

COCO-Search18: A Dataset for Predicting Goal-directed Attention Control

Yupei Chen¹, Zhibo Yang², Seoyoung Ahn¹, Dimitris Samaras², Minh Hoai², and Gregory Zelinsky^{1,2,*}

¹Department of Psychology, Stony Brook University

²Department of Computer Science, Stony Brook University

*gregory.zelinsky@stonybrook.edu

ABSTRACT

Attention control is a basic behavioral process that has been studied for decades. The currently best models of attention control are deep networks trained on free-viewing behavior to predict bottom-up attention control—saliency. We introduce COCO-Search18, the first dataset of laboratory-quality *goal-directed behavior* large enough to train deep-network models. We collected eye-movement behavior from 10 people searching for each of 18 target-object categories in 6202 natural-scene images, yielding ~300,000 search fixations. We thoroughly characterize COCO-Search18, and benchmark it using three machine-learning methods: a ResNet50 object detector, a ResNet50 trained on fixation-density maps, and an inverse-reinforcement-learning model trained on behavioral search scanpaths. Models were also trained/tested on images transformed to approximate a foveated retina, a fundamental biological constraint. These models, each having a different reliance on behavioral training, collectively comprise the new state-of-the-art in predicting goal-directed search fixations. Our expectation is that future work using COCO-Search18 will far surpass these initial efforts, finding applications in domains ranging from human-computer interactive systems that can anticipate a person’s intent and render assistance to the potentially early identification of attention-related clinical disorders (ADHD, PTSD, phobia) based on deviation from neurotypical fixation behavior.

Keywords: Goal Dataset, Attention Dataset, Fixation Dataset, Gaze Dataset, Visual Search, Inverse-Reinforcement Learning

1 The control of visual attention comes broadly in two forms.
2 One is bottom-up, where control is exerted purely by the
3 visual input^{1,2}. This is the form of attention predicted by
4 saliency models, which exploded in popularity in the behav-
5 ior fixation-prediction and computer-vision literatures^{1,3-5}.
6 The other form of control is top-down, where behavioral goals
7 rather than bottom-up salience control the allocation of visual
8 attention. Goal-directed attention control underlies all the
9 things that we *try* to do, and this diversity makes its prediction
10 vastly more challenging than predicting bottom-up saliency,
11 and more important. In addition to its basic research value, a
12 better understanding of goal-directed attention could lead to
13 the development of biomarkers for neurotypical attention behav-
14 ior against which clinical conditions can be quantitatively
15 compared, and to advances in intelligent human-computer
16 interactive systems that can anticipate a user’s visual goals
17 and render real-time assistance⁶⁻⁸.

18 Goal-directed attention has been studied for decades⁹⁻¹⁶,
19 largely in the context of visual search. Search is arguably the
20 most basic of goal-directed behaviors; there is a target object
21 and the goal is to find it, or conclude its absence. Goals are ex-
22 tremely effective in controlling the allocation of gaze. Imagine
23 two encounters with a kitchen, first with the goal of learning
24 the time from a wall clock and again with the goal of warming
25 a cup of coffee. These “clock” and “microwave” searches
26 would yield two very different patterns of eye movement, as
27 recently demonstrated in a test of this gedanken experiment¹⁷,
28 and understanding this goal-directed control has been a core
29 aim of search theory. The visual search literature is itself volu-

minous (see reviews¹⁸⁻²⁰). Here we focus on the prediction of
image locations that people fixate as they search for objects,
and how the selection of these fixation locations depends on
the target goal.

The visual search literature is not only mature in its empiri-
cal work, it is also rich with many hugely influential theories
and models^{12-16,21}. Yet despite this success, over the last
years progress has stalled. Our premise is that this is due to
the absence of a dataset of search behavior sufficiently large to
train deep network models. Our belief is based on observation
of what occurred in the bottom-up attention-control literature
during the same time. The prediction of fixations during free
viewing, the task-less cousin of visual search, has become
an extremely active research topic, complete with managed
competitions and leaderboards for the most predictive mod-
els²² (<http://saliency.mit.edu/>). The best of these saliency
models are all deep networks, and to our point, all of them
were trained on large datasets of labeled human behavior²³⁻²⁷.
For example, one of the best of these models, DeepGaze II²³,
is a deep network pre-trained on SALICON²⁵. SALICON is
a crowd-sourced dataset consisting of images that were an-
notated with mouse-based data approximating the attention
shifts made during free viewing. This model of fixation pre-
diction during free viewing was therefore trained on a form
of free-viewing behavior. Without SALICON, DeepGaze II,
and models like it²⁴⁻²⁷, would not have been possible, and
our understanding of free-viewing behavior, widely believed
to reflect bottom-up attention control, would be greatly di-
minished. For the task of visual search, there is nothing

30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58

remotely comparable to SALICON²⁵. Here we describe in detail COCO-Search18, the largest dataset of goal-directed search fixations in the world. COCO-Search18 was recently introduced at CVPR2020²⁸, and our aim in this paper is to elaborate on the richness of this dataset so as to increase its usefulness to researchers interesting in modeling top-down attention control.

Methods

Behavioral Data Collection

COCO-Search18 is built from Microsoft COCO, Common Objects in Context²⁹. COCO consists of over 200,000 images of scenes that have been hand-segmented into 80 object categories. This ground-truth labeling of objects in images makes COCO valuable for training computer vision models of object detection^{29–33}. However, in order for COCO to be similarly valuable for training models of goal-directed attention, these images would also need to be labeled with the locations fixated by people searching for different target-object goals. COCO-Search18 fills this niche by providing these training labels of search behavior.

The dataset consists of a large-scale annotation of a subset of COCO, 18 of its 80 object categories, with goal-directed search fixations. Each of 10 participants searched for each of 18 target-object categories (blocked) in 6,202 COCO images, mostly of indoor scenes. This effort required an average of 12 hours per participant, distributed over 6 days. This substantial behavioral commitment makes it possible to train models of individual searchers²⁸, although our focus here is on group behavior. The eye position of each participant was sampled every millisecond using a high-quality eye-tracker under controlled laboratory conditions and procedure, resulting in ~70,000,000 gaze-position samples in total. These raw gaze samples were clustered into 299,037 search fixations (~30,000 per participant), which dropped to 268,760 fixations after excluding those from incorrect trials. Figure 1 shows representative images and fixation behavior for each target category. See SM1 for details about: selection criteria (for images, target categories, and fixations), the eye tracker and eye tracking procedure, participant instruction, and a comparison between COCO-Search18 and existing datasets of search behavior.

Search-Relevant Image Statistics

Figure 2A shows three search-relevant characterizations of the COCO-Search18 images. The left panel shows the distribution of target-object sizes, based on bounding-box COCO labels. This distribution skewed toward smaller targets, with the range constrained by image selection to be between 1% and 10% of the image size (see SM1). The mean visual angle of the targets, based on the square root of bounding-box size, was 8.4°, about the size of a clenched fist at arm's length. The middle panel shows the distribution of initial target eccentricities, which is how far the target appeared in peripheral vision, based on center fixation at the start of search. Target eccentricities

ranged from 10° to 25° of visual angle, with a mean of ~15° eccentricity. The right panel shows the distribution of the number of “things” in each image, again based on the COCO object and stuff labels³⁴. Some images depicted only a handful of objects, whereas others depicted 20 or more (keeping in mind that this labeling was coarse). We report this statistic because search efficiency is known to degrade with the number of items in a search array³⁵, and a similar relationship has been suggested for the feature and object clutter of scenes^{36–38}. Figure 2B again shows these measures, now grouped by the 18 target categories. Target size and initial target eccentricity varied little across target categories, while the measure of set size varied more. See SM2 for analyses showing how each of these three measures correlated with search efficiency, for each target category.

Search Procedure and Metrics

The paradigm used for data collection was speeded categorical search^{39–41}. The participant's task was to indicate whether an exemplar of a target category appeared in an image of a scene (Figure S3). They did this by making a target present/absent judgment as quickly as possible while maintaining accuracy. The target category was designated at the start of a block of trials. Half of the search images depicted an exemplar of a target (target-present, TP), and the other half did not (target-absent, TA).

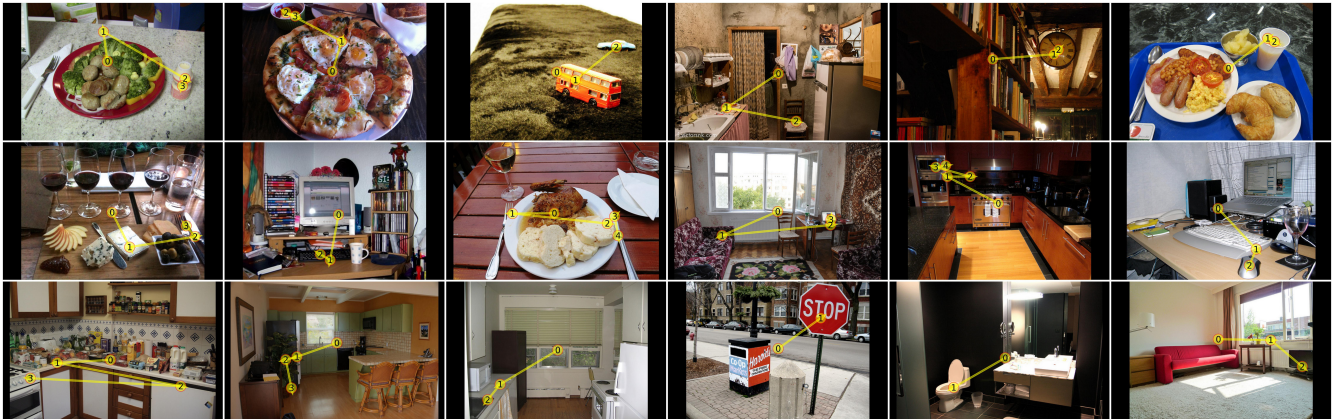
We measure goal-directed attention control as the efficiency in which gaze moves to the search target. Because the target was an object category, the term used for this measure of search efficiency is *categorical target guidance*^{39,41}, defined as the controlled direction of gaze to a target-category goal. We consider multiple measures of target guidance in Figure 3, but here we focus on the cumulative probability of fixating the target after each search saccade^{42–45}. A target category that can successfully control gaze will be fixated in fewer eye movements compared to one that has less capacity for target guidance. A desirable property of the target-fixation-probability (TFP) function (Figure 4) is that it is meaningful to compute the area under the TFP curve (TFP-*auc*), which we suggest as a new metric for evaluating search guidance across target categories and models.

Model Comparison

Now that COCO-Search18 exists, what can we do with it? To answer this question we conducted benchmarking to determine how well current state-of-the-art methods, using COCO-Search18, can predict categorical search fixations. To create a context for this model comparison we considered three very different modeling approaches, which all shared a common backbone model architecture, a ResNet50 pre-trained on ImageNet⁴⁶.

Our first approach predicted search fixations using object detectors trained for each of the target categories. We did this by re-training the pre-trained ResNet50 on just the 18 target categories using the COCO labels. Standard data augmentation methods of re-sizing and random crops were used

A



B

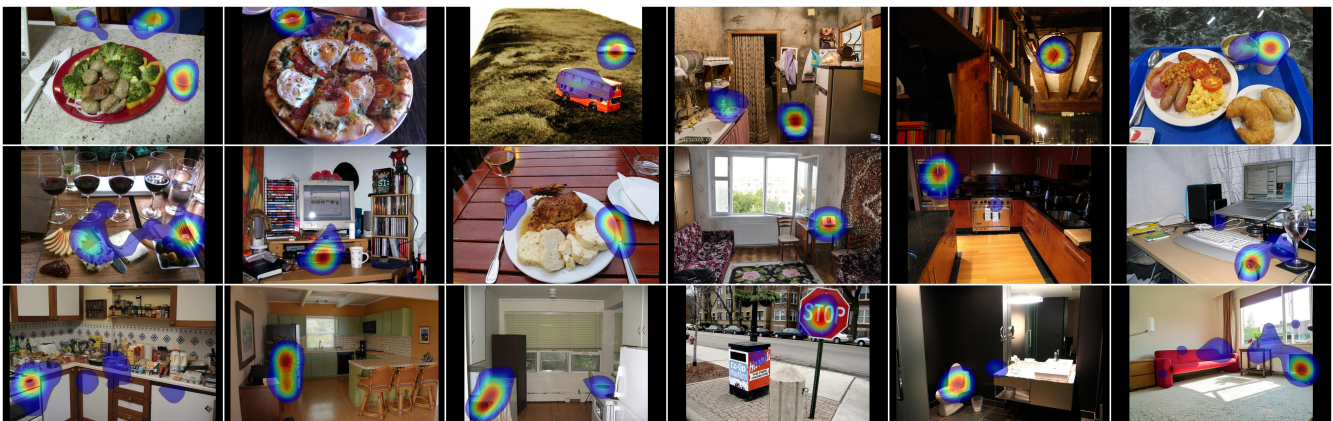


Figure 1. (A). Examples of target-present images for each of the 18 target categories. Yellow lines and numbered discs indicate a representative search scanpath from a single participant. From left to right, top to bottom: bottle, bowl, car, chair, (analog) clock, cup, fork, keyboard, knife, laptop, microwave, mouse, oven, potted plant, sink, stop sign, toilet, tv. (B). Examples of fixation density maps (excluding initial fixations at the center) computed over participants for the same scenes.

