

Target-absent Human Attention

Zhibo Yang, Sounak Mondal, Seoyoung Ahn,
Gregory Zelinsky, Minh Hoai, and Dimitris Samaras

Stony Brook University, Stony Brook, NY 11794, USA

Abstract. The prediction of human gaze behavior is important for building human-computer interaction systems that can anticipate the user’s attention. Computer vision models have been developed to predict the fixations made by people as they search for target objects. But what about when the target is not in the image? Equally important is to know how people search when they cannot find a target, and when they would stop searching. In this paper, we propose a data-driven computational model that addresses the search-termination problem and predicts the scanpath of search fixations made by people searching for targets that do not appear in images. We model visual search as an imitation learning problem and represent the internal knowledge that the viewer acquires through fixations using a novel state representation that we call *Foveated Feature Maps (FFMs)*. FFMs integrate a simulated foveated retina into a pretrained ConvNet that produces an in-network feature pyramid, all with minimal computational overhead. Our method integrates FFMs as the state representation in inverse reinforcement learning. Experimentally, we improve the state of the art in predicting human target-absent search behavior on the COCO-Search18 dataset. Code is available at: <https://github.com/cvlab-stonybrook/Target-absent-Human-Attention>.

Keywords: Visual Search, Human Attention, Inverse Reinforcement Learning, Scanpath Prediction, Termination Prediction, Target Absent

1 Introduction

The attention mechanism used by humans to prioritize and select visual information [37,36,35] has attracted the interest of computer vision researchers seeking to reproduce this selection efficiency in machines [43,8,44,7,38]. The most often-used paradigm to study this efficiency is a visual search task, where efficiency is measured with respect to how many attention shifts (gaze fixations) are needed to detect a target in an image. But what about when the target is not there? Understanding gaze behavior during target-absent search (including search termination) would serve applications in human-computer interaction while addressing basic questions in attention research. No predictive model of human search fixations would be complete without addressing the unique problems arising from target-absent search.

The neuroanatomy of the primate foveated retina is such that visual acuity decreases with increasing distance from the high-resolution central fovea. When searching for a target, this foveated retina drives people to move their eyes selectively to image locations most likely to be the target, thereby providing the highest-resolution visual input to the target-recognition task, with each fixation movement guided by low-resolution input from peripheral vision. Recognizing the fact that the human visual input is filtered through a foveated retina is crucial to understanding and predicting human gaze behavior, and this is especially true for target-absent search where there is no clear target signal and gaze is driven instead by contextual relationships to other objects and the spatial cues that might provide about the target’s location.

To simulate a foveated retina for predicting human search fixations, Zelinsky *et al.* [44] directly applied a pretrained ResNet [16] to foveated images [34] to extract feature maps for the state representation. Yang *et al.* [43] proposed DCBs that approximate a high-resolution fovea and a low-resolution periphery by using the segmentation maps of a full-resolution image and its blurred version, respectively, predicted by a pretrained Panoptic-FPN [22]. Like other models for predicting human attention [31,25,26,7,46], both approaches rely on pretrained networks to extract image features and train much smaller networks for the downstream tasks using transfer learning, usually due to the lack of human fixation data for training. Also noteworthy is that these approaches apply networks pretrained on full-resolution images (e.g., ResNets [16] trained on ImageNet [39]) on blurred images, expecting the pretrained networks to approximate how humans perceive blurred images. However, Convolutional Neural Networks (ConvNets) are highly vulnerable to image perturbation [17,13] and the visual features extracted from the model on blurred images are hardly meaningful in the context of object recognition (contrary to human vision that actively seeks guidance from low-resolution peripheral vision for target recognition).

To better represent the degraded information that humans have available from their peripheral vision and can therefore use to guide their search, we exploit the fact that modern ConvNets have an inherent hierarchical architecture such that deeper layers have progressively larger receptive fields, corresponding to the greater blurring that occurs with increasing visual eccentricity. We propose combining the feature maps at different layers in a manner that is contingent upon the human fixation locations, approximating the information available from a foveated retina¹. We name this method *Foveated Feature Maps (FFMs)*. FFMs are computed on full-resolution images, so they can be readily applied to a wide range of pretrained ConvNets. Moreover, FFMs are a lightweight modification of modern ConvNets that capable of representing the subtle transition from fovea to periphery and are thus better suited for predicting human gaze movement. We find that our FFMs, when combined with inverse reinforcement learner (IQ-Learn [12]), significantly outperforms DCBs [43] and other baselines (see Sec. 4.3) in predicting both target-absent and target-present fixations.

¹ Note that it is not our aim to perfectly approximate the information extracted by a human foveated retina.

In short, our paper makes the following contributions: (1) we introduce a data-driven computational model applicable to both the target-present and target-absent search prediction problems; (2) we propose a new state representation that dynamically integrates knowledge collected via a foveated retina, similar to humans; (3) we predict target-absent search fixations at the ceiling of human performance, and achieve superior performance in predicting target-present scanpaths compared to previous methods; and (4) We propose a novel evaluation metric called semantic sequence score that measures the object-level consistency between human scanpaths. Compared to the traditional sequence score [4], it better captures the contextual cues that people use to guide their target-absent search behavior.

2 Related Work

Visual search is one of the fundamental human goal-directed gaze behaviors that actively scan the visual environment to find any exemplar of a target-object category [42,45,11]. There is an emerging interest in modeling and predicting human gaze during visual search [43,44,7,38,9]. Yang *et al.* [43] first used inverse reinforcement learning to model target-present search fixations spanning 18 target categories. Most recently, [7] directly applied reinforcement learning to predict scanpaths in various visual tasks including target-present search. However, their generalizability has never been interrogated for the prediction of target-absent search scanpaths, where no strong target signal is available in the images. Early work showed that target-absent search is not random behavior [10,2] but greatly influenced by target-relevant visual features to such an extent that the target category being searched for can be decoded from one’s scanpaths [47]. However, that study used only two target categories and the task was to search through only four non-targets. In this work, we study target-absent gaze behavior from a data-driven perspective.

