

# What am I Searching for: Zero-shot Target Identity Inference in Visual Search

Mengmi Zhang<sup>1,2</sup> and Gabriel Kreiman<sup>1,2</sup>

{mengmi.zhang@childrens, gabriel.kreiman@tch}.harvard.edu

<sup>1</sup>Children’s Hospital, Harvard Medical School

<sup>2</sup>Center for Brains, Minds and Machines

## Abstract

Can we infer intentions from a person’s actions? As an example problem, here we consider how to decipher what a person is searching for by decoding their eye movement behavior. We conducted two psychophysics experiments where we monitored eye movements while subjects searched for a target object. We defined the fixations falling on non-target objects as “error fixations”. Using those error fixations, we developed a model (InferNet) to infer what the target was. InferNet uses a pre-trained convolutional neural network to extract features from the error fixations and computes a similarity map between the error fixations and all locations across the search image. The model consolidates the similarity maps across layers and integrates these maps across all error fixations. InferNet successfully identifies the subject’s goal and outperforms competitive null models, even without any object-specific training on the inference task.

## 1. Introduction

Inferring goals from actions is a central challenge for human-machine interactions and modeling human decisions. We consider the problem of inferring goals during eye movement behavior. Eye movements reflect rich information about the complex cognitive states of the brain, including thought processes and goals [3, 5, 1, 7].

In a visual search task, “error fixations” prior to locating the target (Fig. 1) are more likely to be on objects similar to the target [4, 9]. Therefore, it is possible to decode target information from the eye movements [2, 12, 8]. Existing target decoding methods are limited in using elementary search statistics [8], or handcrafted features [2, 12]. Moreover, existing approaches have only been tested with pre-defined classes and limited object set sizes. Those models do not generalize to infer targets from arbitrary classes. In contrast, here we propose a *zero-shot model*, the Inference Network (InferNet). InferNet applies knowledge

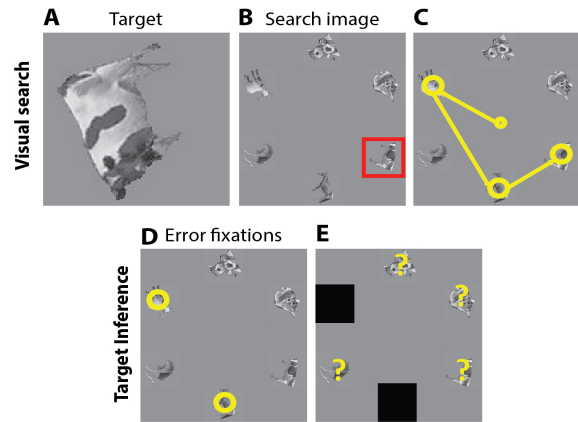


Figure 1. **Illustration of the target inference problem.** Human subjects were instructed to move their eyes to search for a given target (A) in the search image (B) irrespective of changes in size, rotation angles, or other format changes [13]. The red box shows the target location in the search image (not shown in the actual experiment). The visual search task resulted in a sequence of fixations (C, yellow circles). In this example, the subject required 3 fixations to find the target. We defined the 2 fixations on *non-target* objects as “error fixations”. In the target inference task, given the error fixations recorded from the psychophysics visual search task (D, yellow circles), the model infers what the subject was searching for (E, yellow question marks).

from object recognition to target inference. We tested InferNet on two visual search tasks [13] and show that it can predict human goals *without any re-training on new tasks*.

## 2. InferNet

We provide an overview of the proposed zero-shot deep network (InferNet, Fig. 2). Error fixations share more visual feature similarities with the target than with distractors [4, 9, 11]. Our model hypothesizes that locations showing higher feature similarity to error fixations are more likely to represent the subject’s goal. Given  $T$  error fixations with coordinates  $(x_i, y_i)$  where  $1 \leq i \leq T$ , the task is to predict a probabilistic map  $M_f$  of where the search target is most