166 to increase variability in the training samples. We then used
167 these trained detectors to predict search fixations on the test
168 images. For a given target and test image, we obtained a con-
169 fidence map from the target detector and used it to sample a
170 sequence of fixation locations based on the level of confidence.
171 Note that this approach is pure computer vision, meaning that
172 it uses the image pixels solely and knows nothing about behav-
173 ior.

174 With COCO-Search18, however, it is possible to also train
175 on the search behavior. There are multiple ways of doing this.
176 In our second approach we re-trained the same ResNet50,
177 only this time using labels as the fixations made by searchers
178 viewing the training images. Specifically, fixation-density
179 maps (FDMs) were obtained for each TP training image for
180 a given category, and these were used as labels for model
181 training. This model is in a sense a search version of mod-
182 els like DeepGaze II²³ in the free-viewing fixation-prediction
183 literature, which are also trained to predict FDMs. We there-
184 fore refer to this model as Deep Search. Deep Search differs
185 from the Target Detector model in that it is trained on search

fixation density to predict search behavior.

186
187 For our third modeling approach we used inverse-
188 reinforcement learning (IRL)⁴⁷⁻⁴⁹, an imitation-learning
189 method from the machine-learning literature, to simply mimic
190 the search scanpaths observed during training. We chose IRL
191 over other imitation-learning methods because it is based on
192 reward, known to be a powerful driver of behavior⁵⁰⁻⁵², but
193 we think it is likely that other imitation-learning methods
194 would perform similarly. The IRL model we used⁴⁹ works
195 by learning, through an adversarial process playing out over
196 many iterations, how to make model-generated behavior, ini-
197 tially random, become more like human behavior. It does this
198 by rewarding behavior that happens to be more human-like.
199 IRL is therefore very different from a Target Detector, but
200 also different from Deep Search, which also gets to use search
201 behavior in its training. The IRL model learns to imitate the
202 search scanpath, meaning the sequence of fixations made to
203 the search target, whereas Deep Search uses only FDMs that
204 do not represent the temporal order of fixations. Because the
205 IRL model used the most search behavior for training, we

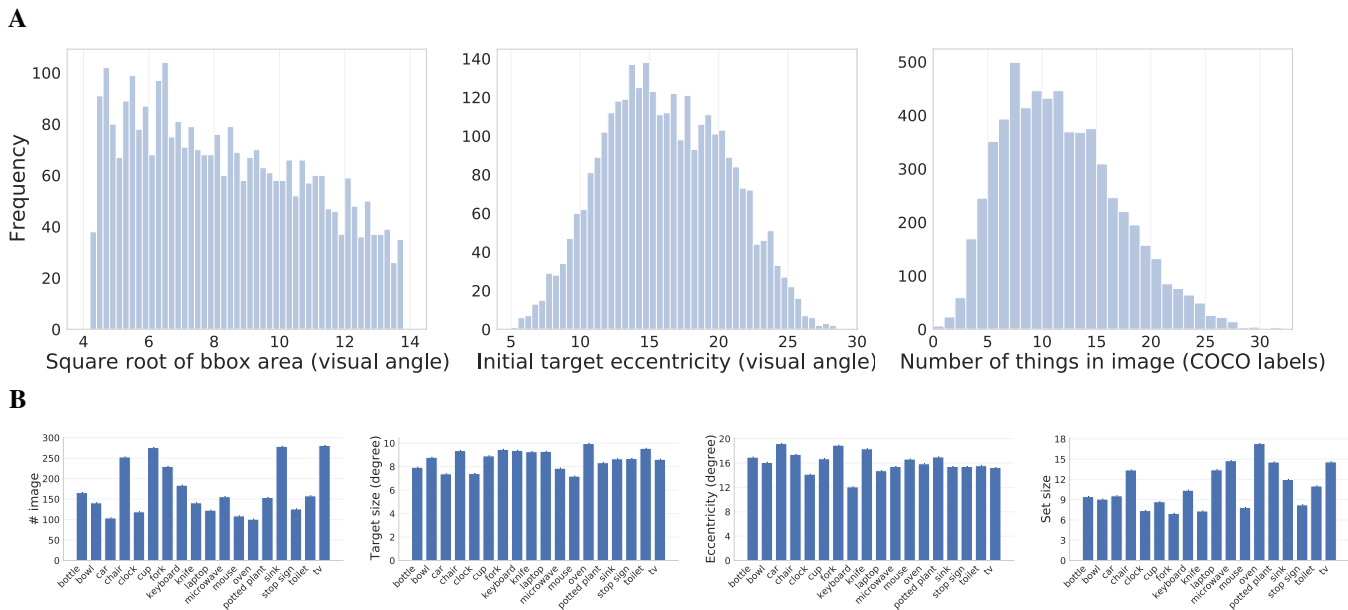


Figure 2. (A). Distributions of target sizes, based on the visual angle of their bounding-box areas (left), and initial target eccentricities (middle), both for the target-present images. The number of “things” (objects and “stuff” categories, both based on COCO-stuff labels) appearing in the search images (right). (B). Image statistics from COCO-Search18, grouped by the 18 target categories. The left plot shows the number of images, followed by three analyses paralleling those presented in (A): averaged target-object size in degrees of visual angle, initial target eccentricity based on bounding-box centers, and the average number of things in an image (a proxy for set size).

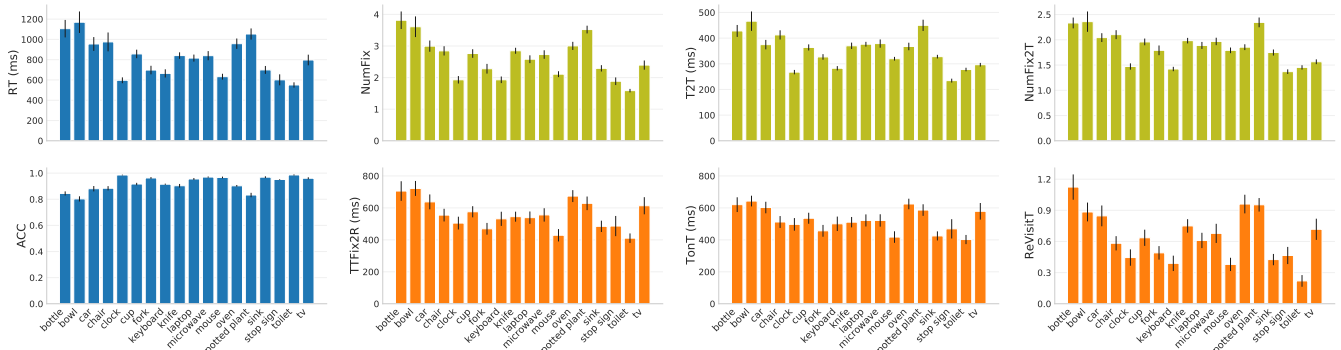


Figure 3. Basic behavioral analyses of the target-present data from COCO-Search18, grouped by the 18 target categories. Blue plots (left two) show the manual measures of reaction time (RT) and response accuracy (ACC). Olive plots (top row) show gaze-based analyses of categorical guidance efficiency: number of fixations made before the button press (NumFix), time until first target fixation (T2T), and number of fixations made until first target fixation (NumFix2T). Orange plots (bottom row) show gaze-based measures of target verification: time from first target fixation until response (TTFix2R), total time spent fixating the target (TonT), and the number of re-fixations on the target (ReVisitT). Values are means over 10 participants, and error bars represent standard errors.

206 hypothesized that it would best predict search behavior in our
207 model comparison. See SM3 for additional details about IRL.

208 State Comparison

209 In addition to the model comparison, we also compared several
210 state representations used by the models. In the current
211 context, the state is the information that is available to control
212 search behavior, and essential to this are the features extracted
213 from each search image. We refer to the original images as
214 high-resolution (Hi-Res), in reference to the fact that they
215 were not blurred to reflect retina constraints. Extracting fea-

216 tures from a Hi-Res image produces a Hi-Res state, and it
217 is this state that is used by most object-detection models in
218 computer vision where the goal is to maximize detection suc-
219 cess. Primate vision, however, is profoundly degraded from
220 this Hi-Res state by virtue of the fact that we have a foveated
221 retina. A foveated retina means that high-resolution visual
222 inputs exist only for a small patch of the image at the current
223 fixation location, and blurred everywhere else. Given our
224 goal to model the fixation behavior of the COCO-Search18
225 searchers, each of whom had a foveated retina, we included

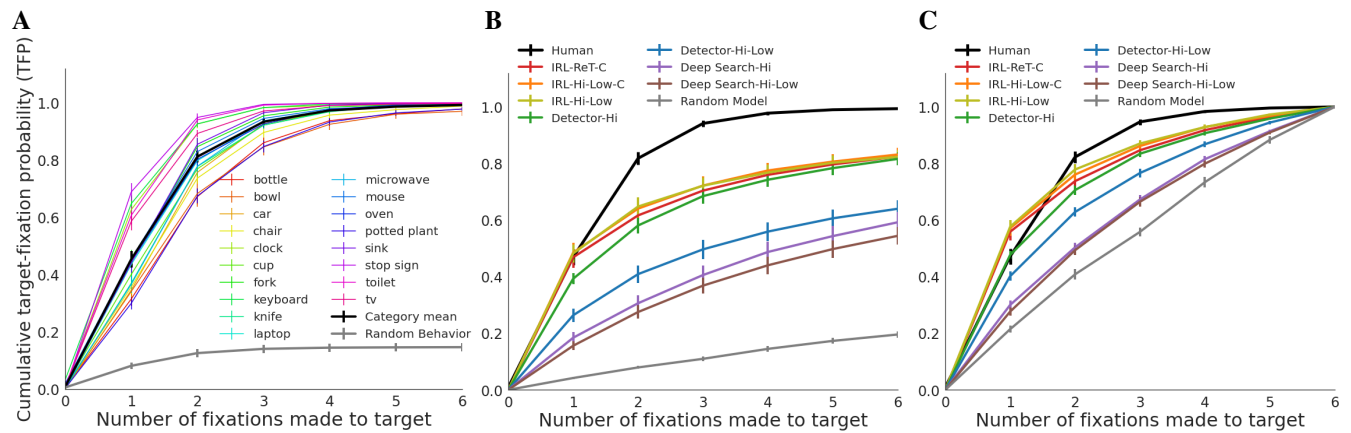


Figure 4. (A). Cumulative probability of fixating the target (y-axis; target-fixation probability or TFP) as a function of fixation serial position (x-axis; 0-6), shown individually for the 18 target categories (color lines) and averaged over target types (bold black line). The bottom-most function is a Random Behavior baseline obtained by computing target-fixation probability using a scanpath from the same participant searching for the same target category but in a different image. For the 18 target functions, means were computed by first averaging over images and then over participants, and standard errors were computed over participants. For the averaged behavioral data and the Random Behavior baseline (black and gray lines), means were computed by first averaging over images and then over categories, and standard errors were computed over categories. (B). TFP functions generated from model predictions on the test images. Names designate a model type (IRL, Detector, Deep Search) and a state representation (ReT, Hi-Low, Hi, C), separated by hyphens. Average behavioral TFP is again plotted in bold black, this time for just the test data (which explains the small differences from the corresponding function in A, which included the training and testing data). The Random Model baseline was obtained by making six movements of the Hi-Low foveated retina, with ISTs after each, and determining whether any of these movements brought the high-resolution central window to the target. Means were first computed over images and then over categories, and standard errors were computed over categories. (C). A re-plot of B, but only including data from trials in which the target was successfully fixated within the first six fixations (i.e., search scanpaths that succeed in locating the target.)

226 this basic biological constraint in the state to determine its
 227 importance in model training and prediction of search behav-
 228 ior (see also⁵³). Relatedly, and as fundamentally, each
 229 new fixation changes the state by allowing high-resolution
 230 information to be obtained from the vantage of a new image
 231 location. Capturing these fixation-dependent spatio-temporal
 232 state changes in the context of search was a core goal in the
 233 development of COCO-Search18.

234 We considered two fovea-inspired states. In the first we
 235 used the method from Perry and Geisler⁵⁴ to compute a Retina-
 236 Transformed (ReT) image. A ReT image is a version of the
 237 Hi-Res image that is blurred to approximate the gradual loss
 238 in visual acuity that occurs when viewing at increasing ec-
 239 centricities in peripheral vision. Second, we implemented an
 240 even more simplified foveated retina consisting of just a high-
 241 resolution central patch ($7^\circ \times 7^\circ$ visual angle) surrounded by
 242 low-resolution “peripheral” vision elsewhere, with the critical
 243 difference from the ReT image being that only a single level of
 244 blur (Gaussian filter with $\sigma = 2$) was used to approximate the
 245 low-resolution periphery. Computing the gradual blur used in
 246 the ReT image was computationally very demanding, and the
 247 inclusion of the simpler Hi-Low state was motivated largely to
 248 reduce these computational demands (ReT requires $\sim 15\times$ the
 249 processing time per image). However, having this condition
 250 also enabled a needed initial evaluation of how veridically
 251 low-level visual-system constraints need to be followed when
 252 training deep-network models of human goal-directed behav-

ior.

