings of the recognizability and informativeness of isolated scene patches shown to raters independently of the scenes from which they are derived and independently of any task or goal besides the rating itself. Other types of meaning may be of interest in future studies. For example, we can consider contextualized meaning in which meaning is determined based on how important a scene patch is with respect to its global scene context. Additionally, the role of task may affect meaning as well. For example, meaning within a scene may change depending on a viewer’s current tasks or goals. Because meaning can be defined in so many ways, it is necessary that we understand how these variants influence attentional guidance. The meaning map approach provides a method for pursuing these important questions.

Conclusion

We investigated the relative importance of meaning and image salience on attentional guidance in scenes using tasks that do not require semantic analysis and in which salience plays a critical role. Overall, the results strongly suggested that viewers can’t help but attend to meaning (Greene & Fei-Fei, 2014). These findings are most consistent with cognitive control theories of scene viewing in which attentional priority is assigned to scene regions based on semantic properties rather than image properties.

Acknowledgements

We thank Conner Daniels, Praveena Singh, and Diego Marquez for assisting in data collection. This research was supported by the National Eye Institute of the National Institutes of Health under award number R01EY027792. The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health.

References

- Anderson NC, Donk M, & Meeter M (2016). The influence of a scene preview on eye movement behavior in natural scenes. *Psychonomic Bulletin and Review*, 23, 1794–1801. 10.3758/s13423-016-1035-4 [PubMed: 27073087]
- Anderson NC, Ort E, Kruijne W, Meeter M, & Donk M (2015). It depends on when you look at it: Saliency influences eye movements in natural scene viewing and search early in time. *Journal of Vision*, 15(5), 1–22. 10.1167/15.5.9
- Antes JR (1974). The time course of picture viewing. *Journal of Experimental Psychology*, 103(1), 62–70. 10.1037/h0036799 [PubMed: 4424680]
- Benjamini Y, & Hochberg Y (1995). Controlling the false discovery rate: A practical and powerful approach to multiple testing. *Royal Statistical Society B*, 57(1), 289–300.
- Borji A, Parks D, & Itti L (2014). Complementary effects of gaze direction and early saliency in guiding fixations during free viewing. *Journal of Vision*, 14(13), 3–3.
- Borji A, Sihite DN, & Itti L (2013). Quantitative analysis of human-model agreement in visual saliency modeling: A comparative study. *IEEE Transactions on Image Processing*, 22(1), 55–69. 10.1109/TIP.2012.2210727 [PubMed: 22868572]
- Buswell GT (1935). *How people look at pictures: a study of the psychology and perception in art* Oxford, England.
- Carmi R, & Itti L (2006). The role of memory in guiding attention during natural vision. *Journal of Vision*, 6(9), 898–914. doi:10.1.1.79.1508 [PubMed: 17083283]
- Greene MR, & Fei-Fei L (2014). Visual categorization is automatic and obligatory: Evidence from Stroop-like paradigm. *Journal of Vision*, 14(1), 14–14. 10.1167/14.1.14
- Harel J, Koch C, & Perona P (2006). Graph-based visual saliency. *Advances in Neural Information Processing Systems*
- Hayhoe M, & Ballard D (2005). Eye movements in natural behavior. *Trends in Cognitive Sciences*, 9(4), 188–194. 10.1016/j.tics.2005.02.009 [PubMed: 15808501]
- Hayhoe MM, Shrivastava A, Mruczek R, & Pelz JB (2003). Visual memory and motor planning in a natural task. *Journal of Vision*, 3(6), 49–63. 10.1167/3.1.6 [PubMed: 12678625]
- Henderson JM (2007). Regarding scenes. *Current Directions in Psychological Science*, 16(4), 219–222. 10.1111/j.1467-8721.2007.00507.x
- Henderson JM (2017). Gaze control as prediction. *Trends in Cognitive Sciences*, 21(1), 15–23. 10.1016/j.tics.2016.11.003 [PubMed: 27931846]
- Henderson JM, Brockmole JR, Castelano MS, & Mack M (2007). Visual saliency does not account for eye movements during visual search in real-world scenes. In Gompel R. P. G. v., Fischer MH, Murray WS, & Hill RL (Eds.), *Eye Movements: A Window on Mind and Brain* (pp. 537–562): Elsevier Ltd.
- Henderson JM, Hayes TR, Rehrig G, & Ferreira F (2018). Meaning guides attention during real-world scene description. *Scientific Reports*, 8(13504). doi: 10.1038/s41598-018-31894-5
- Henderson JM, & Ferreira F (2004). Scene perception for psycholinguists *The interface of language, vision, and action: Eye movements and the visual world* (pp. 1–58). New York, NY, US: Psychology Press.
- Henderson JM, & Hayes TR (2017). Meaning-based guidance of attention in scenes as revealed by meaning maps. *Nature Human Behavior*, 1, 743–747. 10.1038/s41562-017-0208-0
- Henderson JM, & Hayes TR (2018). Meaning guides attention in real-world scene images: Evidence from eye movements and meaning maps. *Journal of Vision*, 18(6), 1–18.