Several recent studies have attempted to model a foveated representation of the input image for predicting human gaze behavior [43,44,38] or solving other visual tasks (e.g., object detection; [1,19]). Yang *et al.* [43] approximated a foveated retina by having a high-resolution center (full-resolution image) surrounded by a degraded visual periphery (a slightly blurred version of the image) at each fixation. A pretrained Panoptic-FPN [22] was applied on the full-resolution and blurred images separately to obtain the panoptic segmentation maps that were finally combined into the final state representation. Instead of approximating the foveated retina as a central-peripheral pairing of high- and low-resolution images, Zelinsky *et al.* [44] used a pretrained ResNet-50 [16] directly to extract feature maps from foveated images [34] for the state representation. Notably, both methods apply pretrained networks on blurred images whereas our FFMs are extracted from the full-resolution images, for which the pretrained networks are more robust. Rashidi *et al.* [38] proposed a method to directly estimate the foveated detectability of a target object [32] from eye tracking data. However, this approach cannot be easily extended to a larger number of target categories

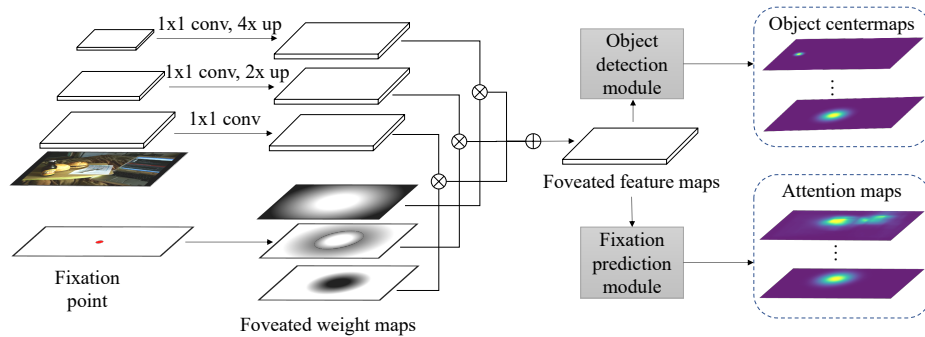


Fig. 1. Overview of the foveated feature maps (FFMs). FFMs are a set of multi-resolution feature maps constructed by combining the in-network feature pyramid produced by a pretrained ConvNet using the foveated weight maps computed based on previous fixations. The object detection module and the fixation prediction module map FFMs to a set of object center heatmaps (80 object categories in COCO [30]) and a set of attention maps for the 18 targets in COCO-Search18 [8], respectively.

because it requires training multiple detectors for each target and manually creating specialized datasets by showing each target at multiple scales against different textured backgrounds. In contrast, our model is able to jointly learn the foveation process at feature level and the networks that predict human scanpaths through back-propagation from human gaze behavior.

3 Approach

Following the model of Yang *et al.* [43] for target-present data, we also propose to model visual search behavior for target absent data using Inverse Reinforcement Learning (IRL). Specifically, we assume a human viewer is a reinforcement learning agent trying to localize the target object on a given *target-absent* image (the human viewer does not know if the image contains the target or not). The viewer acquires knowledge through a sequence of gaze fixations and allocates their next gaze point based on this knowledge to search for the target (Sec. 3.1). The search is terminated when the viewer confirms there is no target in the given image (Sec. 3.2). In this framework, we assume access to ground-truth human scanpaths (expert demonstrations), and the goal is to learn a policy that mimics or predicts human gaze behavior given an image and the target (Sec. 3.2).

3.1 Foveated Feature Maps (FFMs)

To capture the information a person acquires from an image through a sequence of fixations, we propose a novel state representation, called Foveated Feature Map (FFMs). Fig. 1 shows an overview of how our FFMs are constructed. FFMs take advantage of pretrained ConvNets, which produce a pyramid of feature

maps with progressively larger receptive fields. By treating deeper feature maps as information obtained at larger eccentricity (lower-resolution) in the peripheral vision, we construct FFMs as a set of multi-resolution feature maps, which is a weighted combination of different levels of feature maps using foveated weight maps generated based on previous fixations. Similar to image foveation [44], in FFMs, deeper feature maps with a lower resolution (corresponding to a larger eccentricity) are more weighted at locations with increasing distance from the fixation points. Below we discuss FFMs in greater detail.

Relative resolution map. The human vision system is known to be foveated, meaning the visual information in the view is not processed at a uniform resolution. Rather, high spatial details are only obtained around the fixation point (i.e., the fovea) and the resolution outside of the fovea drops off as the distance between the peripheral pixels and the fovea increases. To simulate this, Perry and Geisler [34] proposed an image foveation method, which has been used in both free viewing [20] and visual search [44] tasks. Here, we extend image foveation to produce multi-resolucional feature maps to represent the foveated view of an image at the level of image features. Specifically, given a fixation $f = (x_f, y_f)$, we first define a relative resolution map contingent on f as

$$R(x, y|f) = \frac{\alpha}{\alpha + \frac{\sqrt{(x-x_f)^2+(y-y_f)^2}}{p}}. \quad (1)$$

Here, p is the number of pixels in one degree of visual angle, which depends on the distance between the viewer and the display. α is a learnable parameter that controls the decreasing speed of resolution as the pixel (x, y) moves away from the fixation point.