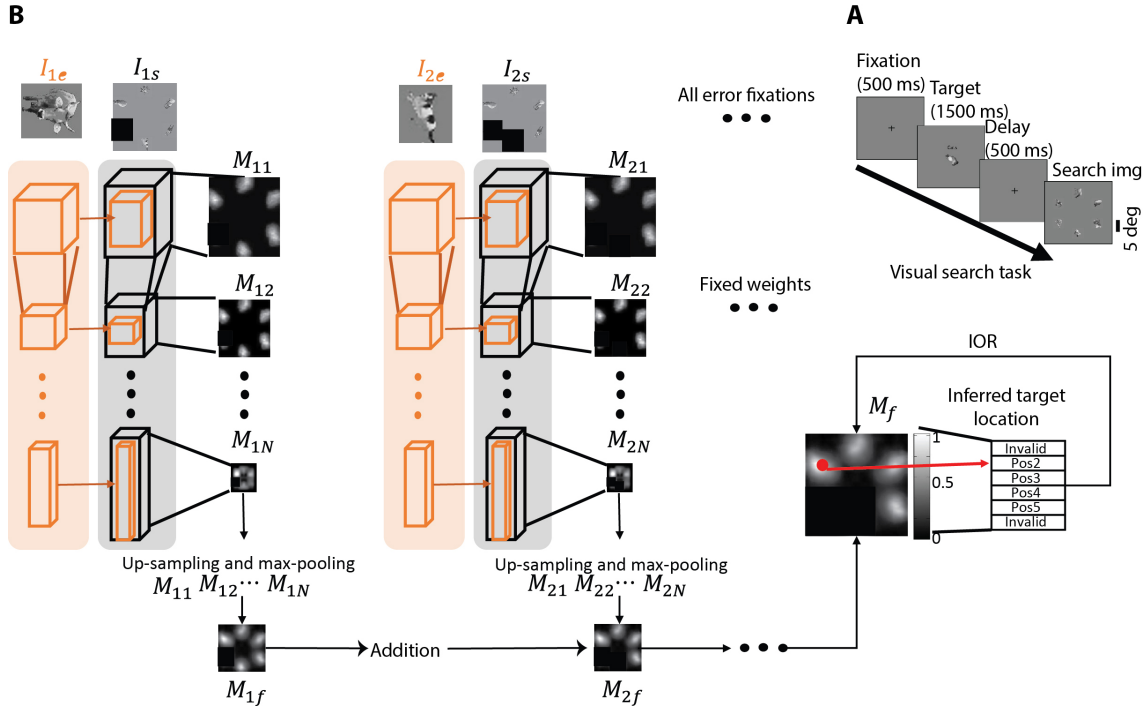


Figure 2. **Architecture of InferNet.** **A** Schematic of the visual search task [13]. **B** InferNet takes two inputs: the object  $I_{ie}$  at the error fixation  $i$  and the search image  $I_{is}$  with the object at the error fixation covered with a black mask. The error fixations are recorded from human subjects in the visual search task (A). See Sec. 2 for model details.

likely to be (Fig. 2). We take the maximum on  $M_f$  as the current guess location. If the current guess location overlaps with the target object location, the inference is deemed successful. Otherwise, after each incorrect guess, the map is updated by removing the erroneous inference location on  $M_f$  and the process continues iteratively.

The model uses a pre-trained deep convolutional network applied to the error fixations (**Prior Network (PN)**) and to the search image (**Likelihood Network (LN)**). **PN** takes the cropped area  $I_{ie}$  ( $28 \times 28$  pixels) centered at error fixation  $i$  as input and outputs feature maps across layers. We define  $I_{is}$  as the search image, with the objects at past error fixations  $1, \dots, i$  covered in black. **LN** modulates the feature maps from  $I_{is}$ , generating likelihood maps ( $M_{i1}, M_{i2}, \dots, M_{ij}, \dots, M_{iN}$ ) where  $j$  denotes the layer ( $1 \leq j \leq N$ ). These maps are concatenated and max-pooled to produce the final likelihood map  $M_{if}$  for error fixation  $i$  which tracks the parts of the image that are most similar between  $I_{ie}$  and  $I_{is}$ . InferNet adds the maps  $M_{if}$  across all  $T$  error fixations to build the final inference map  $M_f$ .

### 3. Experiments

We tested InferNet on data from two invariant visual search tasks using object arrays and natural images[13]. We filtered out those fixations on targets, and only used error

fixations (Fig. 1). The appearance of the target object in the search image was different from that in the target image.