- Henderson JM, & Hollingworth A (1999). High-level scene perception. *Annual Review of Psychology*, 50(243–271).
- Henderson JM, & Hollingworth A (1999). The role of fixation position in detecting scene changes across saccades. *Psychological Science*, 10(5), 438–443.
- Henderson JM, Malcolm GL, & Schandl C (2009). Searching in the dark: Cognitive relevance drives attention in real-world scenes. *Psychonomic Bulletin and Review*, 16(5), 850–856. 10.3758/PBR.16.5.850 [PubMed: 19815788]
- Itti L, & Koch C (2001). Computational modelling of visual attention. *Nature Reviews*, 2(3), 1–11.
- Itti L, Koch C, & Niebur E (1998). A model of saliency-based visual attention for rapid scene analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20(11).
- Koch C, & Ullman S (1985). Shifts in selective visual attention: Towards the underlying neural circuitry. *Human Neurobiology*, 4(4), 219–227. [PubMed: 3836989]
- Land MF, & Hayhoe M (2001). In what ways do eye movements contribute to everyday activities? *Vision Research*, 41(25–26), 3559–3565. 10.1016/S0042-6989(01)00102-X [PubMed: 11718795]
- Mackworth NH, & Morandi AJ (1967). The gaze selects informative details within pictures. *Perception and Psychophysics*, 2(11), 547–552.
- Navalpakkam V, & Itti L (2005). Modeling the influence of task on attention. *Vision Research*, 45(2), 205–231. 10.1016/j.visres.2004.07.042 [PubMed: 15581921]
- Navalpakkam V, & Itti L (2007). Search goal tunes visual features optimally. *Neuron*, 53(4), 605–617. doi: 10.1016/j.neuron.2007.01.018 [PubMed: 17296560]
- Nuthmann A, & Henderson JM (2010). Object based attentional selection in scene viewing. *Journal of Vision*, 10(8): 20, 1–19.
- Parkhurst D, Law K, & Niebur E (2002). Modeling the role of salience in the allocation of overt visual attention. *Vision Research*, 42, 107–123. [PubMed: 11804636]
- Rahman S, & Bruce N (2015). Visual saliency prediction and evaluation across different perceptual tasks. *PlosOne* 10.1371/journal.pone.0138053
- Rothkopf CA, Ballard DH, & Hayhoe MM (2016). Task and context determine where you look. *Journal of Vision*, 7(16), 1–20. 10.1167/7.14.16
- Tatler BW, Hayhoe MM, Land MF, & Ballard DH (2011). Eye guidance in natural vision: Reinterpreting salience. *Journal of Vision*, 11(5), 5–5.
- Torralba A, Oliva A, Castelano MS, & Henderson JM (2006). Contextual guidance of eye movements and attention in real-world scenes: The role of global features in object search. *Psychological Review*, 113(4), 766–786. 10.1037/0033-295X.113.4.766 [PubMed: 17014302]
- Yarbus AL (1967). Eye movements during perception of complex objects *Eye Movements and Vision* (pp. 171–211). Boston, MA: Springer.

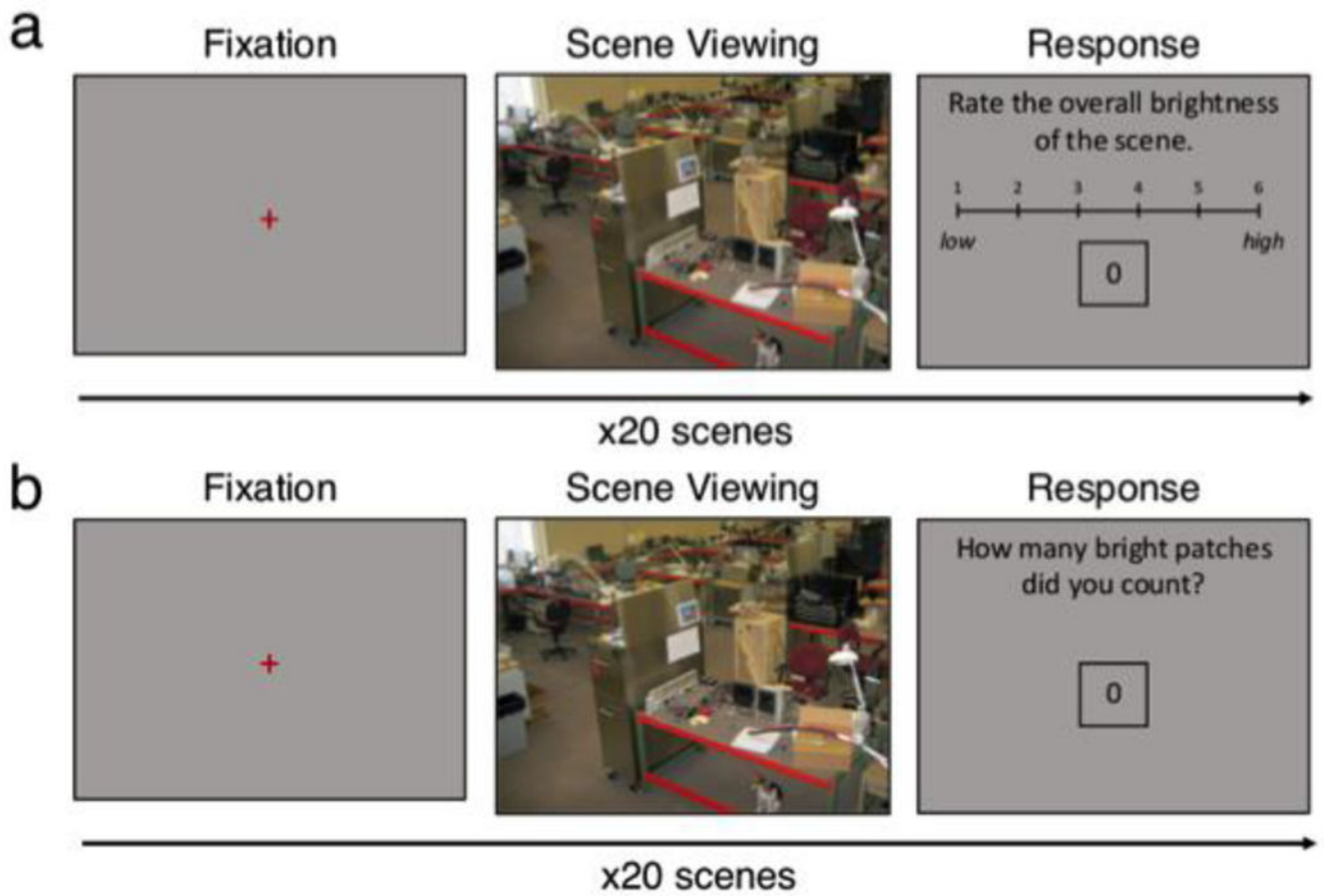


Figure 1. Trial structure for the two tasks.