For multiple fixations $\{f_1, \dots, f_n\}$, we compute the combined resolution map by taking the maximum at every location: $R(x, y|\{f_1, \dots, f_n\}) = \max_i R(x, y|f_i)$. In contrast to [34], which creates a Gaussian pyramid of the given image I to produce the multi-resolucional version of I , we take inspiration from the Feature Pyramid Network [28] and use the in-network feature pyramid produced by existing pretrained ConvNets and blend the feature maps at each level of the feature pyramid to construct multi-resolucional feature maps (i.e., FFMs) based on the relative resolution map $R(x, y|f)$. For brevity, we will write $R(x, y|f)$ as $R(x, y)$ in the following text.

Foveated feature maps (FFMs). We use a ResNet-50 [16] as the backbone (the method can be easily extended to other ConvNet backbones such as VGG nets [40]), and let the feature pyramid from the ResNet be $\{C_1 \dots, C_5\}$, which represents the feature activation outputs from the last residual block at each stage of ResNet-50, namely the outputs of conv1, conv2, conv3, conv4, and conv5. Similar to the Gaussian pyramid of an image, a lower level of the feature pyramid contains more spatial details, while a higher-level feature map is stronger in semantics. To reduce the semantic discrepancy among different levels, we apply an 1×1 convolutional layer on every C_i to project them to the same embedding space. Then, we upsample $\{C_i\}_{i=1}^5$ to the same spatial dimensions of C_1 , yielding

3D tensors of the same size, denoted as $\{P_1, \dots, P_5\}$. We then compute a spatial weight map W_i for each P_i and produce a set of multi-resolution feature maps M as the weighted combination of W_i and P_i : $M = \sum_i W_i \odot P_i$, where \odot denotes the element-wise multiplication at the spatial axes. We call these multi-resolution feature maps FFMs. Below we describe how to compute W_i based on the relative resolution map $R(x, y)$.

Each level of the feature pyramid P_i represents a certain eccentricity, corresponding to a fixed spatial resolution, which we denote as R_i^* . It is defined as the relative resolution where a transfer function $T_i(\cdot)$ is at its half maximum, i.e., $T_i(R_i^*) = 0.5$ [34]. The transfer function $T_i(\cdot)$ is the function that maps relative resolution r to relative amplitude, and it is defined as:

$$T_i(r) = \exp(-(2^{i-3}r/\sigma)^2/2). \quad (2)$$

It can be shown that $R_1^* > R_2^* > R_3^* > R_4^* > R_5^*$, forming four resolution bins whose boundaries are defined by R_i^* and R_{i-1}^* ($i \in \{2, 3, 4, 5\}$). To compute the weights at location (x, y) , we first determine which bin pixel (x, y) falls in, according to its relative resolution $R(x, y)$ (see the supplementary material for more details). Assume pixel (x, y) falls in between layer j and $j - 1$, i.e., $R_{j-1}^* \geq R(x, y) > R_j^*$. Then, we set the weights at layer j and $j - 1$ to be the ratio of the distance between pixel (x, y) and the corresponding layer to the distance between the layer j and $j - 1$ at (x, y) in relative amplitude space:

$$W_i(x, y) = \begin{cases} \frac{0.5 - T_j(R(x, y))}{T_{j-1}(R(x, y)) - T_j(R(x, y))} & \text{if } i = j - 1, \\ 1 - \frac{0.5 - T_j(R(x, y))}{T_{j-1}(R(x, y)) - T_j(R(x, y))} & \text{if } i = j, \\ 0 & \text{otherwise.} \end{cases} \quad (3)$$

Apparently, $\sum_i W_i(x, y) = 1$ and at location (x, y) only features from layer j and layer $j - 1$ are integrated into the final FFMs. In [34], α is tuned to match human perception via physiological experiments. Here we learn the parameters of FFMs, α and σ , together with the policy from human gaze data directly.

3.2 Reward and Policy Learning

Using FFMs as our state representation, we train a policy that mimics human gaze behavior using the IRL framework [43]. However, we found that the GAIL [18] IRL algorithm used in [43] is too sensitive to its hyper-parameters, due to its adversarial learning design, which is also shown in [23]. We therefore use IQ-Learn [12] as our IRL algorithm instead. Based on soft Q-Learning [14], IQ-Learn encodes both the reward and the policy in a single Q-function, and thus is able to optimize both reward and policy simultaneously.

Let $Q(s, a)$ be the Q-function, which maps a state-action pair (s, a) to a scalar value representing the amount of future reward gained by taking action a under state s . We want to find a reward function that maximizes the expected amount of cumulative rewards that the expert policy obtains over all other possible

policies. Hence, IQ-Learn trains the Q-function by minimizing the following loss:

$$\mathcal{L}_{\text{irl}} = -\mathbb{E}_{\rho_E} [Q(s, a) - \gamma \mathbb{E}_{s' \sim \mathcal{P}(s, a)} V(s')], \quad (4)$$

where $V(s) = \log \sum_a \exp(Q(s, a))$, ρ_E and \mathcal{P} denote the occupancy measure of the expert policy [18] and the dynamics, respectively. We do not apply the χ^2 -divergence proposed in [12] on the reward function since it did not lead to any notable improvement on our task. Given the learned Q-function Q , we can compute the reward as a function of the state and action:

$$r(s, a) = Q(s, a) - \gamma \mathbb{E}_{s' \sim \mathcal{P}(s, a)} V(s'), \quad (5)$$

and the policy as a function of the state:

$$\pi(a|s) = \frac{\exp(Q(s, a)/\tau)}{\sum_{a'} \exp(Q(s, a')/\tau)}. \quad (6)$$

τ is the temperature coefficient, controlling the entropy of the action distribution.

Action space. Our task is to predict the next fixation given the previous fixations, the input image, and the categorical target. To predict fixations on an image, we follow [43] and discretize the image space into a 20×32 grid (action space). At each time step, the policy samples one cell out of 640 grid cells according to the predicted categorical action distribution $\pi(\cdot|s)$. For the selected grid cell, we set the predicted fixation to be the center of the cell.