253 We also considered two spatio-temporal state representa-
 254 tions for how information is accumulated with each new fixa-
 255 tion in a search scanpath. A behavioral consequence of having
 256 a foveated retina is that we make saccadic eye movements,
 257 and the order in which these eye movements are made cor-
 258 respond to different visual states. Our first spatio-temporal
 259 state assumed a high-resolution foveal window that simply
 260 moves within a blurred image. This means that each change in
 261 fixation brings peripherally blurred visual inputs into clearer
 262 view, and causes previously clear visual inputs to become
 263 blurred. This spatio-temporal state representation is aligned
 264 most closely with the neuroanatomy of the oculomotor system,
 265 so we will consider this to be the default state. However, this
 266 default state representation assumes that foveally-obtained
 267 information on fixation n is completely lost by fixation $n+1$,
 268 and indeed something like this is true for high-resolution in-
 269 formation about visual detail⁵⁵. However, this state fails to
 270 capture any memory for the fixated objects that persists over
 271 eye movement, which is also known to exist⁵⁶. To address
 272 the potential for an object context to build over fixations, we
 273 therefore also used a state that accumulates the high-resolution
 274 foveal views obtained at each fixation in the search scanpath,
 275 a state we refer to as Cumulative (-C). Over the course of
 276 multi-fixation search, the Hi-Low-C state would therefore
 277 accumulate high-resolution foveal snapshots with each new
 278 fixation, progressively de-blurring what would be an initially
 279

280 moderately-blurred version of the image. We explore these
281 two extremes of information preservation during search so as
282 to inform future uses of a fovea-inspired spatio-temporal state
283 representation to train deep network models.

284 General Model Methods

285 All of the models followed the same general pipeline. Each
286 1050×1680 image input was resized to 320×512 pixels
287 to reduce computation. This is what we refer to as the Hi-Res image
288 (or just Hi); the ReT and Hi-Low images were computed
289 from this. These images were passed through the ResNet50
290 backbone to obtain 20×32 feature map outputs, with the fea-
291 tures extracted from these images now reflecting either Hi-Res,
292 ReT, or Hi-Low states, respectively. Different models were
293 trained using these features and others, as described in the
294 Model Comparison section, and all model evaluations were
295 based on a 70% training, 10% validation, and 20% testing,
296 random split of COCO-Search18 within each target category.
297 See SM3 for additional details about the training and testing
298 separation, and Figure S15 for how the two compare on search
299 performance measures.

300 The trained models were used to obtain model-specific pri-
301 ority maps for the purpose of predicting the search fixations in
302 each test image. The priority map for the Target Detector was
303 a map of detector confidence values at each pixel location, and
304 fixations were sampled probabilistically from this confidence
305 map. The priority map for Deep Search is its prediction of the
306 FDM, given the input image and the model's learned mapping
307 between image features and the FDM ground-truth during
308 training. The priority map for the IRL model is the reward
309 map recovered during its training, which recall occurred dur-
310 ing its learning to mimic search behavior. Because this search
311 behavior was itself reward driven, the priority map for the
312 IRL model is therefore a map of the total reward expected by
313 making a sequence of search fixations to different locations
314 in a test image. The IRL model was additionally constrained
315 to have an action space discretized into a 20×32 grid, which
316 again was done to reduce computation time. A given action,
317 here a change in fixation, is therefore a selection of 1 from 640
318 possible grid cells, a sort of limitation imposed on the spatial
319 resolution of the model's oculomotor system. The selected
320 cell was then mapped back into 320×512 image space by
321 upscaling, and the center of this cell became the location of
322 the model's next fixation. The non-IRL models made their
323 action selection directly in the 320×512 image space, with
324 higher priority values selected with higher probability.

325 All of the model \times state combinations in our comparison
326 were required to make six changes in fixation for each test
327 image. This number was informed by the behavioral data
328 showing that the probability of target fixation was clearly at
329 ceiling by the sixth eye movement (Figure 4A). To produce
330 these 6-fixation scanpaths, we iterated the fixation generation
331 procedure using inhibitory spatial tagging (IST), which is a
332 mechanism serving the dual functions of (1) breaking current
333 fixation, thereby enabling gaze to move elsewhere, and (2)

discouraging the refixation of previously searched locations. 334
IST has long been used by computational models of free 335
viewing^{57,58} and search^{59,60}. Here we enforce IST by setting 336
the priority map to zero after each fixation over a region having 337
a radius of 2.5° visual angle (based on a 3×3 grid within the 338
 20×32 action space). IST was applied identically after each 339
fixation made by all of the models. This was true even for 340
models that did not have a foveated retina, such as a Target 341
Detector with a Hi-Res state, in which case IST was applied 342
to the image locations selected for "fixation". See SM3 for 343
additional details. 344

The nomenclature that we adopted for the model compari- 345
son consists of the model type as the base and the state repre- 346
sentation as a suffix. If the spatio-temporal state is cumulative, 347
there is a second suffix of -C. For example, the IRL-ReT-C 348
model accumulates graded-resolution foveal views of an im- 349
age with each reward-driven eye movement. Although our 350
aim is to explore as systematically as possible each state for 351
every model, for some models a given state representation is 352
not applicable. For example, it makes no sense for the IRL 353
model to use the Hi-Res state. Because that state representa- 354
tion does not change from one search fixation to the next it 355
would be impossible to learn fixation-dependent changes in 356
state, thereby defeating the purpose of using the IRL method. 357
Similarly, it makes no sense to have a cumulative state for 358
anything but the IRL model, as the others would be unable to 359
use this information. However, it does make sense to test a 360
Target Detector and Deep Search on a Hi-Low state as well as 361
a Hi-Res state, and these models are included in the Table 1 362
model evaluation. 363

364 Results

365 Behavioral Performance

366 We interrogated COCO-Search18 using multiple performance 367
measures. Figure 3 reports these analyses for each of the 368
target categories. Analyses can be conceptually grouped into 369
manual measures (accuracy and response time; blue plots), 370
gaze-based measures of categorical guidance (number of fixa- 371
tions before the button press, and both the time and number 372
of fixations until the first target fixation; olive plots), and mea- 373
sures of target verification time (time from first target fixation 374
until the button press, total time spent fixating the target, and 375
the number of target re-fixations; orange plots). What is clear 376
from these analyses is that, except for accuracy, there is wide 377
variability across target categories in these measures, and this 378
variability creates fertile ground for future model develop- 379
ment. Also clear from Figure 3 is that there is considerable 380
correlation among some of these measures, perhaps most evi- 381
dent among the search guidance measures where the shapes 382
of the plots look similar. We include these different measures, 383
not to suggest their independence, but rather as a courtesy to 384
readers who may be familiar with different measures.

385 Figure 5 is a matrix visualization of these analyses, now 386
with color coding a ranking of search efficiency. In Figure 5A, 387
the deepest red for each measure (row) indicates the least effi-

	TFP-AUC \uparrow	Probability Mismatch \downarrow	Scanpath Ratio \uparrow	Sequence Score \uparrow	MultiMatch \uparrow			
					shape	direction	length	position
Human	5.200	-	0.862	0.489	0.903	0.736	0.880	0.910
Random Model	0.744	4.455	0.392	0.266	0.832	0.579	0.783	0.755
Detector-Hi	4.001	1.209	0.680	0.411	0.877	0.665	0.837	0.872
Detector-Hi-Low	2.975	2.225	0.601	0.370	0.863	0.640	0.820	0.833
Deep Search-Hi	2.519	2.681	0.579	0.348	0.890	0.627	0.867	0.861
Deep Search-Hi-Low	2.282	2.918	0.546	0.333	0.882	0.617	0.859	0.848
IRL-ReT-C	4.170	1.131	0.731	0.418	0.879	0.673	0.842	0.874
IRL-Hi-Low-C	4.262	1.031	0.747	0.419	0.886	0.677	0.849	0.885
IRL-Hi-Low	4.245	1.036	0.753	0.417	0.884	0.677	0.847	0.885

Table 1. Results from fixation-prediction models (rows) using multiple scanpath metrics (columns) applied to the COCO-Search18 test images. Arrows indicate the direction of better prediction success, and values in bold indicate best predictions across the model comparison. In the case of Sequence Score and MultiMatch, “Human” refers to an oracle method whereby one searcher’s scanpath is used to predict another searcher’s scanpath; “Human” for all other metrics refers to observed behavior. See the main text for additional details about the scanpath-comparison metrics, and SM3 for purely spatial comparisons using the AUC, NSS, and CC metrics.

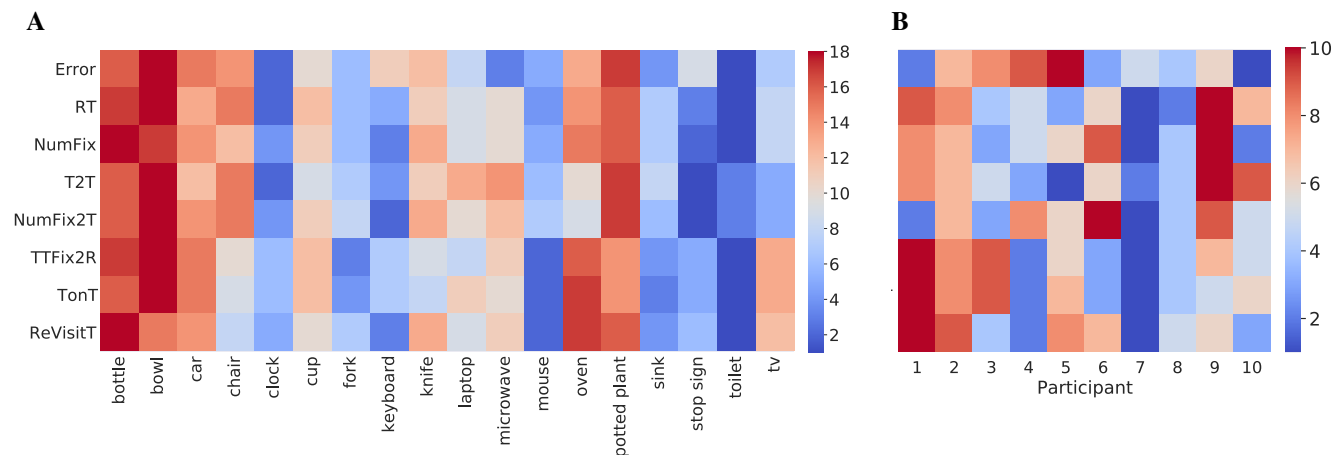


Figure 5. (A). Ranked target-category search efficiency [1-18], averaging over participants. Redder color indicates higher rank and harder search targets, bluer color indicates lower rank and easier search. Target category is grouped (columns) and shown for multiple performance measures (rows). These measures include: response Error, reaction time (RT), number of fixations (NumFix), time to target (T2T), number of fixations to target (NumFix2T), time from first target fixation until response (TTFix2R), time spent fixating the target (TonT), and the number of target re-fixations (ReVisitT). (B). A similar ranking of the target-present data, only now for participant efficiency (columns 1-10), averaged over target category. Performance measures and color coding are the same as in panel A.

388 cient (or most difficult) search over the 18 target categories, 402
 389 and the deepest blue indicates the most efficient (or easiest) 403
 390 search. The appearance of columns in this visualization cap- 404
 391 tures the agreement among the measures. More subtle patterns 405
 392 in the data can also be seen. For example, the two predomi- 406
 393 nately red columns at the left indicate agreement in that the 407
 394 bottle and bowl objects were difficult targets, speculatively 408
 395 because these target categories have particularly high variabil- 409
 396 ity in their image exemplars. Relatedly, appearing near the 410
 397 right are two of the consistently easiest targets, stop signs and 411
 398 toilets, both having relatively well-defined category member-
 399 ship. Figure 5B shows a similar plot, only now performance is
 400 averaged over target categories and plotted for individual par-
 401 ticipants. Search accuracy and efficiency clearly differ among

the participants in this ranking. Participants 7 and 8 were
 better searchers than Participants 2 and 9, meaning that they
 tended to find the target faster and with fewer fixations while
 keeping a low error rate. Differences in search strategy can
 also be seen from this visualization. Participant 1 searched
 the display carefully, resulting in few missed targets, but this
 person’s search was not very efficient. In contrast, Participant
 4 was quick to find and verify targets, but had relatively low
 accuracy. See SM2 for parallel analyses of the target-absent
 data from COCO-Search18.