The trial structure for (a) the brightness rating task and (b) the brightness search task.

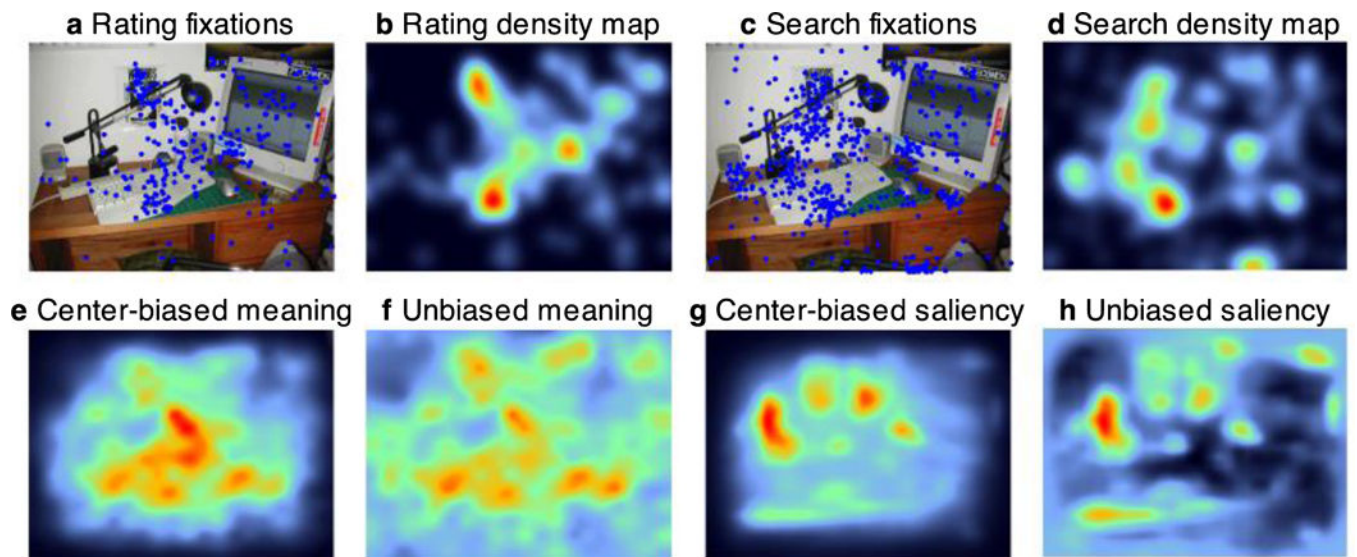


Figure 2. An example scene with the associated maps for each task.

(a) is an example scene with fixation locations from all participants in the rating task aggregated and overlaid. (b) is the fixation density map representing the example scene and fixation locations for the rating task. (c) is the example scene with fixations from the search task overlaid and (d) is the fixation density map representing the example scene and fixation locations in the search task. (e) is the center-biased meaning map and (f) is the unbiased meaning map for the example scene. (g) is the center-biased saliency map and (h) is the unbiased saliency map for the example scene.

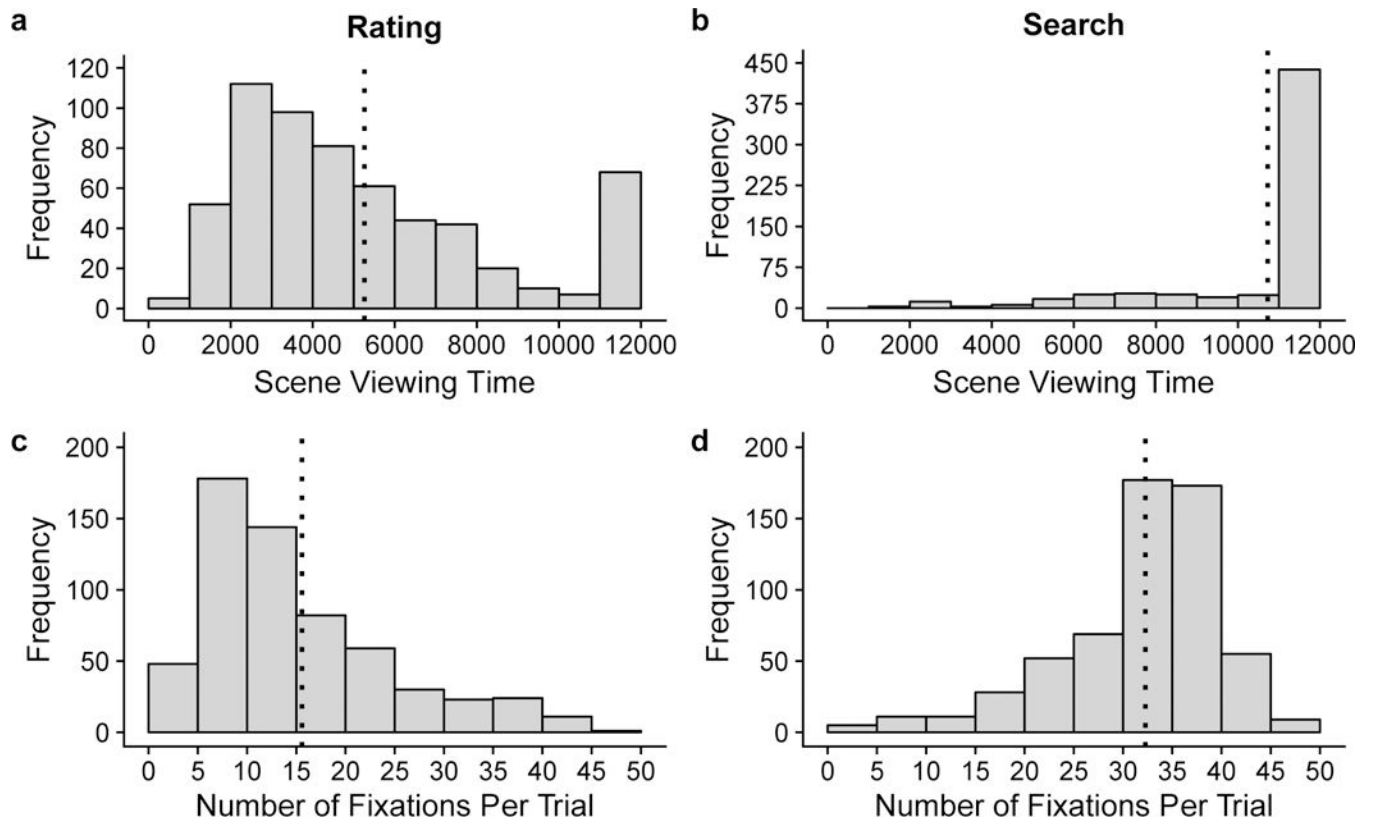


Figure 3. Scene-viewing times and number of fixations per trial for the brightness rating and brightness search tasks.