Auxiliary detection task. A visual search task is essentially a detection task, so it is important for the state representation to capture features of the target object. Moreover, in target-absent search where the target object is absent, human behavior is driven by the expected location of the target in relation to other commonly co-occurring objects. In contrast to [43] which directly uses the output of a pretrained panoptic segmentation network, we train the Q-function with an auxiliary task of predicting the center maps of the objects. Specifically, we add a detection network module on top of FFMs. This module outputs 80 heatmaps \hat{Y} for the 80 object categories in the COCO dataset [30]. Let \hat{Y}_{xyc} denote the value of the c -th heatmap at location (x, y) . Following CenterNet [48], we use pixel-wise focal loss [29] as an additional loss to train the whole network:

$$\mathcal{L}_{\text{det}} = -\frac{1}{N} \sum_{x, y, c} \begin{cases} (1 - \hat{Y}_{xyc})^\kappa \log(\hat{Y}_{xyc}) & \text{if } Y_{xyc} = 1, \\ (1 - Y_{xyc})^\lambda (\hat{Y}_{xyc})^\kappa \log(1 - \hat{Y}_{xyc}) & \text{otherwise,} \end{cases} \quad (7)$$

where Y is the ground-truth heatmap created by an object size dependent Gaussian kernel [27]. We set $\kappa = 2$ and $\lambda = 4$ as in [48]. Note that we do not predict the exact heights and widths of the objects in the image because we think rough estimates of the locations of different objects are sufficient to help predict the target-absent fixations. We learn the Q-function using both the IRL loss and the auxiliary detection loss:

$$\mathcal{L} = \mathcal{L}_{\text{irl}} + \omega \mathcal{L}_{\text{det}}, \quad (8)$$

where ω is a weight to balance the two loss terms.

Termination Prediction. When a person will stop searching is a question intrinsic to target-absent search. Different from [7], which formulates termination as an extra action to fixation prediction in policy learning, we treat termination prediction as an additional task that occurs every step after a new fixation has been made. We found that if we treat termination as an extra action, the policy would overfit to the termination action as it appears much more frequent than other actions.

To this end, we train a binary classifier on top of the Q-function (see Sec. 3.2) for termination prediction using binary cross entropy loss. We weigh the loss computed on the termination and non-termination actions inversely proportionally to their frequencies. In addition, psychology studies [10,41] have suggested that time could be an important ingredient in predicting stopping. However, we do not predict the duration of fixations in our model. Instead, we use the number of previous fixations as an approximation of time and concatenate it with the Q-values from the Q-function as input to train the termination classifier.

4 Experiments

We train and evaluate the proposed method and other models by using COCO-Search18 [8], which contains both target-present and target-absent human scanpaths in searching for 18 different object categories. COCO-Search18 has 3101 target-present images and 3101 target-absent images, each viewed by 10 subjects. In this paper, we mainly focus on the target-absent gaze behavior prediction. All models are only trained with target-absent images and fixations unless otherwise specified. For all models, we predict one scanpath for each testing image in a greedy manner (i.e., always selecting the action with the largest probability mass from the predicted action distribution as the next fixation) and compare them with the ground-truth scanpaths.

4.1 Semantic Sequence Score

The sequence score (SS) has often been used to quantify the success of scanpath prediction [4,43]. The sequence score is computed by an existing string matching algorithm that compares the two fixation sequences [33] after transforming them into strings of fixation cluster IDs. The fixation clusters are computed based on the fixation locations. However, we argue that the sequence score does not capture the semantic meaning of fixations which plays an important role in analyzing goal-directed attention: it only captures “where” a person is looking at but not “what” is being looked at. To this end, we propose the *Semantic Sequence Score (SemSS)*, which transforms a fixation sequence into an *object category* sequence by leveraging the segmentation annotation provided in COCO [30]. Then, we apply the same string matching algorithm used in the traditional sequence score to measure the similarity between two scanpaths. Using the “things” versus “stuff” paradigm [6], we do not distinguish between object instances. In this paper, we

focus on “thing” categories only, as we are interested in how non-target objects collectively affect human gaze behavior in visual search tasks. “Stuff” categories can be easily integrated into the semantic sequence score.

Other metrics. We also report other scanpath prediction metrics including the traditional sequence score and conditional priority maps [24], which measure how well the model predicts a fixation when given the previous fixations using saliency metrics including information gain (IG) and normalized scanpath saliency (NSS) [5]. For clarity, we denote them by cIG and cNSS where “c” represents “conditional”. cIG measures the amount of information gain the model prediction has over a task-specific fixation density map computed using the training fixations. cNSS measures the correspondence between the predicted fixation probability map and the ground-truth fixation. In addition, to measure termination prediction accuracy, we report the Mean Absolute Error (MAE) between predicted and ground-truth scanpath lengths. To compare fairly with models that do not terminate automatically such as IRL [43], we also report the truncated sequence score by truncating predicted and ground-truth scanpaths at the first 2 and 4 new fixations, denoted as SS(2) and SS(4), respectively.

4.2 Implementation Details

Network structure. Following [43], we resize the input images to 320×512 for computational efficiency. As shown in Fig. 1, our model has three components: a set of 1×1 convolutional layers that project the feature maps in the feature pyramid to the same dimension (i.e., the number of channels in FFMs); an object detection module; and a fixation prediction module. We set the number of FFMs channels to 128. The fixation prediction module and the object detection module share the same ConvNet consisting of three consecutive convolutional blocks which reduce the spatial resolution of the input foveated feature maps (FFMs) by a factor of 8 (from 160×256 to 20×32). In between two consecutive convolutional layers of a convolutional block, we apply Layer Normalization [3] and a ReLU activation function. Finally, the fixation prediction module uses two convolutional layers to map the output of the shared ConvNet into 18 attention maps (one for each target in COCO-Search18 [8]). The object detection module has a similar structure, but outputs 80 center maps (one for each object category in COCO [30]). Note that the backbone networks of all models in this paper are kept fixed during training. Detailed network parameters are in supplementary.