However, arguably the gold-standard measure of attention
 control is the cumulative probability of fixating the target
 after each saccade made during search, the target-fixation
 probability (TFP). Figure 4A shows TFP functions for the first

416 six search saccades, averaged over participants and plotted for
417 individual target categories. The mean behavior over targets is
418 indicated by the bold black function. The noteworthy pattern
419 is that the slope of the group function is far steeper than that
420 of the chance baseline, obtained by computing TFP using a
421 scanpath from a different image but from the same participant
422 and target category. On average, about half of the targets were
423 fixated with the very first saccade. By the second saccade
424 TFP jumped to .82, and by the third saccade it reached a
425 performance ceiling of .94, which increased only slightly after
426 saccades 4–6. This high degree of attention control means
427 that, although we aimed to create a search dataset having a
428 moderate level of difficulty, COCO-Search18 skews easy.
429 This is due in part to unexpectedly large practice effects (see
430 SM2 for details). However, it is fortuitous that such a strong
431 attention-control signal exists in the behavioral data, given the
432 challenge faced by even start-of-the-art models in predicting
433 this simple search behavior.

434 Model Evaluation

435 Fixation prediction models broadly fall into two groups, mod-
436 els that predict the spatial distribution of locations fixated by
437 participants viewing an image (i.e., the FDM), and models
438 that predict both the location and order of the fixations made
439 by a person viewing an image (i.e., the scanpaths). In a search
440 task, fixation behavior changes dramatically over the first few
441 eye movements⁶¹, making it important to consider the spatio-
442 temporal fixation order. For this reason, we will focus on
443 spatio-temporal fixation prediction here and defer discussion
444 of purely spatial FDM prediction to SM4, and especially Ta-
445 ble S1. Both types of prediction were based on the 6-saccade
446 sequences that each model was required to make for each test
447 image. Specifically, 10 6-fixation scanpaths (excluding the
448 initial fixation) were predicted for each test image by sam-
449 pling probabilistically from the generated priority map, and
450 for each of these search scanpaths the model behavior was
451 analyzed up to first fixation on the target, or six changes in
452 fixation, whichever came first.

453 Predicting the spatio-temporal order of search fixations
454 can also take two forms. One has been to make fixation
455 predictions with respect to the search target. For example,
456 predicting the probability of the target being fixated by the
457 first search saccade, the second, etc. These target-based pre-
458 dictions capture the efficiency of search, where the goal is
459 to find the target, and models making this type of prediction
460 have been the more common in the search literature^{45,62}.
461 Here we use three metrics to evaluate the success of these
462 predictions. Two of these metrics were derived from the TFP
463 function (Figure 4A): TFP-*auc*, which is the area under the
464 cumulative target-fixation-probability curve, and Probability
465 Mismatch, which sums over each fixation in a scanpath the
466 absolute differences between the behavioral and model TFP.
467 The third metric, Scanpath Ratio, is the Euclidean distance
468 between the initial fixation location (roughly the center of
469 the image) and the location of the target (center of bounding

box) divided by the summed Euclidean distances between
470 the fixation locations in the search scanpath⁴². It is a search
471 efficiency metric because an initial saccade that lands directly
472 on the target would yield a Scanpath Ratio of 1, and all less
473 efficient searches would be < 1 . An alternative to predicting
474 target guidance over the spatio-temporal search scanpath is to
475 predict the scanpath itself. This approach assumes that any
476 target guidance would be reflected in the sequence of fixated
477 image locations leading up to the target decision. We con-
478 sidered two metrics for comparing behavioral and predicted
479 search scanpaths: Sequence Score, which clusters scanpaths
480 into strings and uses a string matching algorithm for com-
481 parison⁶³, and MultiMatch, which takes a multi-dimensional
482 approach to computing scanpath similarity^{64,65}. Both metrics
483 capture properties of the spatio-temporal search scanpath and
484 place less importance on the fact that there is a search target.
485 SM4 should be consulted for additional details about these
486 metrics.
487

488 Table 1 provides an evaluation of how each model \times state
489 combination fared in fair comparison using these metrics.
490 As we hypothesized, the three IRL models generally outper-
491 formed the others (see Table S2 for statistical tests). They did
492 so for every metric except MultiMatch, where all the models
493 performed similarly. The only other model that was compara-
494 bly predictive was Detector-Hi, but this model has no fovea
495 and is therefore the least biologically plausible. A perhaps
496 clearer picture of this model comparison can be obtained by
497 comparing the behavioral TFP function to ones computed
498 for each model. Figure 4B shows this evaluation of search
499 efficiency for each of the model \times state combinations (in color)
500 and for the mean search behavior (in black), limited to the
501 TP test data. Focusing first on state comparisons, we did not
502 find large differences between the states tested. Whether blur
503 was graded or binary appeared not to matter, as indicated by
504 the very similar TFP functions for the ReT and the Hi-Low
505 states using the IRL model. This pattern also appeared in
506 Table 1, where the IRL models differed by tiny margins. For
507 this reason, and its far greater computational efficiency, we
508 adopted only the Hi-Low state in the other model comparisons
509 (therefore, there are no Deep Search-ReT or Detector-ReT
510 models). Similarly, but specific to the IRL model, it made
511 little difference whether or not the state accumulated high-
512 resolution visual information with each fixation in a search
513 scanpath. The fact that the IRL model seemed not to use this
514 accumulated visual information is broadly consistent with the
515 view that very little high-resolution information is preserved
516 across saccades⁵⁵. However, it did matter whether the state
517 included a foveated retina or not, as exemplified by the dif-
518 ference between Hi-Res and Hi-Low states for the Detector
519 model. This state comparison suggests that future work may
520 want to avoid manipulations of fine-grained retinal blur and
521 assumptions about intersaccadic visual memory, and focus on
522 adding more basic limitations on human visual perception to
523 a model's pipeline, with the inclusion of a Hi-Low foveated
524 retina being one example.

525 All of the tested models made reasonable predictions of
526 search behavior in this challenging benchmark, where “rea-
527 sonable” is liberally defined as bearing greater resemblance
528 to the human behavior than the chance baseline. However, the
529 Deep Search models and the Detector-Hi-Low model were
530 clearly less efficient in their search behavior than either hu-
531 man behavior or any of the IRL models. This poor relative
532 performance is likely caused by these models not capturing
533 the serial order of search fixations, and that this order mat-
534 ters. A corollary finding is that the IRL models, because they
535 learned these spatio-temporal sequences of search fixations,
536 better predicted search behavior. This was true for all the IRL
537 models, which all predicted the efficiency of the first search
538 fixation almost perfectly (IRL models vs. Human at fixation 1
539 with post-hoc t-tests, all $p_{S_{bonferroni}} = 1.0$). Also interesting
540 is the degree that an object detector (Detector-Hi) can pre-
541 dict search behavior, supporting previous speculation⁶⁶. If an
542 application’s goal is to predict a person’s early fixation be-
543 havior during search without regard for biological plausibility,
544 a simple object detector will work well based on our testing
545 with COCO-Search18. Another finding from Figure 4B is that
546 none of the models achieved the high level of successful target
547 fixation exhibited in human performance. Performance ceil-
548 ings after six saccades (termed *fixated-in-6 accuracy*) ranged
549 from .54 (Deep Search-Hi-Low) to .83 (IRL-Hi-Low-C), all
550 well below the near perfect fixated-in-6 accuracy (.99) from
551 human searchers (post-hoc t-tests with all $p_{S_{bonferroni}} < .001$).
552 These lower performance plateaus, undoubtedly reflecting
553 limitations in current object detection methods, means that
554 the models tended either to fixate the target efficiently in the
555 first one or two eye movements (like people), or tended not
556 to fixate the target at all (unlike people). If a model cannot
557 represent the features used for target guidance as robustly as
558 people, there may be images for which there is essentially
559 no guidance signal, and on these inefficient search trials the
560 number of eye movements needed to fixate the target would
561 often be greater than six, hence the performance plateaus.

562 These different performance ceilings are problematic in that
563 they conflate limitations arising from object detection with
564 limitations in effective target prioritization, as measured by
565 search efficiency. For example, a strength of the TFP-auc met-
566 ric is that it is grounded in the TFP functions from Figure 4B,
567 but this means that it includes the different performance ceil-
568 ings in its measure and this weakens it as a pure measure of
569 attention control. To address this concern, in Figure 4C we
570 again plot TFP functions, but now only for trials in which the
571 target was successfully fixated within the first six saccades.
572 By restricting analysis to only trials having perfect fixated-
573 in-6 accuracy, the metric becomes more focused on search
574 efficiency. By this measure, and keeping in mind that the data
575 are now skewed toward easier searches, the IRL-Hi-Low-C
576 and IRL-Hi-Low models remain the most predictive overall,
577 although now all IRL models overestimate slightly the effi-
578 ciently of the first search saccade. But perhaps the biggest
579 winner in this comparison is the Detector-Hi model, which

now predicts TFP almost perfectly after the first fixation, and
has generally improved performance for subsequent fixations.
We tentatively conclude that simple prioritization of fixations
by an object detector predicts reasonably well the prioritiza-
tion of behavioral fixations in visual search. The losers in this
comparison were the Deep Search models, which remained
less efficient than human behavior even after normalization
for fixated-in-6 accuracy.

588 Discussion

589 Recent years taught us the importance of large datasets for
590 model prediction, and this importance extends to models of
591 attention control. COCO-Search18 is currently the largest
592 dataset of goal-directed search fixations, having sufficient
593 number to be used as labels for training deep network mod-
594 els. We conducted a systematic (but still incomplete) explo-
595 ration of models and state representations to provide some
596 initial context for the types of model predictions that are pos-
597 sible using COCO-Search18, given current state-of-the-art
598 (or nearly so). This model comparison focused on the de-
599 gree that search behavior was used during training, ranging
600 from none (Detector), to some (Deep Search), to entire search-
601 fixation scanpaths (IRL). With respect to the IRL model, its
602 use with COCO-Search18 is the first attempt to predict the
603 spatio-temporal movements of goal-directed attention by train-
604 ing on human search behavior. We found that the IRL model
605 was far more predictive of search efficiency than the Detector-
606 Hi-Low model or either of the Deep Search models, despite
607 the Deep Search models using methods considered to be state-
608 of-the-art in the fixation-prediction literature on free-viewing
609 behavior. In our state comparison we focused on the different
610 ways that a primate foveated retina, and its movement, might
611 be represented and used to train fixation prediction models.
612 We also extensively benchmarked COCO-Search18, both in
613 terms of the search behavior that it elicited, analyzed using
614 multiple behavioral measures and metrics, and in terms of
615 the predictive success of models ranging in their degree of
616 training on the COCO-Search18 behavior. All this means that
617 COCO-Search18 can be used immediately to start generating
618 new testable hypotheses. But likely the greatest contribution
619 of this work is yet to come. With a dataset the size and quality
620 of COCO-Search18, opportunities exist to explore new poli-
621 cies and reward functions for predicting goal-directed control
622 that have never before been possible²⁸. Our hope is that
623 COCO-Search18 will strengthen the bridge that human atten-
624 tion has built between the machine learning and behavioral
625 science literatures.

626 COCO-Search18 is now part of the MIT/Tuebingen
627 Saliency Benchmark, previously the MIT Saliency Bench-
628 mark but renamed to reflect the group that is now man-
629 aging the competition. The training, validation, and test
630 images in COCO-Search18 are already freely available as
631 part of COCO²⁹. Researchers are also free to see and use
632 COCO-Search18’s training and validation search fixations,
633 but the fixations on the test images are withheld. As part

634 of a managed benchmark, in a separate track it will be
635 possible to upload predictions and have them evaluated
636 on this test dataset. We invite you to participate in this
637 good-natured adversarial competition, and we hope that you
638 enjoy using COCO-Search18: [https://github.com/
639 cvlab-stonybrook/Scanpath_Prediction](https://github.com/cvlab-stonybrook/Scanpath_Prediction).