Distributions are shown for the scene viewing times of (a) the brightness rating and (b) brightness search tasks, and the number of fixations per trial of (c) the brightness rating and (d) brightness search tasks. Black dotted vertical lines represent the mean for each task.

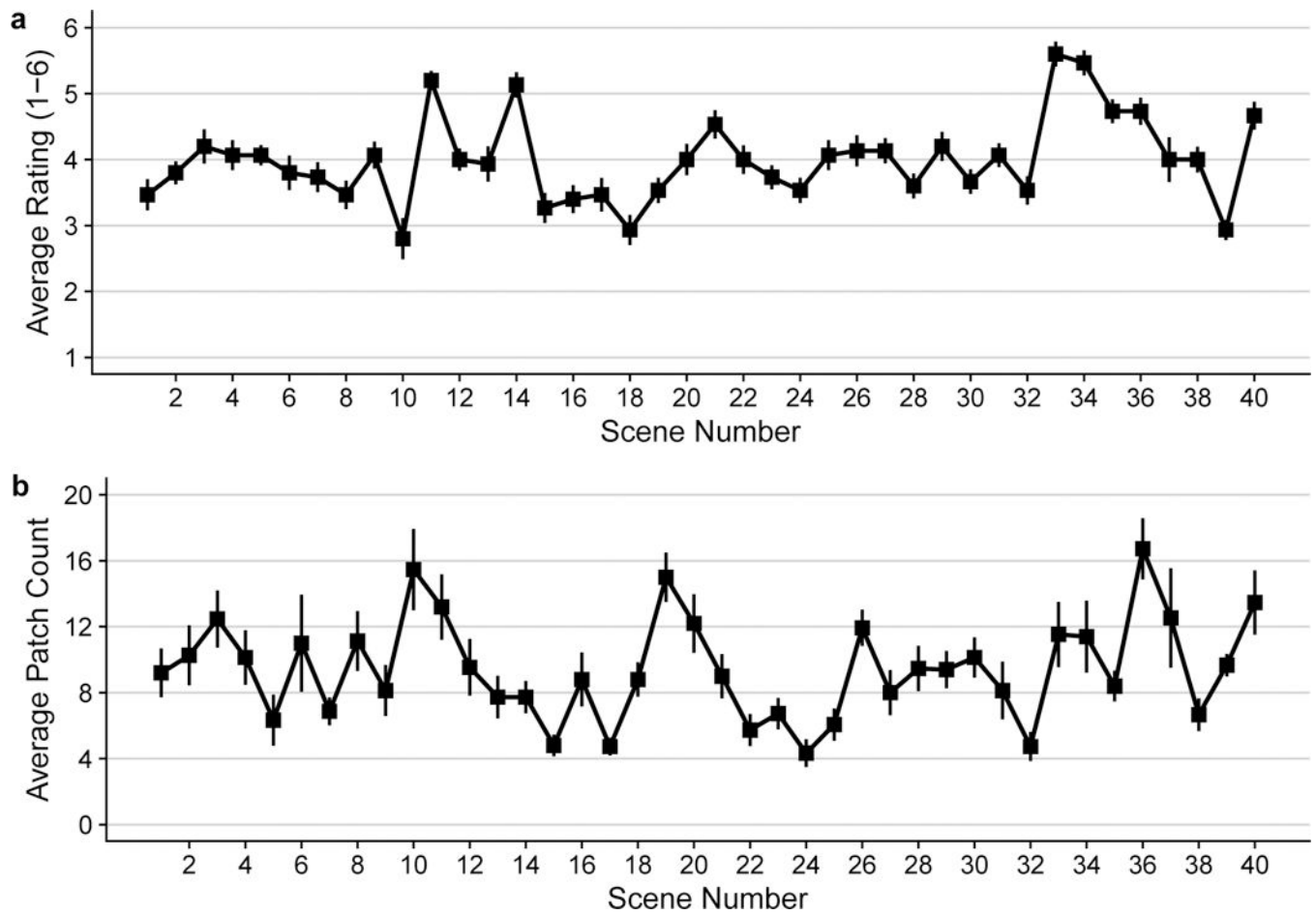


Figure 4. Response variability as a function of scene.

The average participant response and standard error of responses as a function of scene for (a) the brightness rating task and (b) the patch count task.