Hyperparameters. We train the models in this paper by using the Adam [21] optimizer with learning rate 10^{-4} . The weight for the auxiliary detection loss ω in Eq. (8) is 0.1. In COCO-Search18 [8], the number of pixels in one degree of visual angle $p = 9.14$. We scale it according to the spatial resolution of P_1 and set $p = 4.57$. For models with a termination predictor, we set the maximum length of each predicted scanpath to 10 (excluding the initial fixation) during training and testing. For models that do not terminate automatically, we set the length of the scanpath to 6 which is approximately the average length of the target-absent scanpaths in COCO-Search18. For the IQ-Learn algorithm, the

Table 1. Comparing target-absent scanpath prediction algorithms (rows) using multiple scanpath metrics (columns) on the target-absent test set of COCO-Search18. The best results are highlighted in bold.

	SemSS	SS	cIG	cNSS	SS(2)	SS(4)
Human consistency	0.542	0.381	-	-	0.561	0.478
Detector	0.497	0.321	-0.516	0.446	0.497	0.402
Fixation heuristic	0.484	0.298	-0.599	0.405	0.492	0.379
IRL [43]	0.476	0.319	0.032	1.202	0.508	0.407
Chen <i>et al.</i> [7]	0.484	0.331	-	-	0.516	0.434
Ours	0.516	0.372	0.729	1.524	0.537	0.441

reward discount factor is set to 0.8. Following [15,12], we use target updates and a replay buffer in IQ-Learn to stabilize the training. The temperature coefficient τ in Eq. (5) is set to 0.01. We update the target Q networks for four iterations using exponential moving average with a 0.01 coefficient. The replay buffer can hold 8000 state-action pairs and is updated online during training.

4.3 Comparing Scanpath Prediction Methods

We compare our model with the following baselines: 1) *human consistency*, an oracle method where one searcher’s scanpath is used to predict another searcher’s scanpath; 2) *detector*, a ConvNet trained on target-present images of COCO-Search18 to output a target detection confidence map, from which we sample fixations sequentially with inhibition of return (IOR); 3) *fixation heuristic*, similar to detector, but trained to predict human fixation density maps using target-absent data; and more recent approaches including 4) *IRL* [43], and 5) *Chen et al.’s model* [7]. Note that Chen *et al.*’s model used a finer action space 30×40 . For fair comparison, we rescale its predicted fixations to our action space 20×32 .

As can be seen from Tab. 1, our method outperforms all other methods across all metrics in target-absent scanpath prediction². Our method is the closest to human consistency which is regarded as the ceiling of any predictive model. In the sequence score case, our method is only inferior to human consistency by 0.09, leading the second best (Chen *et al.* [7]) by 0.41. Excluding the effect of the termination predictor, the sequence scores of the first 2 and 4 fixations also show that even without terminating the scanpaths our method is still the best compared to all other computational models. Moreover, comparing the sequence scores of truncated scanpaths and full scanpaths, we see a trend of decreasing performance as the scanpath length increases for all methods, i.e., $SS(2) > SS(4) > SS$, and this pattern is particularly pronounced in target-absent search (there is no significant difference between SS and SS(4) for target-present search, see Tab. 3). The fact that later fixations during target-absent search are harder to

² Both cIG and cNSS can only be computed for auto-regressive probabilistic models (our method, IRL, detector and fixation heuristic).

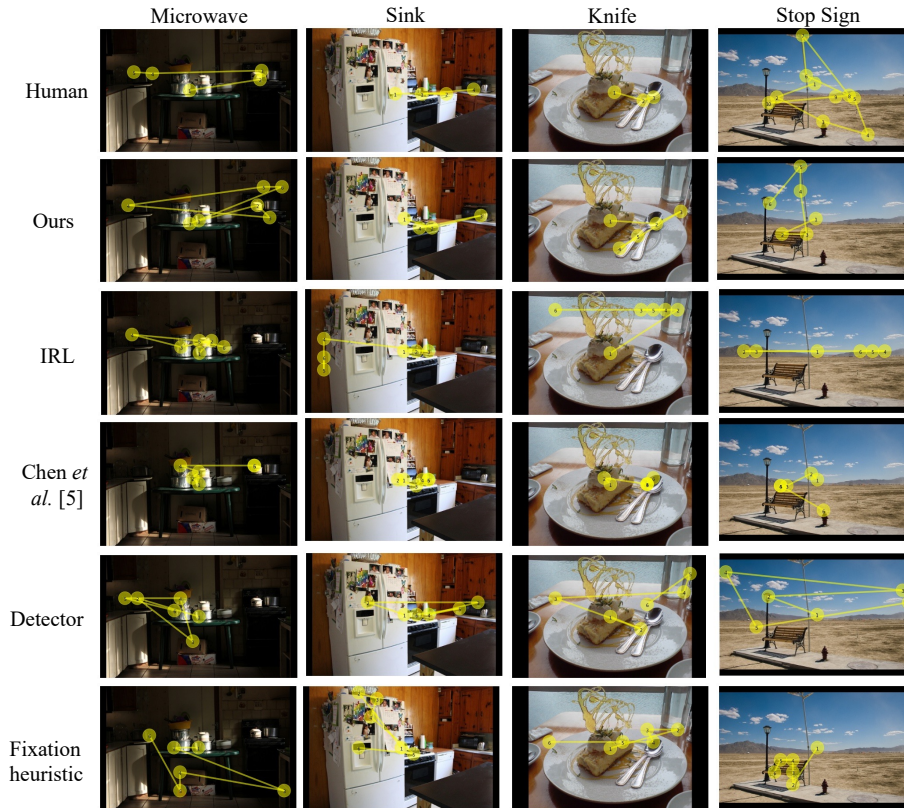


Fig. 2. Visualizing the **predicted scanpaths** of different methods (rows) for different search targets (columns). The top row shows the ground-truth human scanpaths and the other rows are predicted scanpaths from different models.

predict suggests that human eye-movements behave more randomly at the later stage of search especially when there is no target in the scene.