We evaluate the model by the number of guesses required to infer the target as a function of the number of error fixations  $T$ . Because the target inference difficulty varies over trials, we report the relative performance  $P_r$ :  $P_r(T) = \frac{A_c(T) - A_m(T)}{A_c(T)}$ , where  $A_m(T)$  is the average number of guesses required by InferNet and  $A_c(T)$  is the average number of guesses required by chance. The larger  $P_r(T)$ , the more efficient the inference process is.

We compared InferNet with alternative null models, including chance (random guess from the remaining possible locations), template matching (slide error fixation patterns over the search image), bottom-up saliency (Itti-Koch, [6]), and RanWeight (InferNet with random weights without training). The null models proposed an inference map, and target selection procedure was the same as with InferNet. We made all the data (images, eye movements, and source code) publicly available online<sup>1</sup>.

#### 3.1. Object arrays

Figure 3 shows examples illustrating how the model efficiently inferred the target location given only 1-3 error fixations on object arrays and natural images. In Fig. 3A, a

<sup>1</sup><https://github.com/kreimanlab/HumanIntentionInferenceZeroShot>

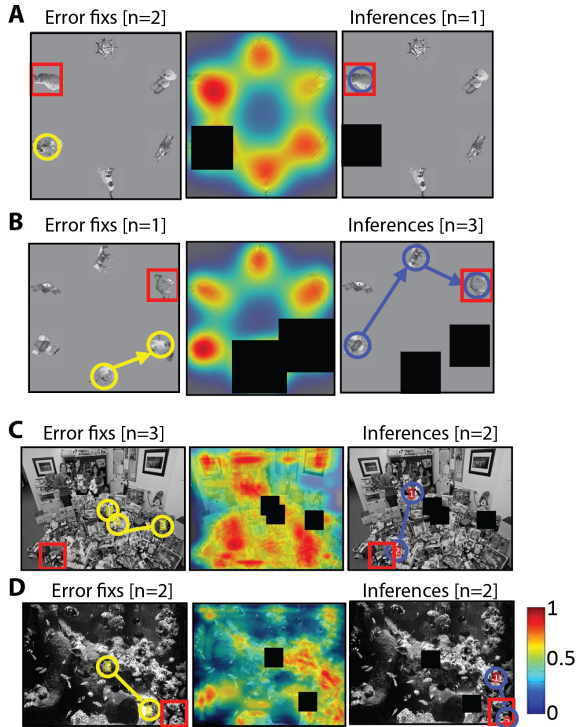


Figure 3. **Example target inference trials.** **A-B** Two object array trials. **C-D** Two natural image trials. Given the “error fixations” (yellow circles, column 1), InferNet predicts the probabilistic map  $M_f$  overlaid on the stimuli (Column 2, scale on the right). The red box (Column 1, 3) denotes the target location. The black boxes show the error fixation locations in  $M_f$ . The blue circles show the inferred locations (Column 3).

subject made one error fixation on the cow (yellow circle) which looks similar to the target (red square) before finding the target sheep. Given this single error fixation, InferNet determined that the subject was probably looking for a sheep (blue circle) instead of the 4 remaining distractors.

InferNet showed an overall improvement of  $3.8 \pm 3\%$  with respect to the chance model over all error fixations (Fig. 4A, blue). Even with a single error fixation, InferNet could infer the target 6.9% faster than chance. While random guessing would correctly land on the target within 3 guesses, InferNet only required  $2.80 \pm 0.01$  guesses.