640 References

- 641 1. Itti, L., Koch, C. & Niebur, E. A model of saliency-
642 based visual attention for rapid scene analysis. *PAMI* **20**,
643 1254–1259 (1998).
- 644 2. Itti, L. & Koch, C. Computational modelling of visual
645 attention. *Nat. reviews neuroscience* **2**, 194–203 (2001).
- 646 3. Harel, J., Koch, C. & Perona, P. Graph-based visual
647 saliency. In *NIPS*, 545–552 (2007).
- 648 4. Borji, A., Sihite, D. N. & Itti, L. Quantitative analysis of
649 human-model agreement in visual saliency modeling: A
650 comparative study. *IEEE Transactions on Image Process.*
651 **22**, 55–69 (2012).
- 652 5. Borji, A. & Itti, L. State-of-the-art in visual attention
653 modeling. *PAMI* **35**, 185–207 (2012).
- 654 6. Kurylo, U. & Wilson, J. R. Using human eye gaze pat-
655 terns as indicators of need for assistance from a socially
656 assistive robot. In *International Conference on Social
657 Robotics*, 200–210 (Springer, 2019).
- 658 7. Admoni, H. & Srinivasa, S. Predicting user intent through
659 eye gaze for shared autonomy. In *2016 AAAI Fall Sympo-
660 sium Series* (2016).
- 661 8. Krishna Sharma, V., Saluja, K., Mollyn, V. & Biswas, P.
662 Eye gaze controlled robotic arm for persons with severe
663 speech and motor impairment. In *ACM Symposium on
664 Eye Tracking Research and Applications*, 1–9 (2020).
- 665 9. Buswell, G. T. *How people look at pictures: a study
666 of the psychology and perception in art.* (Univ. Chicago
667 Press, 1935).
- 668 10. Yarbus, A. L. Eye movements during perception of com-
669 plex objects. In *Eye Movements and Vision*, 171–211
670 (Springer, 1967).
- 671 11. Treisman, A. M. & Gelade, G. A feature-integration
672 theory of attention. *Cogn. Psychol.* **12**, 97–136 (1980).
- 673 12. Duncan, J. & Humphreys, G. W. Visual search and stim-
674 ulus similarity. *Psychol. Rev.* **96**, 433 (1989).
- 675 13. Chelazzi, L., Miller, E. K., Duncan, J. & Desimone, R. A
676 neural basis for visual search in inferior temporal cortex.
677 *Nature* **363**, 345–347 (1993).
- 678 14. Wolfe, J. M. Guided search 2.0 a revised model of visual
679 search. *Psychon. Bull. & Rev.* **1**, 202–238 (1994).
- 680 15. Najemnik, J. & Geisler, W. S. Optimal eye movement
681 strategies in visual search. *Nature* **434**, 387–391 (2005).
- 682 16. Torralba, A., Oliva, A., Castelhan, M. S. & Henderson,
683 J. M. Contextual guidance of eye movements and atten-
684 tion in real-world scenes: the role of global features in
685 object search. *Psychol. Rev.* **113**, 766 (2006).
- 686 17. Zelinsky, G. *et al.* Benchmarking gaze prediction for
687 categorical visual search. In *CVPR Workshops* (2019).
- 688 18. Eckstein, M. P. Visual search: A retrospective. *J. Vis.* **11**,
689 14,1–36 (2011).
- 690 19. Hollingworth, A. Guidance of visual search by memory
691 and knowledge. In *The Influence of Attention, Learning,
692 and Motivation on Visual Search*, 63–89 (Springer, 2012).
- 693 20. Wolfe, J. M. Visual search. In *The Handbook of Attention*,
694 27–56 (2015).
- 695 21. Treisman, A. & Souther, J. Search asymmetry: A diag-
696 nostic for preattentive processing of separable features. *J.
697 Exp. Psychol. Gen.* **114**, 285 (1985).
- 698 22. Judd, T., Ehinger, K., Durand, F. & Torralba, A. Learning
699 to predict where humans look. In *ICCV*, 2106–2113
700 (2009).
- 701 23. Kummerer, M., Wallis, T. S., Gatys, L. A. & Bethge,
702 M. Understanding low-and high-level contributions to
703 fixation prediction. In *ICCV*, 4789–4798 (2017).
- 704 24. Jia, S. & Bruce, N. D. Eml-net: An expandable multi-
705 layer network for saliency prediction. *Image Vis. Comput.*
706 103887 (2020).
- 707 25. Jiang, M., Huang, S., Duan, J. & Zhao, Q. Salicon:
708 Saliency in context. In *CVPR*, 1072–1080 (2015).
- 709 26. Liu, N. & Han, J. A deep spatial contextual long-term
710 recurrent convolutional network for saliency detection.
711 *IEEE Transactions on Image Process.* **27**, 3264–3274
712 (2018).
- 713 27. Cornia, M., Baraldi, L., Serra, G. & Cucchiara, R. Pre-
714 dicting human eye fixations via an lstm-based saliency
715 attentive model. *IEEE Transactions on Image Process.*
716 **27**, 5142–5154 (2018).
- 717 28. Yang, Z. *et al.* Predicting goal-directed human attention
718 using inverse reinforcement learning. In *CVPR*, 193–202
719 (2020).
- 720 29. Lin, T.-Y. *et al.* Microsoft coco: Common objects in
721 context. In *ECCV*, 740–755 (2014).
- 722 30. Redmon, J., Divvala, S., Girshick, R. & Farhadi, A. You
723 only look once: Unified, real-time object detection. In
724 *CVPR*, 779–788 (2016).
- 725 31. Liu, W. *et al.* Ssd: Single shot multibox detector. In
726 *ECCV*, 21–37 (2016).
- 727 32. Zhao, H., Shi, J., Qi, X., Wang, X. & Jia, J. Pyramid
728 scene parsing network. In *CVPR*, 2881–2890 (2017).
- 729 33. He, K., Gkioxari, G., Dollár, P. & Girshick, R. Mask
730 r-cnn. In *ICCV*, 2961–2969 (2017).

- 731 **34.** Caesar, H., Uijlings, J. & Ferrari, V. Coco-stuff: Thing
732 and stuff classes in context. In *CVPR*, 1209–1218 (2018).
- 733 **35.** Wolfe, J. M. What can 1 million trials tell us about visual
734 search? *Psychol. Sci.* **9**, 33–39 (1998).
- 735 **36.** Rosenholtz, R., Li, Y. & Nakano, L. Measuring visual
736 clutter. *J. Vis.* **7**, 17–17 (2007).
- 737 **37.** Neider, M. B. & Zelinsky, G. J. Cutting through the
738 clutter: Searching for targets in evolving complex scenes.
739 *J. Vis.* **11**, 7, 1–16 (2011).
- 740 **38.** Wolfe, J. M., Alvarez, G. A., Rosenholtz, R., Kuzmova,
741 Y. I. & Sherman, A. M. Visual search for arbitrary objects
742 in real scenes. *Attention, Perception, & Psychophys.* **73**,
743 1650 (2011).
- 744 **39.** Schmidt, J. & Zelinsky, G. J. Search guidance is propor-
745 tional to the categorical specificity of a target cue. *Q. J.*
746 *Exp. Psychol.* **62**, 1904–1914 (2009).
- 747 **40.** Castelhana, M. S., Pollatsek, A. & Cave, K. R. Typicality
748 aids search for an unspecified target, but only in identifica-
749 tion and not in attentional guidance. *Psychon. Bull. &*
750 *Rev.* **15**, 795–801 (2008).
- 751 **41.** Maxfield, J. T., Stalder, W. D. & Zelinsky, G. J. Effects
752 of target typicality on categorical search. *J. Vis.* **14**, 1,
753 1–11 (2014).
- 754 **42.** Henderson, J. M., Weeks Jr, P. A. & Hollingworth, A.
755 The effects of semantic consistency on eye movements
756 during complex scene viewing. *J. Exp. Psychol. Hum.*
757 *Percept. Perform.* **25**, 210 (1999).
- 758 **43.** Brockmole, J. R. & Henderson, J. M. Prioritizing new
759 objects for eye fixation in real-world scenes: Effects of
760 object–scene consistency. *Vis. Cogn.* **16**, 375–390 (2008).
- 761 **44.** Mills, M., Hollingworth, A., Van der Stigchel, S., Hoff-
762 man, L. & Dodd, M. D. Examining the influence of task
763 set on eye movements and fixations. *J. Vis.* **11**, 17,1–15
764 (2011).
- 765 **45.** Zhang, M. *et al.* Finding any waldo with zero-shot invari-
766 ant and efficient visual search. *Nat. communications* **9**,
767 1–15 (2018).
- 768 **46.** He, K., Zhang, X., Ren, S. & Sun, J. Deep residual learn-
769 ing for image recognition. In *CVPR*, 770–778 (2016).
- 770 **47.** Ng, A. Y., Russell, S. J. *et al.* Algorithms for inverse
771 reinforcement learning. In *ICML*, vol. 1, 663–670 (2000).
- 772 **48.** Abbeel, P. & Ng, A. Y. Apprenticeship learning via
773 inverse reinforcement learning. In *ICML*, vol. 1 (2004).
- 774 **49.** Ho, J. & Ermon, S. Generative adversarial imitation
775 learning. In *NIPS*, 4565–4573 (2016).
- 776 **50.** Schultz, W. Multiple reward signals in the brain. *Nat.*
777 *Rev. Neurosci.* **1**, 199–207 (2000).
- 778 **51.** Watanabe, K., Lauwereyns, J. & Hikosaka, O. Neural
779 correlates of rewarded and unrewarded eye movements
780 in the primate caudate nucleus. *J. Neurosci.* **23**, 10052–
781 10057 (2003).
- 52.** Montague, P. R., Hyman, S. E. & Cohen, J. D. Computa- 782
tional roles for dopamine in behavioural control. *Nature* 783
431, 760–767 (2004). 784
- 53.** Akbas, E. & Eckstein, M. P. Object detection through 785
search with a foveated visual system. *PLoS Comput. Biol.* 786
13, e1005743 (2017). 787
- 54.** Perry, J. S. & Geisler, W. S. Gaze-contingent real-time 788
simulation of arbitrary visual fields. In *Human Vision and* 789
Electronic Imaging, vol. 4662, 57–69 (2002). 790
- 55.** Irwin, D. E. Integrating information across saccadic eye 791
movements. *Curr. Dir. Psychol. Sci.* **5**, 94–100 (1996). 792
- 56.** Hollingworth, A. & Henderson, J. M. Accurate visual 793
memory for previously attended objects in natural scenes. 794
J. Exp. Psychol. Hum. Percept. Perform. **28**, 113 (2002). 795
- 57.** Parkhurst, D., Law, K. & Niebur, E. Modeling the role 796
of salience in the allocation of overt visual attention. *Vis.* 797
Res. **42**, 107–123 (2002). 798
- 58.** Navalpakkam, V. & Itti, L. Modeling the influence of 799
task on attention. *Vis. Res.* **45**, 205–231 (2005). 800
- 59.** Wang, Z. & Klein, R. M. Searching for inhibition of 801
return in visual search: A review. *Vis. Res.* **50**, 220–228 802
(2010). 803
- 60.** Zelinsky, G. J. A theory of eye movements during target 804
acquisition. *Psychol. Rev.* **115**, 787 (2008). 805
- 61.** Zelinsky, G. J., Rao, R. P. N., Hayhoe, M. M. & Ballard, 806
D. H. Eye movements reveal the spatiotemporal dynamics 807
of visual search. *Psychol. Sci.* **8**, 448–453 (1997). 808
- 62.** Zelinsky, G. J., Adeli, H., Peng, Y. & Samaras, D. Mod- 809
elling eye movements in a categorical search task. *Phi-* 810
los. Transactions Royal Soc. B: Biol. Sci. **368**, 20130058 811
(2013). 812
- 63.** Needleman, S. & Wunsch, C. A general method ap- 813
plicable to the search for similarities in the amino acid 814
sequence of two proteins. *Mol. Biol.* **48**, 443–153 (1970). 815
- 64.** Dewhurst, R. *et al.* It depends on how you look at it: 816
Scanpath comparison in multiple dimensions with multi- 817
match, a vector-based approach. *Behav. Res. Methods* **44**, 818
1079–1100 (2012). 819
- 65.** Anderson, N. C., Anderson, F., Kingstone, A. & Bischof, 820
W. F. A comparison of scanpath comparison methods. 821
Behav. Res. Methods **47**, 1377–1392 (2015). 822
- 66.** Zelinsky, G. J., Peng, Y., Berg, A. C. & Samaras, D. 823
Modeling guidance and recognition in categorical search: 824
Bridging human and computer object detection. *J. Vis.* 825
13, 30,1–20 (2013). 826
- 67.** Ehinger, K. A., Hidalgo-Sotelo, B., Torralba, A. & Oliva, 827
A. Modelling search for people in 900 scenes: A com- 828
bined source model of eye guidance. *Vis. cognition* **17**, 829
945–978 (2009). 830

- 831 **68.** Gilani, S. O. *et al.* Pet: An eye-tracking dataset for
832 animal-centric pascal object classes. In *2015 IEEE Inter-*
833 *national Conference on Multimedia and Expo (ICME)*,
834 1–6 (IEEE, 2015).
- 835 **69.** Everingham, M., Van Gool, L., Williams, C. K. I.,
836 Winn, J. & Zisserman, A. The PASCAL Visual
837 Object Classes Challenge 2012 (VOC2012) Results.
838 <http://host.robots.ox.ac.uk/pascal/VOC/index.html>.
- 839 **70.** Maxfield, J. T. & Zelinsky, G. J. Searching through
840 the hierarchy: How level of target categorization affects
841 visual search. *Vis. Cogn.* **20**, 1153–1163 (2012).
- 842 **71.** Papadopoulos, D. P., Clarke, A. D., Keller, F. & Ferrari,
843 V. Training object class detectors from eye tracking data.
844 In *ECCV*, 361–376 (2014).
- 845 **72.** Cerf, M., Harel, J., Einhäuser, W. & Koch, C. Predicting
846 human gaze using low-level saliency combined with face
847 detection. In *NIPS*, 241–248 (2008).
- 848 **73.** Treisman, A. & Gormican, S. Feature analysis in early
849 vision: evidence from search asymmetries. *Psychol. Rev.*
850 **95**, 15 (1988).
- 851 **74.** Neider, M. B. & Zelinsky, G. J. Exploring set size effects
852 in scenes: Identifying the objects of search. *Vis. Cogn.*
853 **16**, 1–10 (2008).
- 854 **75.** Zelinsky, G. J. Tam: Explaining off-object fixations
855 and central fixation tendencies as effects of population
856 averaging during search. *Vis. Cogn.* **20**, 515–545 (2012).
- 857 **76.** Schulman, J., Wolski, F., Dhariwal, P., Radford, A. &
858 Klimov, O. Proximal policy optimization algorithms.
859 *arXiv preprint arXiv:1707.06347* (2017).
- 860 **77.** Bylinskii, Z., Judd, T., Oliva, A., Torralba, A. & Durand,
861 F. What do different evaluation metrics tell us about
862 saliency models? *PAMI* **41**, 740–757 (2018).