We also qualitatively compare different methods by visualizing their predicted scanpaths for four scenes in Fig. 2. When searching for a microwave in this scene, our method alone predicted fixations on all three table and counter-top surfaces in the image where microwaves are often found (similar to how a representative human searched). Similar phenomena are observed for the sink and knife searches. This shows that our method is able to capture the contextual relations between objects that play a role in driving target-absent fixations. When searching for the stop sign, our method was the only one that looked at the top of the centrally-located vertical object, despite heavy occlusion, speculatively because stop signs are usually mounted to the tops of poles. In contrast, IRL, which extracts features from blurred pixels using a pretrained ConvNet, completely failed to capture the vertical objects in this image that seem to be

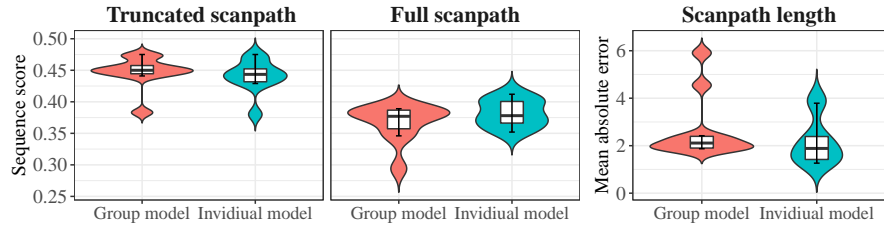


Fig. 3. Comparing group model (red) and individual model (cyan) using: (left) the sequence score of the truncated scanpath (first 4 fixations) without automatic termination, and (middle) the sequence score of the full scanpath including termination, and (right) the mean absolute error for the predicted scanpath length. We perform Wilcoxon signed-rank tests for each experimental setting. The two-sided p values are 0.012, 0.028 and 0.006, respectively.

guiding search. This argues for the value in using our proposed FFMs to capture guiding contextual information extracted from peripheral vision.

4.4 Group Model versus Individual Model

In target-present search, human scanpaths are very consistent due to the strong guidance provided by the target object in the image. Indeed, a model trained with fixations from a group of people generalized well for a new **unseen** person [43]. However, given that there are large individual differences in termination time for target-absent search [8], we expect that individualized modeling may be necessary for target-absent search prediction. To test this hypothesis, we compared the predictive performance of group versus individual modeling of target-absent search fixations. The group model was trained with 9 subjects’ training scanpaths and tested on the testing scanpaths of the remaining subject. The individual model was trained with the training scanpaths of a single subject and tested on the same subject’s testing scanpaths. We did this for all 10 subjects.

Fig. 3 shows the comparison between the group model and the individual model in the sequence score of full scanpaths and truncated scanpaths (first four fixations) and the MAE of the length of the predicted scanpath. Interestingly, despite being trained with less data, the individual model shows better performance than the group model in full scanpath modeling, contrary to the group model being better in the truncated scanpath prediction. A critical difference between the modeling of truncated versus full scanpath is that the latter involves the search termination prediction. The rightmost graph in Fig. 3 also shows that the individual model generates less error (in MAE metric) in scanpath length prediction than the group model. These results altogether suggest that individualized modeling may be more suitable for target-absent search prediction. More experimental results on the termination criterion across different subjects can be found in supplementary.

Table 2. Ablation study. We ablate the loss function (second row) and the state representation (third and fourth rows). All methods are trained using IQ-Learn.

	SemSS	SS	cIG	cNSS	SS(2)	SS(4)
FFMs	0.516	0.372	0.729	1.524	0.537	0.441
FFMs w/o detection loss	0.476	0.350	0.550	1.332	0.545	0.437
DCBs	0.508	0.355	0.212	1.129	0.514	0.426
CFI	0.504	0.352	0.518	1.252	0.506	0.426
FPN	0.508	0.338	0.018	0.881	0.408	0.351
Binary masks	0.510	0.364	0.347	1.148	0.438	0.378

Table 3. Comparing target-present scanpath prediction algorithms using multiple scanpath metrics on the COCO-Search18 test dataset.

	SemSS	SS	cIG	cNSS	SS(2)	SS(4)
Human consistency	0.624	0.478	-	-	0.486	0.480
IRL [43]	0.536	0.419	-9.709	1.977	0.437	0.421
Chen <i>et al.</i> [7]	0.572	0.445	-	-	0.429	0.319
Ours	0.562	0.451	1.548	2.376	0.467	0.450

4.5 Ablation Study

First, we ablate the loss (see Eq. (8)) of our model by removing the auxiliary detection loss. Second, we ablate our proposed foveated feature maps (FFMs) by comparing it with dynamic contextual beliefs (DCBs) [43] and cumulative foveated image (CFI) [44] using the same IRL algorithm (i.e., IQ-Learn). As a finer-grained ablation, we ablate FFMs by using the features extracted by the FPN backbone of a COCO-pretrained Mask R-CNN as the state representation. We use the highest-resolution feature maps of FPN P_2 . We further binarize the FFMs of our model such that the values are one at the fixated locations of the finest level (fovea) and the non-fixated locations of the coarsest level (periphery) and zero elsewhere. As shown in Tab. 2, the proposed auxiliary detection loss improves the performance in 5 out of 6 metrics. The semantic sequence score is increased from 0.476 to 0.516, which indicates that knowing the locations of non-target objects in the image is helpful for predicting the target-absent fixations. Comparing different state representations (i.e., FFMs, DCBs, CFI, FPN and binary masks), we can see that the proposed FFMs are superior to all other state presentations in predicting target-absent fixations. This shows the superiority of FFMs in representing the knowledge a human acquires through fixations compared to DCBs and CFI, which apply pretrained ConvNets on blurred images to simulate the foveated retina.