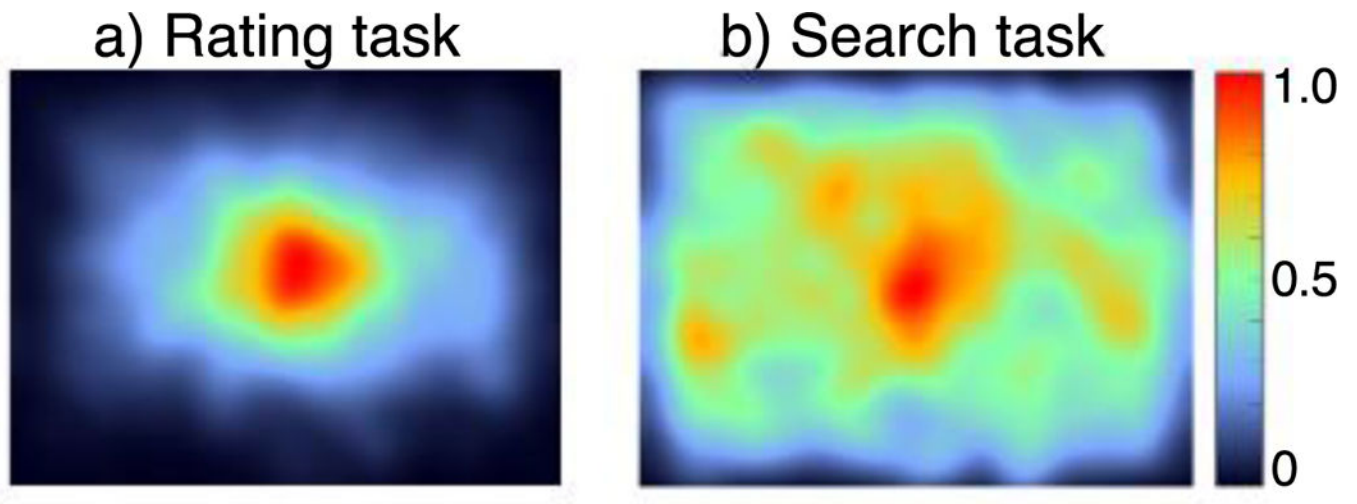


Figure 5. Center biased maps.

Fixation density maps aggregated across subjects and scenes are shown for (a) the brightness rating task and (b) the brightness search task.

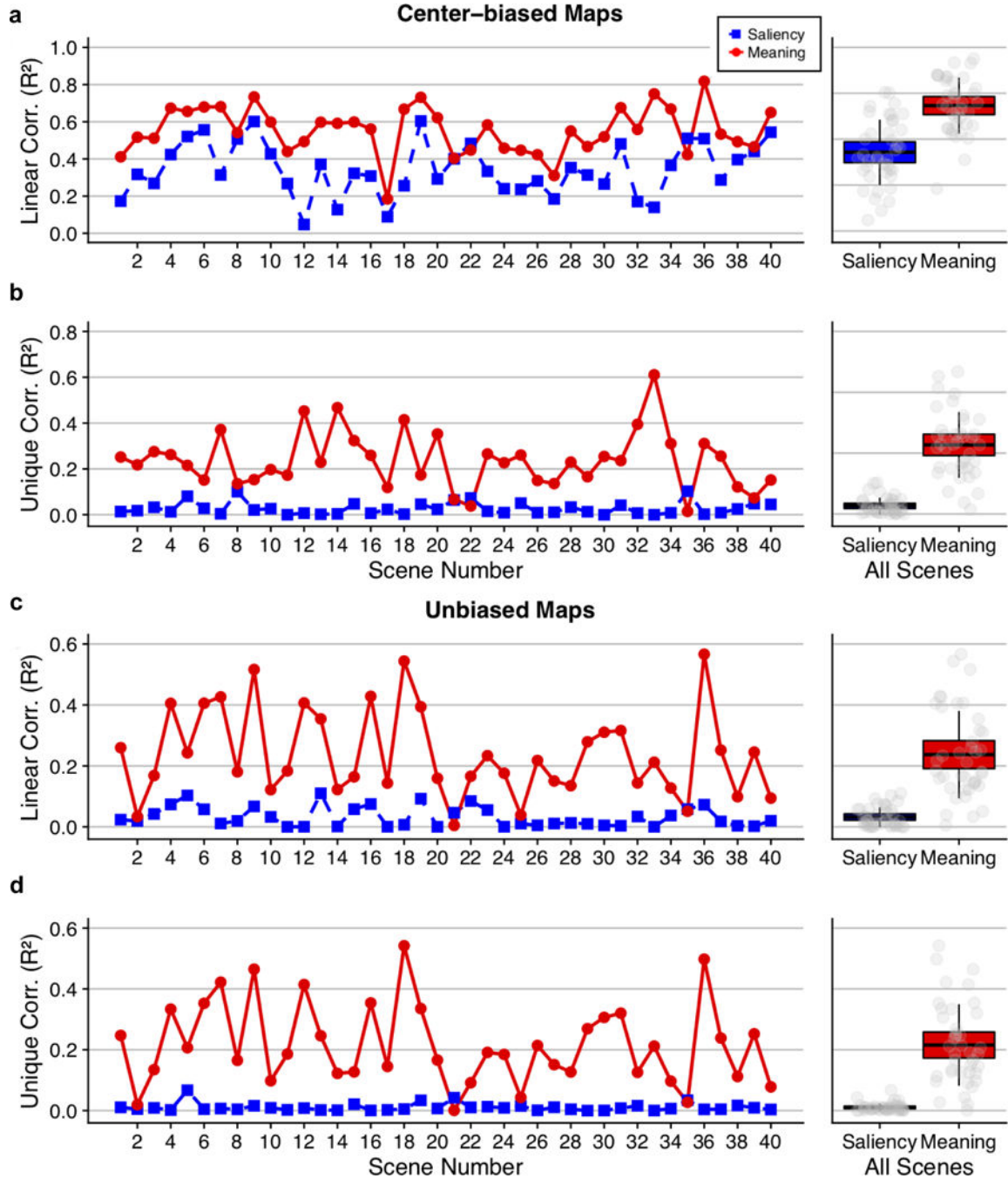


Figure 6. Squared linear correlation and semi-partial correlation by scene for the brightness rating task.

Line plots show (a) the squared linear and (b) semi-partial correlations between the fixation density maps and meaning (red circles) and saliency (blue squares) using center-biased meaning and saliency maps. Line plots also show (c) the squared linear and (d) semi-partial correlations using unbiased meaning and saliency maps. The scatter plots on the right show the grand mean (black horizontal line), 95% confidence intervals (colored boxes), and 1 standard deviation (black vertical line) for meaning and saliency across all 40 scenes for each analysis.