Supplementary Materials

SM1: Behavioral Data Collection

Comparable datasets of search behavior

Figure S1 shows how COCO-Search18 compares to other large-scale datasets of search behavior. To our knowledge, there were only three such image datasets that were annotated with human search fixations^{17,67,68}. In terms of number of fixations, number of target categories, and number of images, COCO-Search18 is far larger. The PET dataset⁶⁸ collected search fixations for six animal target categories in 4,135 images selected from the Pascal VOC 2012 dataset⁶⁹, but the search task was non-standard in that participants were asked to “find all the animals” rather than search for a particular target category. This paradigm is therefore search at the superordinate categorical level, which is far more weakly guided than basic-level search⁷⁰. Gaze fixations were also recorded for only 2 seconds/image, and multiple targets often appeared in each scene. The microwave-clock search dataset (MCS¹⁷) is our own work and a predecessor of COCO-Search18. In collecting data for the 18 target categories in COCO-Search18 we had to start somewhere, and our first two categories were microwaves and clocks (although the datasets differed for even those two categories due to the use of different exclusion criteria). Until recently, perhaps the best dataset of search fixations was from⁶⁷, but it is relatively small, limited to only the search for people in scenes, and is now a decade old. Note that, whereas there are larger datasets with respect to free-viewing fixations (SALICON²⁵) or fixations collected using other visual tasks (POET⁷¹), these tasks were not visual search and therefore these datasets cannot be used to train models of search behavior. These collective inadequacies demanded the creation of a newer, larger, and higher-quality dataset of search fixations, enabling deep network models to be trained on people’s movements of attention as they pursue target-object goals.

Selection of target categories and search images

Here we more fully describe how we selected from COCO’s trainval dataset²⁹ the 18 target categories and the 6,202 images included in COCO-Search18. A goal in implementing our selection criteria was to elicit the behavior that we are trying to measure, namely, the guidance of search fixations by a target category. We also put care into excluding images that might elicit other gaze patterns that would introduce noise with respect to identifying the target-control signal. This sort of attention to detail is uncommon in datasets created for the training of deep network models, where the approach seems to be “the more images the better”. But whereas this is usually true because more images leads to better-trained models, in creating a dataset of human behavior this more-is-better impulse should be tempered with some quality control to be confident that the behavior is of the purported type. In the current context this behavior should be search fixations that are guided to the target, because search fixations that are unguided have less value as training labels. Because a standard

search paradigm collects behavioral responses for both TP and TA images, separate selection criteria were needed. All image selection was based on object labels and/or bounding boxes provided by COCO. On this point, while inspecting the images that were ultimately selected we noticed that exemplars in some categories were mislabeled, probably due to poor rater agreement on that category. For instance, several chair exemplars were mislabeled as couches, and vice versa. Rather than attempting to correct these mislabels, which would be altering COCO, we decided to keep them and tolerate a higher-than-normal error rate for the affected categories. This action seemed best, given our plan to discard error trials from the search performance analyses in our study, but researchers interested in interpreting button press errors in COCO-Search18 should be aware of this labeling issue.

Target-present image selection. Six criteria were imposed on the selection of images to be used for target-present search trials.

- (1) Images were excluded if they depicted people or animals. We did this to avoid the known biases to fixate on these objects when they appear in a scene^{22,72}. Such biases would compete with guidance from target-category features, thereby distorting study of the target-bias that is more central to search.
- (2) Images were excluded if they depicted multiple instances of the target. A scene showing a classroom with many chairs would therefore be excluded from the “chair” target category because one, and only one, instance of a chair would be allowed in an image.
- (3) Images were excluded if the size of the target, measured by the area of its bounding box, was smaller than 1% or larger than 10% of the total image area. This was done to create searches that were not too hard or too easy.
- (4) Images were excluded if the target appeared at the image center, based on a 5×5 grid. We did this because the participant’s gaze was pre-positioned at this central location at the start of each search trial.
- (5) Images were excluded if their width/height ratio fell outside the range of 1.2-2.0 (based on a screen ratio of 1.6). This criterion excluded very elongated images, which we thought might distort normal viewing behavior.
- (6) Images, and entire image categories, were excluded if the above criteria left fewer than 100 images per object category. We did this because fewer than 100 images would likely be insufficient for training and testing a deep network model specific to that object category.

Applying these exclusion criteria left 32 object categories from COCO’s original 80. Given that this left still far too many images for people to practically annotate with search fixations, we decided to attempt exclusion of images where targets were highly occluded or otherwise difficult to recognize. We did this out of concern that such images would largely introduce noise into the search behavior. To do this,

we trained object detectors on cropped views of these 32 categories, and excluded images if the object bounding boxes had a classification confidence $< .99$. Specifically, for these 32 categories we created a validation set consisting of images meeting the selection criteria and a training set consisting of the images that did not. The bounding box of the object, for each of the 32 object classes, was then cropped in the image to obtain the positive training samples. Negative samples were same-sized image patches that had 25% intersection with the target (area of intersection divided by area of target), meaning that they were class-specific hard negatives. All cropped patches (over 1 million) were resized to 224×224 pixels while maintaining the aspect ratio using padding. The classifier was a ResNet50 pre-trained on ImageNet, which we fine-tuned by dilating the last fully-connected layer and re-training on 33 outputs (32+“Negative”). Images were excluded if the cropped object patch had a classification score of less than .99. This procedure resulted in 18 categories with at least 100 images in each category, totaling 3,131 TP images.

Two final exclusion criteria were implemented by manual selection. First, for the clock target category we included only images of analog clocks, meaning that we excluded digital clocks from being clock targets. We did this because the features of analog and digital clocks are highly distinct and very different, and we were concerned that this would introduce variability in the search behavior and reduce data quality. Five images depicting only digital clocks were excluded for this reason. Lastly, images from all 18 of the target categories were screened for objectionable content, which we defined as offensive content or content evoking discomfort or disgust. The “toilet” category had the most images (17) excluded for objectionable content, with a total of 25 images excluded across all target categories. After implementing all exclusion criteria discussed in this section, we obtained 3,101 TP images from 18 categories: bottle, bowl, car, chair, (analog) clock, cup, fork, keyboard, knife, laptop, microwave, (computer) mouse, oven, potted plant, sink, stop sign, toilet, and tv. See Figure 2 for the specific number of images in each category.

Target-absent image selection. To balance the selection of the 3,101 TP images, we selected an equal number of TA images from COCO. To do this, we kept the criteria excluding images depicting people or animals, extreme width/height image ratios, and images with objectionable content, all as described for the TP image selection, but added two more exclusion criteria that were specific to each of the 18 target-object categories.

- (1) Images were excluded if they depicted an instance of the target, a prerequisite for a TA image.
- (2) Images were excluded if they depicted less than two instances of the target category’s siblings, a criterion introduced to discourage searchers from making TA responses purely on the basis of scene type. For example, a person might be biased to make a TA response if they are searching for a toilet target and the image is a street scene.

Because COCO has a hierarchical organization, parent, child, and sibling relationships can be used for image selection. For example, COCO defines the siblings of a microwave to be an oven, toaster, refrigerator, and sink, all under the parent category of appliance. By requiring that the TA scenes for a target category have at least two of that category’s siblings, we impose a sort of scene constraint that minimizes target-scene inconsistency and makes a scene appropriate to use as a TA image. A scene that has an oven and a refrigerator is very likely to be a kitchen, thereby making it difficult to answer on the basis of scene type alone whether a microwave target is present or absent.

These exclusion criteria still left us with many thousands more TA images than we needed, so we sampled randomly within each of the 18 target categories to match the 3,101 TP images.

Order of target-category presentation

Collecting the search behavior for 6,202 images required dividing each participant’s effort into six days of testing. Each testing session was conducted on a different day, lasted about 2 hours, and consisted of about 1000 search trials, evenly divided between TP and TA. Because images from different categories can overlap (e.g., images depicting a microwave may also depict an oven), the presentation order of the target-category blocks was constrained to minimize the repetition of images in consecutive categories and consecutive sessions. For example, because 49 images satisfied the selection criteria for both the sink and microwave target categories, we prevented the microwave and sink categories from appearing in, not only the same session, but the sessions preceding and following. We did this to minimize possible biases resulting from seeing the same scene in different search contexts. A heuristic for maximizing this distance between repeating images resulted in the following fixed target category presentation order across the six sessions:

- (1) tv + sink;
- (2) fork + chair;
- (3) car + bowl + potted plant + mouse;
- (4) knife + keyboard + oven + clock;
- (5) cup + laptop + toilet;
- (6) bottle + stop sign + microwave.

Each participant viewed from Session 1 to Session 6, or from Session 6 to Session 1, with this order counterbalanced across participants.

Data-collection procedure

Participants were 10 Stony Brook University undergraduate and graduate students, 6 males and 4 females, with ages ranging from 18–30 years. All had normal or corrected to normal vision, by self report, were naive with respect to task design and paradigm when recruited, and were compensated with course credit or money for their participation. Informed consent was obtained from each participant at the beginning of testing, in accordance with the Institutional Review Board

1078 responsible for overseeing human-subjects research at Stony
1079 Brook University.

1080 The target category was designated to participants at the
1081 start of each block. This was done using the type of display
1082 shown in Figure S2 for the potted-plant and analog clock
1083 categories. The name of the target category was shown in
1084 text at the top, with examples of objects that would, or would
1085 not, qualify as exemplars of the named category. In selecting
1086 exemplars to illustrate as positive target-category members,
1087 we attempted to capture key categorical distinctions at a level
1088 immediately subordinate to the target category. When needed,
1089 we also gave negative examples by placing a red X through
1090 the object. We did this to minimize potential confusions and
1091 to enable the participant to better define the target category's
1092 boundary.

1093 The procedure (Figure S3) on each trial began with a fixa-
1094 tion dot appearing at the center of the screen. To start a trial,
1095 the participant would press the "X" button on a game-pad con-
1096 troller while carefully looking at the fixation dot. An image
1097 of a scene would then be displayed and the participant's task
1098 would be to answer, "yes" or "no", whether an exemplar of the
1099 target category appears in the displayed scene by pressing the
1100 right or left triggers of the game-pad, respectively. The search
1101 scene remained visible until the manual response. Participants
1102 were told that there were an equal number of TP and TA trials,
1103 and that they should make their responses as fast as possible
1104 while maintaining high accuracy. No accuracy or response
1105 time feedback was provided.

1106 The presentation of images during the experiment was con-
1107 trolled by Experiment Builder (SR research Ltd., Ottawa,
1108 Ontario, Canada). Stimuli were presented to participants on
1109 a 22-inch LCD monitor (1680×1050 pixel resolution) at a
1110 viewing distance of 47cm from the monitor, enforced by chin
1111 and head rests. These viewing conditions resulting in hori-
1112 zontal and vertical visual angles of $54^\circ \times 35^\circ$, respectively.
1113 Participants were asked to keep their gaze on the fixation point
1114 at the start of each trial, but were told that they should feel free
1115 to move their eyes as they searched. Eye movements were
1116 recorded throughout the experiment using an EyeLink 1000
1117 eye-tracker in tower-mount configuration (SR research Ltd.,
1118 Ottawa, Ontario, Canada). Eye-tracker calibrations occurred
1119 before every block or whenever necessary, and these 9-point
1120 calibrations were not accepted unless the average calibration
1121 error was $\leq .51^\circ$ and the maximal error was $\leq .94^\circ$. The ex-
1122 periment was conducted in a quiet laboratory room under dim
1123 lighting conditions.

1124 SM2: Behavioral evaluation of COCO-Search18

1125 *Effects of set size and target eccentricity*

1126 The visual search literature has done excellent work in identi-
1127 fying many of the factors that increase search difficulty (for
1128 reviews, see: [12, 18, 60, 73](#)). Larger set sizes (number of items in
1129 the search display), smaller target size, larger target eccentricity,
1130 and greater target-distractor similarity are all known to
1131 make search more difficult. However, most of this work was

1132 done in the context of simple stimuli, and generalization to
1133 realistic images is challenging. For example, what to consider
1134 an object in a scene is often unclear, making it difficult to de-
1135 fine a set size⁷⁴. Objects in images also do not usually come
1136 annotated with labels and bounding boxes. These problems of
1137 object segmentation and identification, which largely do not
1138 exist for search studies using object arrays, become significant
1139 obstacles to research when scaled up to images of scenes.