None of the null models reached the performance level of InferNet (Fig. 4A,  $P < 10^{-19}$ , two-tailed t-test), except with 4 error fixations where none of the models were above chance. Although we took precautions to normalize *average* low-level features, InferNet can capitalize on shared low-level features between error fixations and the target. Performance for the bottom-up saliency model (IttiKoch) was better than the chance model but still below InferNet, suggesting that there is additional target information embedded in error fixations. The random weights (RanWeight) and pixel template matching (TempMatch) models showed minimal improvements

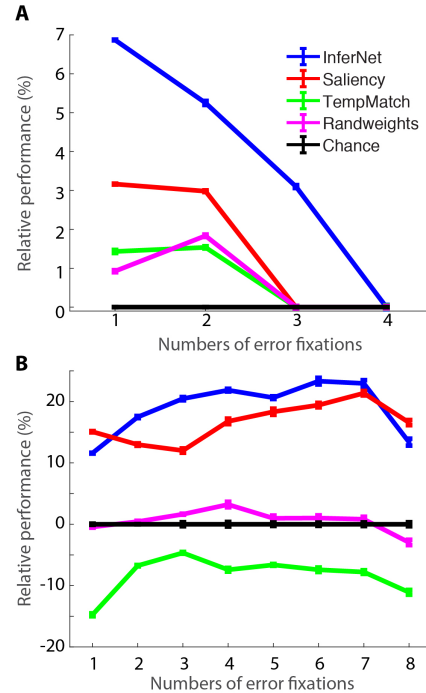


Figure 4. **Evaluation of model inference performance for object arrays (A) and natural images (B).** See Sec. 3 for descriptions of evaluation metrics and comparative models. Error bars are standard error of the mean.

compared to random locations (Fig. 4A), suggesting that the discriminative features learnt for image classification are important for target inference.

### 3.2. Natural scenes

The experiment reported so far focused on images consisting of segmented objects at discrete locations, presented on a uniform background, at fixed positions equidistant from the center of the image. In the real world, visual search takes place in cluttered environments involving non-segmented objects amidst a complex background. Next, we evaluated InferNet in the more challenging continuous domain of natural images (Fig. 3 and Fig. 4B).

Figure 3C shows an example where InferNet successfully determined the subject’s goal in natural images. The error fixations fell on plush toys, such as teddy bears. Based on the features of error fixations, the inference map showed high activations around plush toys regions. InferNet correctly inferred the target within 2 guesses.

InferNet showed significant improvements of  $19 \pm 4\%$  compared to chance (Fig. 4B). InferNet outperformed all the alternative null models (Fig. 4B,  $P < 10^{-26}$ , two-tailed t-test). The bottom-up saliency model (IttiKoch) performed high among all the null models because target objects were salient and relatively large. Template matching performed worse than the chance model.

Error Fixations	Top $N$ category inference accuracy %							
	1	2	4	8	16	32	64	128
1	6	8	11	13	17	29	38	55
2	4	9	14	20	23	33	46	65
3	3	10	16	25	28	38	51	72
4	0	13	20	20	30	39	54	74
5	3	11	23	28	34	42	56	74
6	0	14	23	31	37	45	54	77
7	1	15	25	31	37	47	53	78
8	4	17	28	37	40	48	57	80

Table 1. InferNet top- $N$  target category accuracy across error fixations (rows) where  $N = 1, 2, \dots, 128$  (columns)

#Error Fix.	Object Arrays		Natural Images			
	1	3	1	3	5	7
InferNet	6.87	3.10	12.83	22.48	24.35	28.28
Layer 5	1.88	1.58	9.35	17.11	17.24	20.91
Layer 10	3.98	0.67	14.69	24.82	25.16	26.97
Layer 17	5.96	1.99	16.50	19.28	22.42	26.43
Layer 23	7.46	0.01	13.32	24.72	28.07	23.56
Layer 24	6.60	3.28	18.53	28.04	30.59	30.42
Layer 30	8.21	3.08	-	4.45	6.03	3.36
Layer 31	7.56	2.34	-	6.15	5.00	3.93
Max + Max	6.87	1.13	12.84	21.11	22.96	24.49
Mean + Max	7.01	2.63	8.67	11.97	14.22	16.05
Mean + Mean	7.01	3.68	8.67	9.68	10.61	13.31
Humans	-	-	60.87	67.33	38.18	35.65
Common Fix.	10.96	1.52	20.09	30.35	24.98	25.77

Table 2. Target inference relative performance (%) of ablated models compared with chance  $T$  error fixations (the larger, the better). (-) means not significantly better than chance. Layer number refers to the index in the VGG16 network [10]. The first row (InferNet) corresponds to the full model, whereas the other rows show the predictions using either only one feature similarity map from  $M_{i1}$  to  $M_{i7}$  in Fig. 2 or their combinations.