4.6 Generalization to Target-present Search

Despite being motivated by target-absent search, our method is also directly applicable to target-present fixation prediction. In this section, we compare our model with two competitive models, IRL [43] and Chen *et al.* [7], in target-present scanpath prediction. For fair comparison, we follow [43] and set the maximum scanpath length to be 6 (excluding the first fixation) for all models and automatically terminate the scanpath once the fixation falls in the bounding box of the target. Tab. 3 shows that our method achieves the best performance in 5 out of 6 metrics. Chen *et al.*'s model is slightly better than ours in semantic sequence score. They used a pretrained CenterNet [48] trained on COCO images [30] (about 118K images) to predict the bounding box of the target as input for their model, whereas we only used the target-present images in COCO-Search18 [8] (about 3K images) to train our object detection module (see Fig. 1). Despite being trained with less data, our model still outperforms Chen *et al.* [7] in the other five metrics, especially when evaluated in truncated fixed-length scanpaths (i.e., SS(2) and SS(4)). We further expect our model to perform better when using all COCO training images to train our object detection module. Tab. 1 and Tab. 3 together demonstrate that our proposed method not only excels in predicting target-absent fixations (see Sec. 4.3), but also target-present fixations.

5 Conclusions and Discussion

We have presented the first computational model for predicting target-absent search scanpaths. To represent the internal knowledge that the viewer acquires through fixations, we proposed a novel state representation, *foveated feature maps (FFMs)*. FFMs circumvent the drawbacks of directly applying pretrained ConvNets on blurred images in previous methods [44,43] by integrating the in-network feature pyramid produced by a pretrained ConvNet with a foveated retina. When trained and evaluated on the COCO-Search18 dataset, FFMs outperform previous state representations and achieve state-of-the-art performance in predicting both target-absent and target-present search fixations using the IRL framework. Moreover, we also proposed a new variant of the sequence score for measuring scanpath similarity, called semantic sequence score. It better captures the object-to-object relation used to guide target-absent search.

Future work. Inspired by [43], our future work will involve extending our model and semantic sequence score to include “stuff” categories in COCO [6] to study the impact of background categories to target-absent search gaze behavior, and exploring using semi-supervised learning to address the lack of human gaze data by leveraging the rich annotation in COCO images [30].

Acknowledgements. The authors would like to thank Jianyuan Deng for her help in result visualization and statistical analysis. This project was partially supported by US National Science Foundation Awards IIS-1763981 and IIS-2123920, the Partner University Fund, the SUNY2020 Infrastructure Transportation Security Center, and a gift from Adobe.

References

1. Akbas, E., Eckstein, M.P.: Object detection through search with a foveated visual system. *PLoS computational biology* (2017)
2. Alexander, R.G., Zelinsky, G.J.: Visual similarity effects in categorical search. *Journal of vision* **11**(8), 9–9 (2011)
3. Ba, J.L., Kiros, J.R., Hinton, G.E.: Layer normalization. *arXiv preprint arXiv:1607.06450* (2016)
4. Borji, A., Tavakoli, H.R., Sihite, D.N., Itti, L.: Analysis of scores, datasets, and models in visual saliency prediction. In: *Proceedings of the IEEE international conference on computer vision*. pp. 921–928 (2013)
5. Bylinskii, Z., Judd, T., Oliva, A., Torralba, A., Durand, F.: What do different evaluation metrics tell us about saliency models? *IEEE transactions on pattern analysis and machine intelligence* **41**(3), 740–757 (2018)
6. Caesar, H., Uijlings, J., Ferrari, V.: Coco-stuff: Thing and stuff classes in context. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 1209–1218 (2018)
7. Chen, X., Jiang, M., Zhao, Q.: Predicting human scanpaths in visual question answering. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 10876–10885 (2021)
8. Chen, Y., Yang, Z., Ahn, S., Samaras, D., Hoai, M., Zelinsky, G.: Coco-search18 fixation dataset for predicting goal-directed attention control. *Scientific reports* **11**(1), 1–11 (2021)
9. Chen, Y., Yang, Z., Chakraborty, S., Mondal, S., Ahn, S., Samaras, D., Hoai, M., Zelinsky, G.: Characterizing target-absent human attention. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*. pp. 5031–5040 (2022)
10. Chun, M.M., Wolfe, J.M.: Just say no: How are visual searches terminated when there is no target present? *Cognitive psychology* **30**(1), 39–78 (1996)
11. Eckstein, M.P.: Visual search: A retrospective. *Journal of vision* **11**(5), 14–14 (2011)
12. Garg, D., Chakraborty, S., Cundy, C., Song, J., Ermon, S.: Iq-learn: Inverse soft-q learning for imitation. *Advances in Neural Information Processing Systems* **34** (2021)
13. Geirhos, R., Rubisch, P., Michaelis, C., Bethge, M., Wichmann, F.A., Brendel, W.: Imagenet-trained CNNs are biased towards texture; increasing shape bias improves accuracy and robustness. In: *International Conference on Learning Representations* (2019), <https://openreview.net/forum?id=Bygh9j09KX>
14. Haarnoja, T., Tang, H., Abbeel, P., Levine, S.: Reinforcement learning with deep energy-based policies. In: *International Conference on Machine Learning*. pp. 1352–1361. PMLR (2017)
15. Haarnoja, T., Zhou, A., Abbeel, P., Levine, S.: Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In: *International conference on machine learning*. pp. 1861–1870. PMLR (2018)
16. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 770–778 (2016)
17. Hendrycks, D., Dietterich, T.: Benchmarking neural network robustness to common corruptions and perturbations. In: *International Conference on Learning Representations* (2019), <https://openreview.net/forum?id=HJz6tiCqYm>