1140 With COCO-Search18, we can begin to ask how the search
1141 for targets in images is affected by set size and target eccen-
1142 tricity. Set size is determined based on the COCO object and
1143 stuff labels, which collectively map every pixel in an image
1144 to an object or stuff category. Set size is the count of the
1145 number of these labels for a given image. Figure S4 shows
1146 the relationship between the number of fixations made on an
1147 image, averaged over participants, and the set size of that im-
1148 age, grouped by target category. Some target categories, such
1149 as laptop, oven, microwave, and potted-plant, have significant
1150 positive set size effects ($r = .21$ to $.37$, $ps \leq .01$), indicating
1151 a less efficient search with more objects. A similar pattern is
1152 shown in Figure S5 for the relationship between the number of
1153 fixations on a search image and the initial visual eccentricity
1154 of the target (distance between the image center and the target
1155 bounding-box center), where for these same objects there was
1156 a decrease in search efficiency with increasing target eccen-
1157 tricity. For other target object categories, such as: stop sign,
1158 fork, and keyboard, search efficiency was unaffected by either
1159 set size or target eccentricity ($ps > .05$), possibly because
1160 these objects are either highly salient (stop sign) or highly
1161 constrained by scene context (keyboard).

1162 *Distance between search fixations and the target*

1163 How much closer does each search fixation bring gaze to
1164 the target? We analyzed this measure of search efficiency
1165 and report the results in Figure S6. Plotted is the Euclidean
1166 distance between the target location and the locations of the
1167 starting fixation (0) and the fixation locations after the first six
1168 eye movements (1-6). The most salient pattern is the rapid
1169 decrease in fixation-target distance in the first two new fix-
1170 ations, which dovetails perfectly with the steep increase in
1171 the cumulative probability of target fixation over these same
1172 eye movements reported in Figure 4A. From a starting lo-
1173 cation near the center of the image, these eye movements
1174 brought gaze steadily closer to the target. Note that because
1175 this fixation-target distance is averaged over images and partic-
1176 ipants, the roughly 5 degrees of visual angle at the bottom of
1177 these functions should not be misinterpreted as gaze being this
1178 distance from the target on a given trial. More interpretable
1179 are the overall trends, where a steep drop in distance is fol-
1180 lowed by a plateau, or even a smaller increase in distance with
1181 the 5th and 6th new fixations. This small increase is likely an
1182 artifact of these 5 and 6-fixation trials being the most difficult,
1183 with more idiosyncratic search behavior.

1184 **Target-absent search fixations**

1185 In the main text we focused on the TP data, where the guid- 1239
1186 ance signal is clearer and the modeling goals are better defined, 1240
1187 but we conducted largely parallel analyses of the TA data. Fig- 1241
1188 ure S7A shows representative TA images with fixation data 1242
1189 from one participant, and Figure S7B shows FDMs from all 1243
1190 participants for the same images. Comparing these data with 1244
1191 the TP data from Figure 1, it is clear that people made many 1245
1192 more fixations in the absence of a target. This was expected 1246
1193 from the search literature, but it should also be noted that the 1247
1194 FDMs are still much sparser than what would be hypothesized 1248
1195 by an exhaustive search. Paralleling Figure 3, in Figure S8 we 1249
1196 report applicable analyses of the TA search behavior. These 1250
1197 are grouped by manual accuracy and response time, and the 1251
1198 mean number of fixations made before the target-absent but- 1252
1199 ton press terminating a trial. Note that accuracy was high 1253
1200 (low false positive error rate) for all of the target categories 1254
1201 except chairs and cups, with the reason for the former already 1255
1202 discussed in the context of mislabeling and the reason for the 1256
1203 latter likely reflecting an occasionally challenging category 1257
1204 distinction (e.g., some bottles can look like some cups). Also 1258
1205 note that there was an average of only five fixations made 1259
1206 during search, even on the TA search trials. As in Figure 5, 1260
1207 Figure S9 visualizes the agreement and other patterns among 1261
1208 these measures. The rows show ranked performance, with 1262
1209 dark red indicating more difficult (or least efficient) search 1263
1210 and dark blue indicating relatively easy or efficient search. 1264
1211 The columns in Figure S9A group the measures by target 1265
1212 category. Similar to the TP data, there was again good con- 1266
1213 sistency among the measures. Also consistent is the fact that 1267
1214 bottles and cups were among the most difficult target cate- 1268
1215 gories, whereas the toilet category was the easiest. There was 1269
1216 also evidence in the TA data for a speed-accuracy trade-off 1270
1217 for some target categories. For example, microwaves and stop 1271
1218 signs had relatively low error rates, but these categories were 1272
1219 searched with relatively high effort, as measured by ranked 1273
1220 response time and number of fixations. Figure S9B visualizes 1274
1221 the measures by participant instead of category, where we 1275
1222 again found individual differences between participants in 1276
1223 search efficiency. 1277

1224 **Practice effects**

1225 Each of the participants contributing to COCO-Search18 1278
1226 searched more than 6000 images, making it possible to ana- 1279
1227 lyze how their search efficiency improved with practice. Fig- 1280
1228 ure S10 shows practice effects for both response time (top) 1281
1229 and the number of fixations before the button press (bottom), 1282
1230 where we define practice effects as performance on the first 1283
1231 1/3 of the trials compared to performance on the last 1/3 of the 1284
1232 trials for each target category. Practice effects were larger for 1285
1233 TA trials (right) than for TP trials (left), noting the differences 1286
1234 in y-axis scales, and that considerable differences existed 1287
1235 across categories. Some categories, such as bottles, showed 1288
1236 large practice effects, while other categories, such as analog 1289
1237 clocks, showed none at all. We speculate that this difference is 1290
1238 due to some categories requiring more exemplars to fully learn 1291
1292

1239 compared to others. For example, analog clock was perhaps 1240
1241 the most well defined of COCO-Search18's categories, and 1242
1243 bottle certainly one of the least well defined, creating greater 1244
1245 opportunity to better learn the bottle category with practice 1246
1247 over trials. 1248

1244 **Search fixation durations**

1245 Figures S11 and S12 show density histograms of the search 1246
1247 fixation durations for the TP and TA data, respectively, plot- 1248
1249 ted for each of the target categories. Fixation durations are 1250
1251 plotted across the x-axes with a bin size of 50ms, and y-axes 1252
1253 show the normalized probability density at each fixation. Of 1254
1255 note in the TP data is that the mode initial fixation durations 1255
1256 (blue lines) were a bit longer than the mode duration of the 1256
1257 rest (averaged mode difference = 63ms), consistent with the 1257
1258 very strong guidance observed in the initial eye movements, 1258
1259 and they tended to have more bi-modal distributions. The 1259
1260 main peak was at ~250 ms, with a smaller and very short- 1260
1261 latency peak at ~50 ms that is likely a truncation artifact of 1261
1262 fixation duration being measured relative to the onset of the 1262
1263 search display. In contrast, the distributions of second fixa- 1263
1264 tions (orange lines) were consistently shorter, even relative to 1264
1265 the subsequent fixations. Speculatively, this may be due to 1265
1266 a greater proportion of the first new fixations being “off ob- 1266
1267 ject”⁷⁵, which are often followed by short-latency corrective 1267
1268 saccades that bring gaze accurately to an object. This inter- 1268
1269 pretation is consistent with the high probability of the target 1269
1270 being fixated by the second eye movement (Figure 4A). As 1270
1271 for the subsequent fixations, they tended to be short (~200ms) 1271
1272 and not highly variable in their durations. The TA fixations 1272
1273 showed similar trends, except for the durations of the second 1273
1274 fixations no longer differing from the rest. 1274

1270 **Saccade amplitudes**

1271 We also analyzed the distribution of saccade amplitudes dur- 1271
1272 ing visual search, defined here as the Euclidean distance be- 1272
1273 tween consecutive fixations in visual angle. Figure S13 and 1273
1274 Figure S14 show the distributions of saccade amplitudes in 1274
1275 the TP and TA data, respectively. In the TP data, saccade 1275
1276 amplitudes were larger in some categories (toilet and stop 1276
1277 sign) than others (bottle and potted plant), likely because eas- 1277
1278 ier target categories could be identified from farther in the 1278
1279 visual periphery. There was also evidence for bimodality in 1279
1280 the amplitude distributions, shown most clearly for clocks, 1280
1281 forks, stop signs, and tvs. We speculate that this bimodal- 1281
1282 ity reflects larger-amplitude exploratory saccades mixed with 1282
1283 smaller-amplitude saccades used in the verification of an ob- 1283
1284 ject category. Mean saccade amplitudes in the TA data were 1284
1285 clearly larger than for the TP data ($t(17) = 11.79, p < .001$), 1285
1286 and this difference was consistent across target categories (all 1286
1287 $ps \leq .001$). We attribute this to the relatively large viewing 1287
1288 angle of the search displays (54×35 degrees of visual angle) 1288
1289 creating a greater need for exploration, but this is also specula- 1289
1290 tion. The distributions of saccade amplitudes were also more 1290
1291 consistent across categories in the TA data, with there being 1291
1292 weaker evidence of bi-modality. 1292

1293 **SM3: Model Methods**

1294 ***Training and testing datasets***

1295 Model success depends on the training dataset being an accurate
1296 reflection of the test dataset. When the training dataset
1297 includes a behavioral annotation, as does COCO-Search18, it
1298 is therefore important to know that similar patterns exist in
1299 the training and testing search behavior. The analyses shown
1300 in Figure 5A included images from all of COCO-Search18,
1301 which recall were randomly split into 70% for training, 10%
1302 for validation, and 20% for testing. Figure S15 replots the
1303 data from Figure 5A, but divides it into the training/validation
1304 (left) and testing (right) datasets. Note the high agreement
1305 between the testing and train/val datasets across this battery
1306 of behavioral performance measures.

1307 ***Inverse Reinforcement Learning***

1308 The specific inverse-reinforcement learning (IRL) method
1309 that we used was generative adversarial imitation learning
1310 (GAIL⁴⁹) with proximal policy optimization (PPO)⁷⁶. The
1311 model policy is a generator that aims to create state-action
1312 pairs that are similar to human behavior. The reward function
1313 (the logarithm of the discriminator output) maps a state-action
1314 pair to a numeric value. The generator and discriminator are
1315 trained within an adversarial optimization framework to obtain
1316 the policy and reward functions. The discriminator's task is
1317 to distinguish whether a state-action pair was generated by
1318 a person (real) or by the generator (fake), with the generator
1319 aiming to fool the discriminator by maximizing the similarity
1320 between its state-action pairs and those from people. The
1321 reward function and policy that are learned from the fixation-
1322 annotated images during training are then used to predict new
1323 search fixations in the unseen test images.

1324 **SM4: Performance metrics and model evaluation**

1325 ***Metrics for comparing search efficiency and scanpaths***

1326 We considered five metrics for quantifying search efficiency
1327 and comparing search scanpaths (Table 1). Two metrics for
1328 quantifying search efficiency follow directly from the group
1329 target-fixation probability (TFP) function shown in Figure 4.
1330 The first of these computes the area under the TFP curve, a
1331 metric we refer to as TFP-auc. Second, and relatedly, we
1332 compute the sum of the absolute differences between the hu-
1333 man and model target-fixation-probabilities in a metric that
1334 we refer to as Probability Mismatch. A third metric for quan-
1335 tifying overt search efficiency is Scanpath Ratio. It is the
1336 Euclidean distance between the initial fixation location and
1337 the target divided by the summed Euclidean distances between
1338 the fixation locations in the search scanpath⁴². It is an effi-
1339 ciency metric because an initial saccade that lands directly
1340 on the target would give a Scanpath Ratio of 1, meaning that
1341 the distance between starting fixation and the target would
1342 be the same as the summed saccade distance. These three
1343 metrics emphasize target-fixation efficiency by penalizing ei-
1344 ther the number of fixations or the saccade-distance traveled
1345 to achieve the target goal. The final two metrics focus on
1346 scanpath comparison, and specifically comparing the search

scanpaths between people and the models. The first of these
scanpath-comparison metrics computes a Sequence Score by
first converting a scanpath into a string of fixation cluster IDs,
and then using a string matching algorithm⁶³ to measure the
similarity between the two strings. Figure S16 shows exam-
ples of behavioral and model scanpaths and their sequence
scores to develop an intuition for this metric. Lastly, we use
MultiMatch^{64,65} to measure the scanpath similarity at the
pixel level. MultiMatch measures five aspects of scanpath
similarity: shape, direction, length, position, and duration.
We excluded the duration measure from our use of this metric
because the models in our comparison group did not predict
fixation duration. See Table S2 for the results of statistical
tests comparing predictions from each pair of models.