### 3.3. Target category inference

In addition to inferring the target location, we asked whether InferNet could infer the categorical content of what the subjects were looking for. We selected 100 images (out of 240) where the target categories were in ImageNet. InferNet predicted the probability that the target belonged to each one of 1,000 categories by leveraging on the weights pre-trained on ImageNet and accumulating those probabilities across error fixations. Given even only one error fixation, InferNet surpassed the chance model (1/1000) (Table 1). As expected, the target category inference accuracy increased with  $T$ . Given 8 error fixations, InferNet correctly inferred the target category with top-4 accuracy of 28% out of 1,000 categories. These results show that InferNet can not only infer location goals but also categorical content goal information.

### 3.4. Ablation effect on target inference

To evaluate the contribution of different layers, we tested each feature similarity map  $M_j$  (Table 2). Target inference was better with higher layer feature similarity maps.

Averaging instead of max-pooling the likelihood maps did not yield any significant improvements in object arrays. However, alternative ways of combining feature similarity

maps led to large differences in natural images. InferNet outperforms the alternative models, suggesting that error fixations are guided by sub-patterns of the search target [8].

InferNet disregards error fixation order. Incorporating fixation order did not impact performance.

We asked how well humans can infer another subject’s intentions by conducting additional psychophysics experiments (Table 2). Subjects were not able to solve the inference problem in object arrays but were better than InferNet in natural images, perhaps by using contextual cues [14]. To investigate the between-subject variability, we created a new model using only fixations that are common across subjects. The results (last row) show that in some (but not all) cases, InferNet can overcome the consequences of variability in human scanpath patterns.

## 4. Conclusion

We proposed a computational model (InferNet) to infer intentions from observed eye movement behaviors. We evaluated the model in two visual search tasks. InferNet can predict the subjects’ goals (what they are looking for), in object array images and in natural images, by using the set of non-target error fixations, and without any domain-specific training. This proof-of-principle demonstration provides a possible solution to effectively estimate what the subject is searching for in complex images. The results suggest that computational models can make reasonable conjectures to predict human intentions from behavioral data.

## References

- [1] Betz et al. Investigating task-dependent top-down effects on overt visual attention. *JoV*, 10(3):15–15, 2010. 1
- [2] Borji et al. What do eyes reveal about the mind?: Algorithmic inference of search targets from fixations. *Neurocomputing*, 149:788–799, 2015. 1
- [3] Castelano et al. Viewing task influences eye movement control during active scene perception. *Journal of vision*, 9(3):6–6, 2009. 1
- [4] Eckstein et al. Similar neural representations of the target for saccades and perception during search. *Journal of Neuroscience*, 27(6):1266–1270, 2007. 1
- [5] Henderson et al. Predicting cognitive state from eye movements. *PLoS one*, 8(5):e64937, 2013. 1
- [6] Itti et al. A model of saliency-based visual attention for rapid scene analysis. *TPAMI*, 20(11):1254–1259, 1998. 2
- [7] Iqbal et al. Using eye gaze patterns to identify user tasks. In *The Grace Hopper Celebration of Women in Computing*, pages 5–10, 2004. 1
- [8] Rajashekar et al. Visual search: Revealing the influence of structural cues by gaze-contingent classification image analysis. *JoV*, 6(4):7–7, 2006. 1, 4
- [9] Robert G et al. Visual similarity effects in categorical search. *Journal of Vision*, 11(8):9–9, 2011. 1
- [10] Simonyan et al. Very deep cnn for large-scale image recognition. *CoRR*, abs/1409.1556, 2014. 4
- [11] Wolfe et al. Guided search 4.0. *Integrated models of cognitive systems*, pages 99–119, 2007. 1
- [12] Zelinsky et al. Eye can read your mind: Decoding gaze fixations to reveal categorical search targets. *JoV*, 13(14):10–10, 2013. 1
- [13] Zhang et al. Finding any waldo with zero-shot invariant and efficient visual search. *NC*, 9(1):1–15, 2018. 1, 2
- [14] Zhang et al. Putting visual object recognition in context. *arXiv*, 1911.07349, 2020. 4