18. Ho, J., Ermon, S.: Generative adversarial imitation learning. *Advances in neural information processing systems* **29** (2016)
19. Jaramillo-Avila, U., Anderson, S.R.: Foveated image processing for faster object detection and recognition in embedded systems using deep convolutional neural networks. In: *Conference on Biomimetic and Biohybrid Systems*. pp. 193–204. Springer (2019)
20. Jiang, M., Huang, S., Duan, J., Zhao, Q.: Salicon: Saliency in context. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (June 2015)
21. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: *ICLR (Poster)* (2015)
22. Kirillov, A., Girshick, R., He, K., Dollár, P.: Panoptic feature pyramid networks. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 6399–6408 (2019)
23. Kostrikov, I., Agrawal, K.K., Dwibedi, D., Levine, S., Tompson, J.: Discriminator-actor-critic: Addressing sample inefficiency and reward bias in adversarial imitation learning. In: *International Conference on Learning Representations* (2019), <https://openreview.net/forum?id=Hk4fpoA5Km>
24. Kümmerer, M., Bethge, M.: State-of-the-art in human scanpath prediction. *arXiv preprint arXiv:2102.12239* (2021)
25. Kümmerer, M., Theis, L., Bethge, M.: Deep gaze i: Boosting saliency prediction with feature maps trained on imagenet. *arXiv preprint arXiv:1411.1045* (2014)
26. Kümmerer, M., Wallis, T.S., Bethge, M.: Information-theoretic model comparison unifies saliency metrics. *Proceedings of the National Academy of Sciences* **112**(52), 16054–16059 (2015)
27. Law, H., Deng, J.: Cornernet: Detecting objects as paired keypoints. In: *Proceedings of the European conference on computer vision (ECCV)*. pp. 734–750 (2018)
28. Lin, T.Y., Dollár, P., Girshick, R., He, K., Hariharan, B., Belongie, S.: Feature pyramid networks for object detection. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 2117–2125 (2017)
29. Lin, T.Y., Goyal, P., Girshick, R., He, K., Dollár, P.: Focal loss for dense object detection. In: *Proceedings of the IEEE international conference on computer vision*. pp. 2980–2988 (2017)
30. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: *European conference on computer vision*. pp. 740–755. Springer (2014)
31. Linardos, A., Kümmerer, M., Press, O., Bethge, M.: Deepgaze iie: Calibrated prediction in and out-of-domain for state-of-the-art saliency modeling. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 12919–12928 (2021)
32. Najemnik, J., Geisler, W.S.: Optimal eye movement strategies in visual search. *Nature* **434**(7031), 387–391 (2005)
33. Needleman, S.B., Wunsch, C.D.: A general method applicable to the search for similarities in the amino acid sequence of two proteins. *Journal of molecular biology* **48**(3), 443–453 (1970)
34. Perry, J.S., Geisler, W.S.: Gaze-contingent real-time simulation of arbitrary visual fields. In: *Human vision and electronic imaging VII*. vol. 4662, pp. 57–70. International Society for Optics and Photonics (2002)
35. Petersen, S.E., Posner, M.I.: The attention system of the human brain: 20 years after. *Annual review of neuroscience* **35**, 73–89 (2012)

36. Posner, M.I.: Attention: the mechanisms of consciousness. *Proceedings of the National Academy of Sciences* **91**(16), 7398–7403 (1994)
37. Posner, M.I., Petersen, S.E.: The attention system of the human brain. *Annual review of neuroscience* **13**(1), 25–42 (1990)
38. Rashidi, S., Ehinger, K., Turpin, A., Kulik, L.: Optimal visual search based on a model of target detectability in natural images. *Advances in Neural Information Processing Systems* **33** (2020)
39. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., et al.: Imagenet large scale visual recognition challenge. *International journal of computer vision* **115**(3), 211–252 (2015)
40. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014)
41. Wolfe, J.M.: Guided search 6.0: An updated model of visual search. *Psychonomic Bulletin & Review* **28**(4), 1060–1092 (2021)
42. Wolfe, J.: Visual search. pashler, h.(ed.), *attention* (1998)
43. Yang, Z., Huang, L., Chen, Y., Wei, Z., Ahn, S., Zelinsky, G., Samaras, D., Hoai, M.: Predicting goal-directed human attention using inverse reinforcement learning. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. pp. 193–202 (2020)
44. Zelinsky, G., Yang, Z., Huang, L., Chen, Y., Ahn, S., Wei, Z., Adeli, H., Samaras, D., Hoai, M.: Benchmarking gaze prediction for categorical visual search. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*. pp. 0–0 (2019)
45. Zelinsky, G.J.: A theory of eye movements during target acquisition. *Psychological review* **115**(4), 787 (2008)
46. Zelinsky, G.J., Chen, Y., Ahn, S., Adeli, H., Yang, Z., Huang, L., Samaras, D., Hoai, M.: Predicting goal-directed attention control using inverse-reinforcement learning. *Neurons, behavior, data analysis and theory* **2021** (2021)
47. Zelinsky, G.J., Peng, Y., Samaras, D.: Eye can read your mind: Decoding gaze fixations to reveal categorical search targets. *Journal of vision* **13**(14), 10–10 (2013)
48. Zhou, X., Wang, D., Krähenbühl, P.: Objects as points. *arXiv preprint arXiv:1904.07850* (2019)