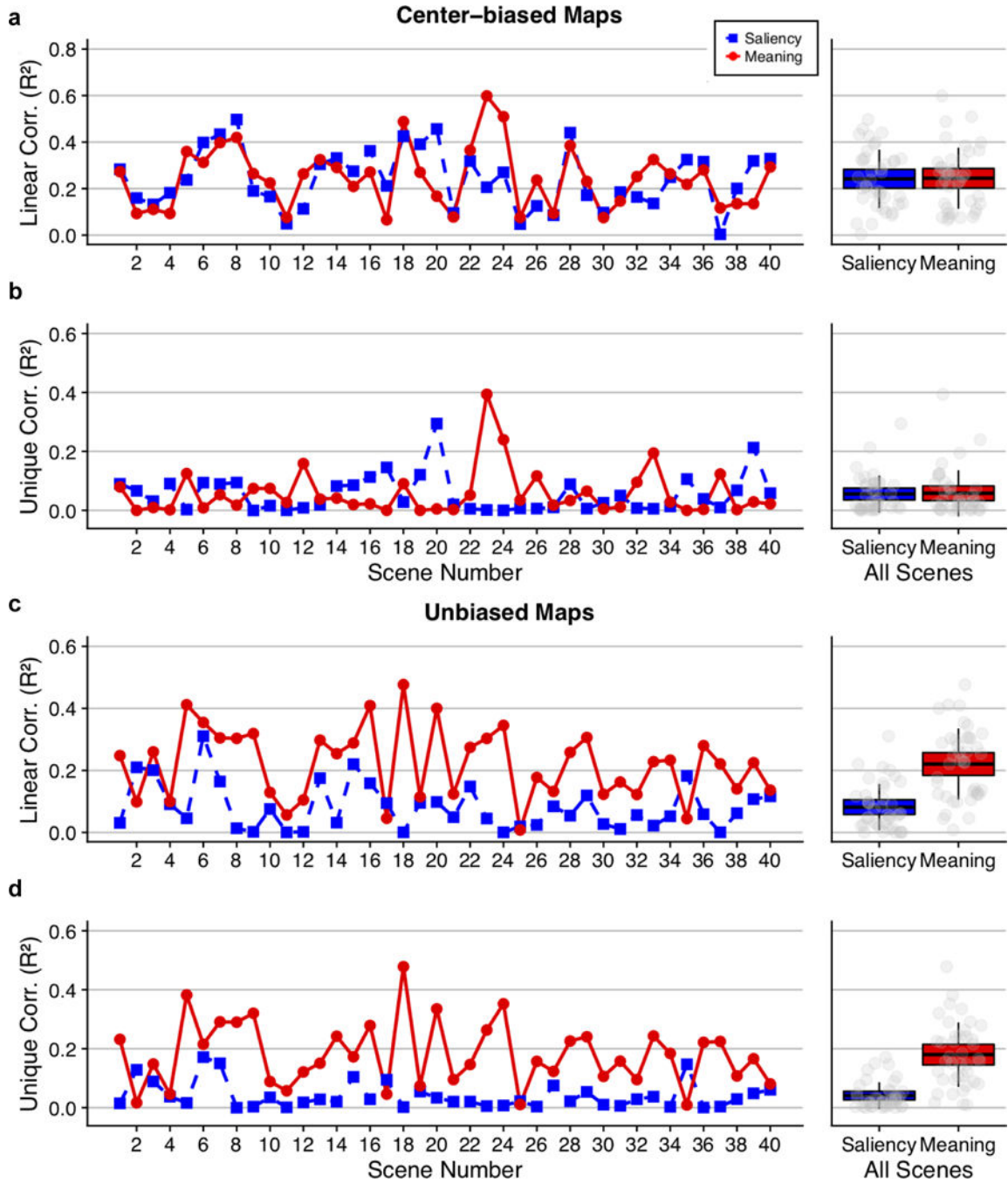


Figure 7. Squared linear correlation and semi-partial correlation by scene for the brightness search task.

The line plots show (a) the linear and (b) semi-partial correlations between fixation density and meaning (red circles) and saliency (blue squares) for the search task using the center-biased meaning and saliency maps. Line plots also show (c) the linear and (d) semi-partial correlations for the search task using the unbiased meaning and saliency maps. The scatter plots on the right show the corresponding grand mean (black line), 95% confidence intervals (colored box), and one standard deviation (black vertical line) for meaning and saliency across all scenes.

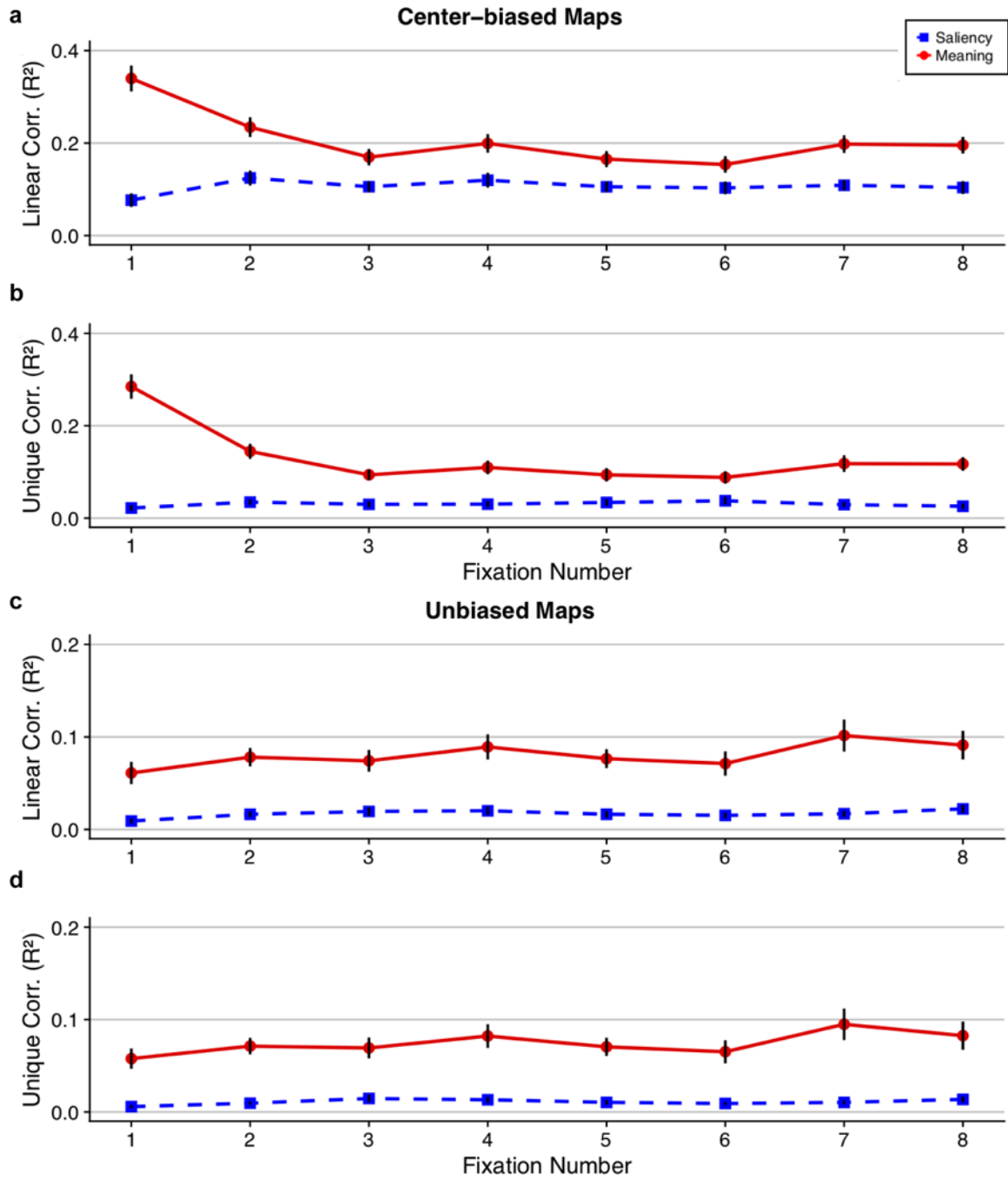


Figure 8. Fixation by fixation time-step analyses for the brightness rating task. The line plots show (a) the squared linear and (b) semi-partial correlations between fixation density and meaning (red circles) and saliency (blue squares) as a function of fixation number collapsed across scenes for the rating task using the center-biased maps. Line plots also show (c) the squared linear and (d) semi-partial correlations between fixation density and meaning (red circles) and saliency (blue squares) as a function of fixation order using the unbiased maps. Error bars represent the standard error of the mean.

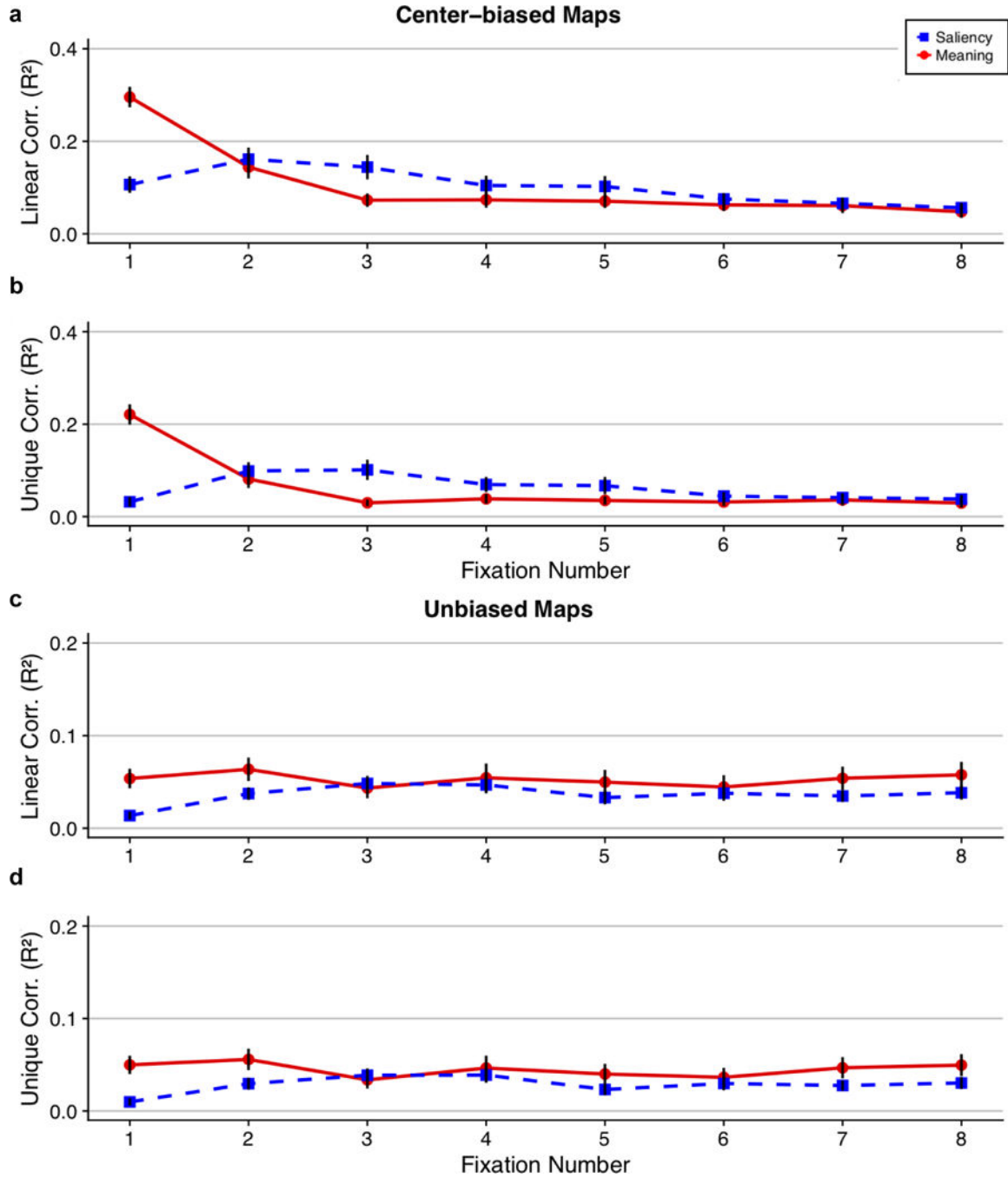


Figure 9. Fixation by fixation time-step analyses for the brightness search task. The line plots show (a) the squared linear and (b) semi-partial correlations between fixation density and meaning (red circles) and saliency (blue squares) as a function of fixation number collapsed across scenes for the search task using the center-biased maps. Line plots also show (c) the squared linear and (d) semi-partial correlations between fixation density and meaning (red circles) and saliency (blue squares) as a function of fixation order using the unbiased maps. Error bars represent the standard error of the mean.

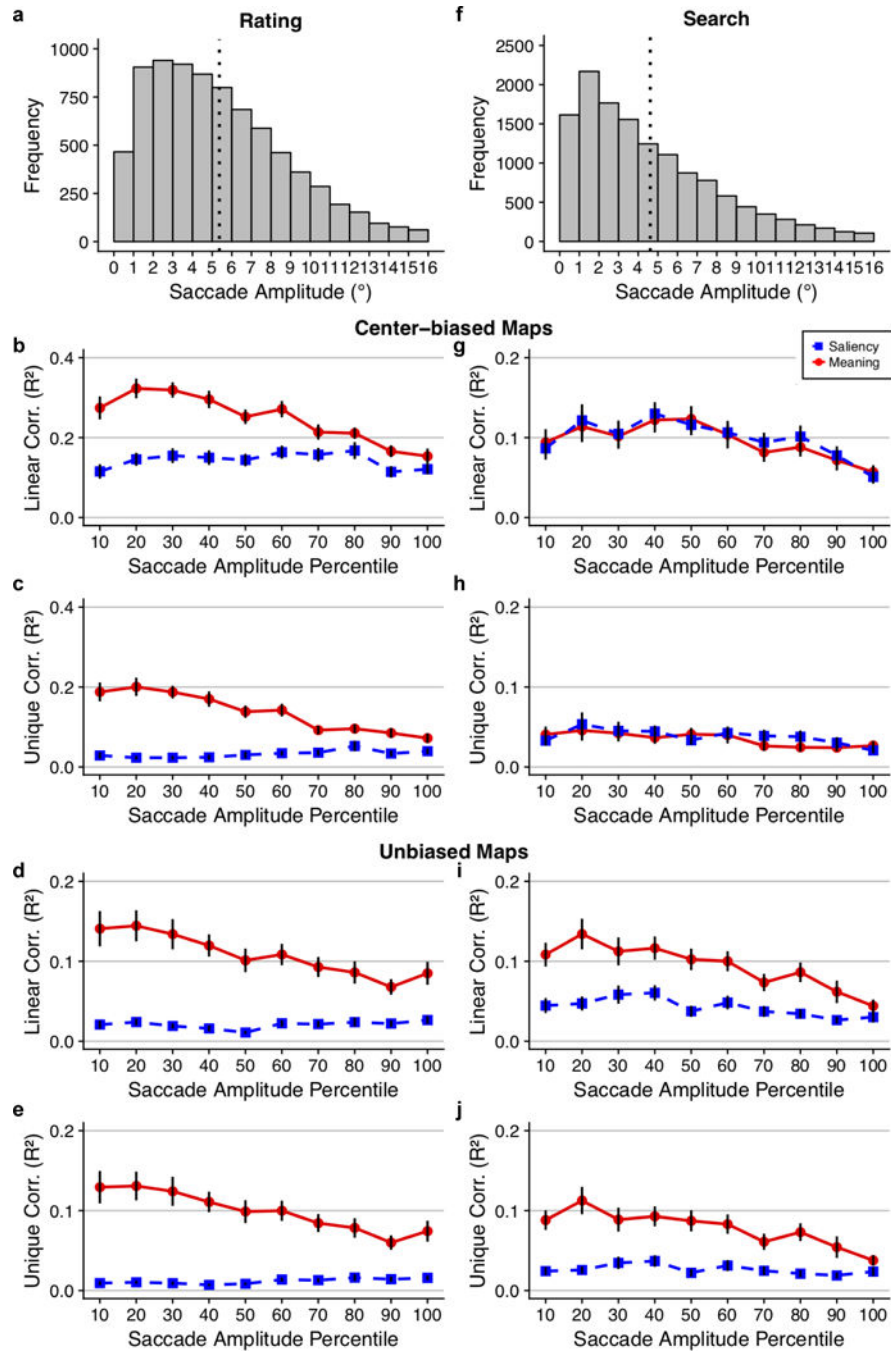


Figure 10. Squared linear correlation and squared semi-partial correlation as a function of saccade amplitude to fixation.

The saccade amplitude results for the rating task are shown in the first column (a through e) in which (a) shows a histogram of saccade amplitude frequencies and average saccade amplitude (black dotted line), (b) and (d) show the squared linear and (c) and (e) show the semi-partial correlations between meaning (red circles) and saliency (blue squares) and fixation density as a function of saccade amplitude percentiles prior to fixation for the center-biased maps (b and d) and the unbiased maps (c and e). The second column (f through j) shows the saccade amplitude results for the search condition in which (d) shows a

histogram of saccade amplitude frequencies and the average saccade amplitude (black dotted line), (g) and (i) show the squared linear correlations and (h) and (j) show the semi-partial correlations between meaning (red circles) and saliency (blue squares) with fixation density as a function of saccade amplitude percentiles using the center-biased maps (g and h) and the unbiased maps (i and j). Data points are averaged across all 40 scenes at each decile. Error bars represent standard error of the mean.