1361 ***Comparing predicted and behavioral fixation-density maps (FDMs)***

1362 Search has a temporal dynamic, making a metric for capturing
1363 the spatio-temporal sequence of fixations preferred over ones
1364 that compare only FDMs, where this temporal component is
1365 disregarded. However, the prediction of FDMs is common
1366 for free-viewing tasks, and because there is no technical rea-
1367 son why FDM metrics cannot be applied to search we do so
1368 here in the hope that the visual saliency literature finds this
1369 comparison useful. Models generated scanpaths having a max-
1370 imum length of 6 new fixations, but FDMs were constructed
1371 only from those fixations leading up to the first fixation on
1372 the target, just as FDMs were constructed from the behav-
1373 ioral fixations. We used three widely accepted metrics for
1374 comparing predicted against observed FDMs. Area Under
1375 the Receiver Operating Characteristic Curve (AUC) uses a
1376 predicted priority map as a binary classifier to discriminate
1377 behavioral fixation locations from non-fixated locations. Nor-
1378 malized Scanpath Saliency (NSS) finds the model predictions
1379 at each of the behavioral fixation locations, then averages and
1380 normalizes these values. Lastly we computed a Pearson's
1381 Correlation Coefficient (CC) between the predicted and be-
1382 havioral FDMs, although this metric reflects only the degree
1383 of linear relationship between predicted and behavioral FDMs
1384 (for additional discussion, see: Borji & Itti⁵; Bylinskii et al.⁷⁷).
1385 Table S1 reports the results of an evaluation comparing model
1386 predictions of search FDMs to behavioral search FDMs using
1387 each of these metrics. The findings that we report in the main
1388 text in the context of scanpath prediction also hold in the case
1389 of FDM prediction. Specifically, the IRL-Hi-Low-C model
1390 outperformed the others, and did so for all three metrics. Ad-
1391 ditionally, the Detector-Hi model also performed relatively
1392 well in all the metrics, supporting our conclusion that a simple
1393 detector does a relatively good job in predicting fixations in
1394 visual search.
1395

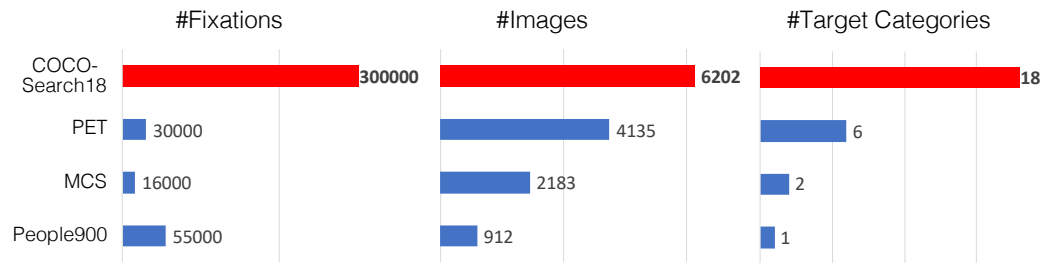


Figure S1. Comparisons between COCO-Search18 and other large-scale datasets of search behavior. COCO-Search18 is the largest in terms of number of fixations (~300,000), number of target categories (18), and number of images (6,202).

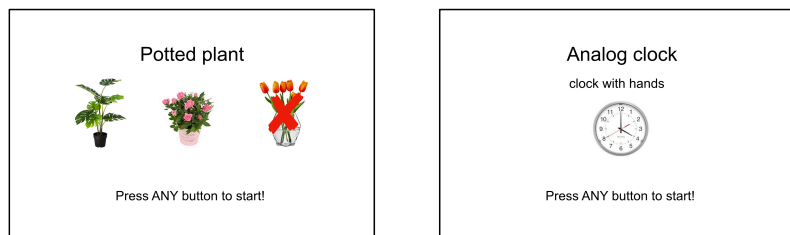


Figure S2. Examples of target-designation displays, shown for the potted-plant and analog clock targets, that preceded the block of trials for a given target category.

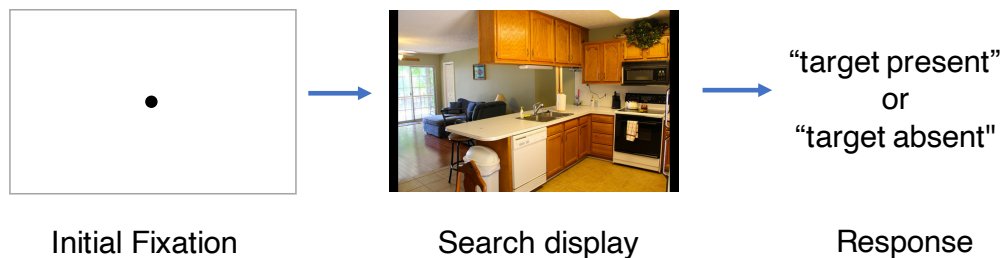


Figure S3. Example of the search procedure. Each trial began with a fixation dot appearing at the center of the screen. Participants would start a trial by pressing a button on a game-pad controller while carefully looking at the fixation dot. An image of a scene would then be displayed and the participant's task was to make a speeded "yes" or "no" target-presence judgment by pressing the right or left triggers, respectively, of a game-pad controller.

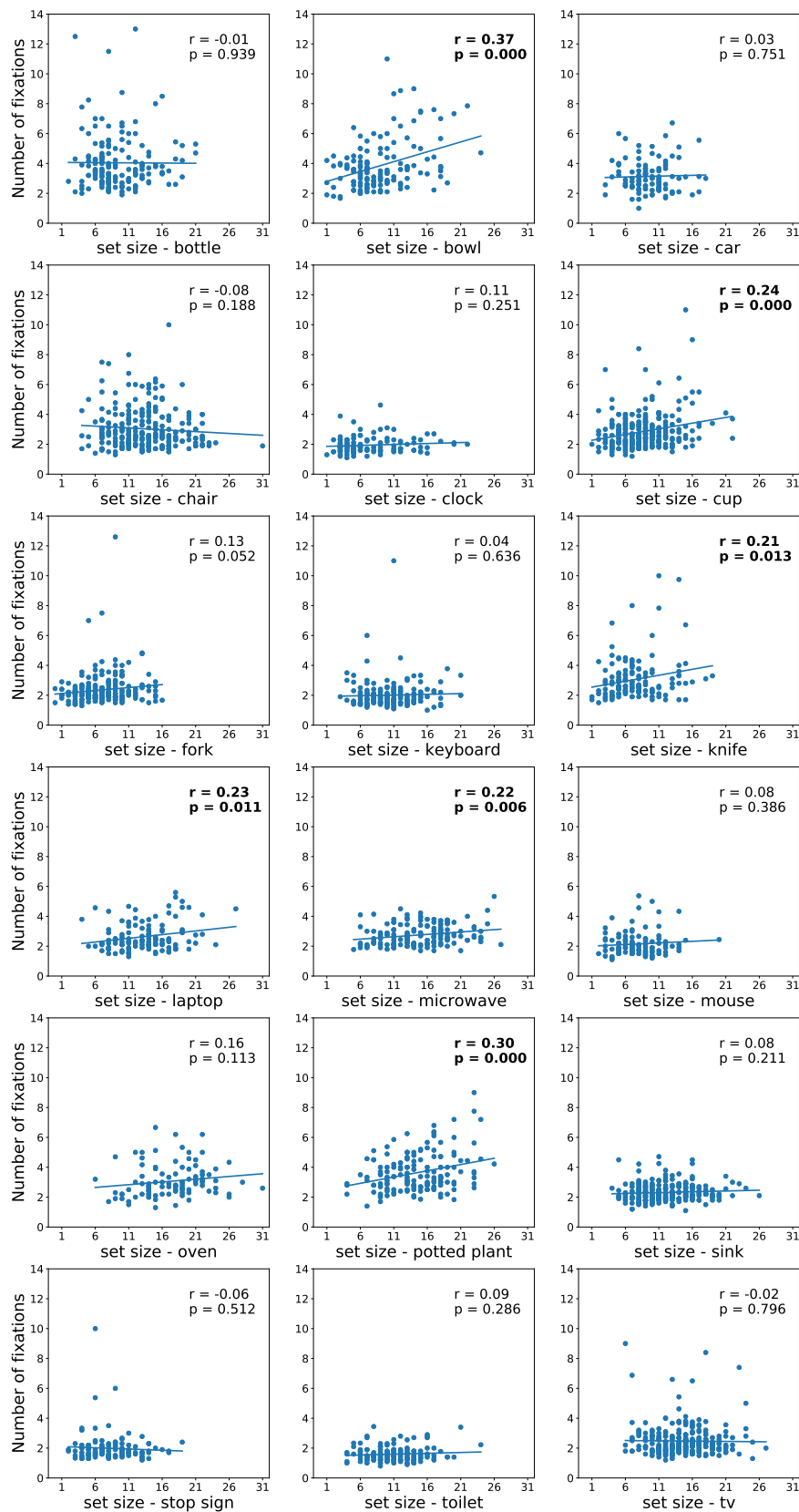


Figure S4. Number of fixations made on the target-present images plotted as a function of the set sizes of those images (using COCO object and stuff labels), averaged over participants and grouped by target category.

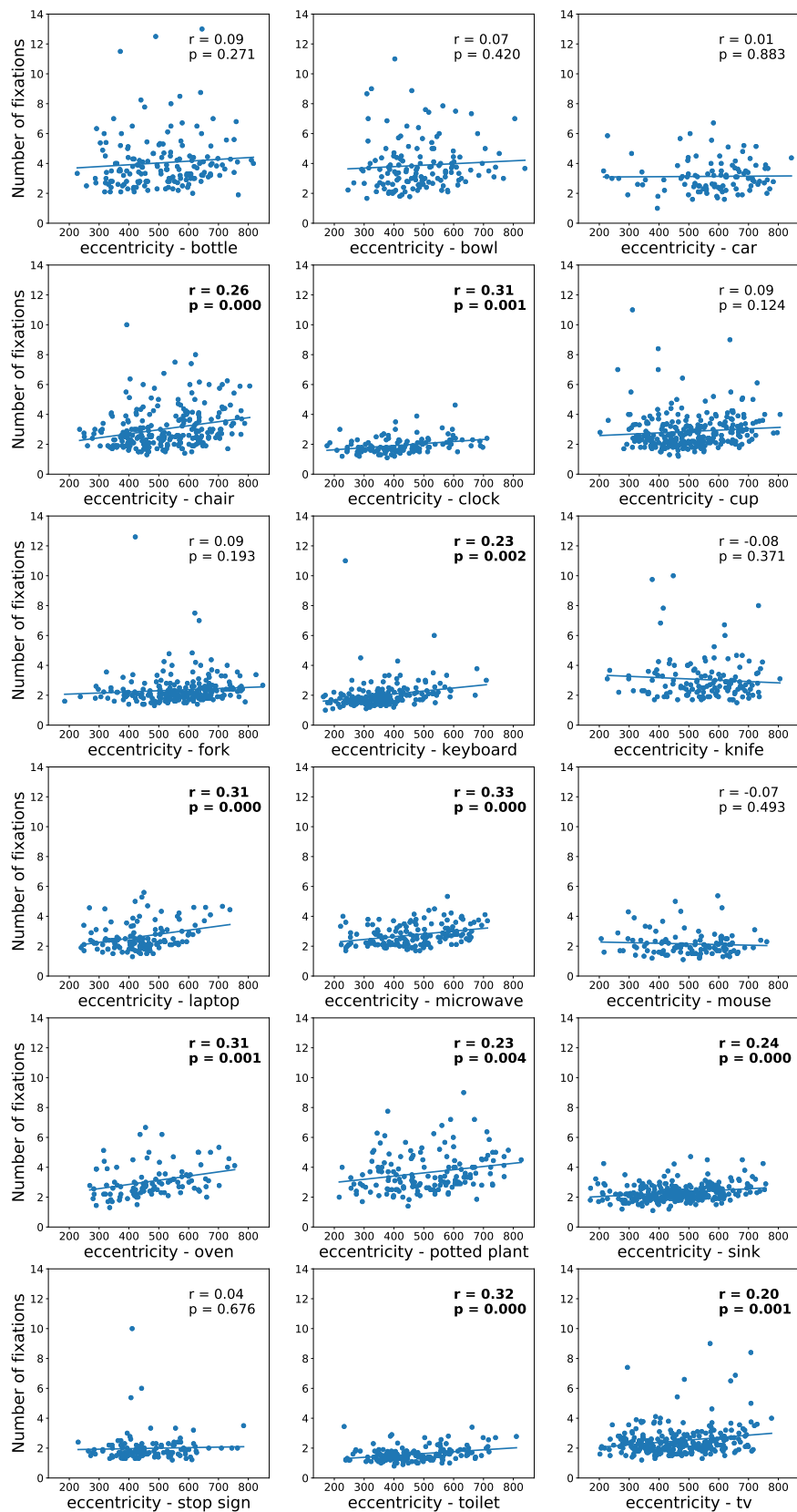


Figure S5. Number of fixations made on the target-present images plotted as a function of initial target eccentricity (using the center of the COCO bounding-box), averaged over participants and grouped by target category.

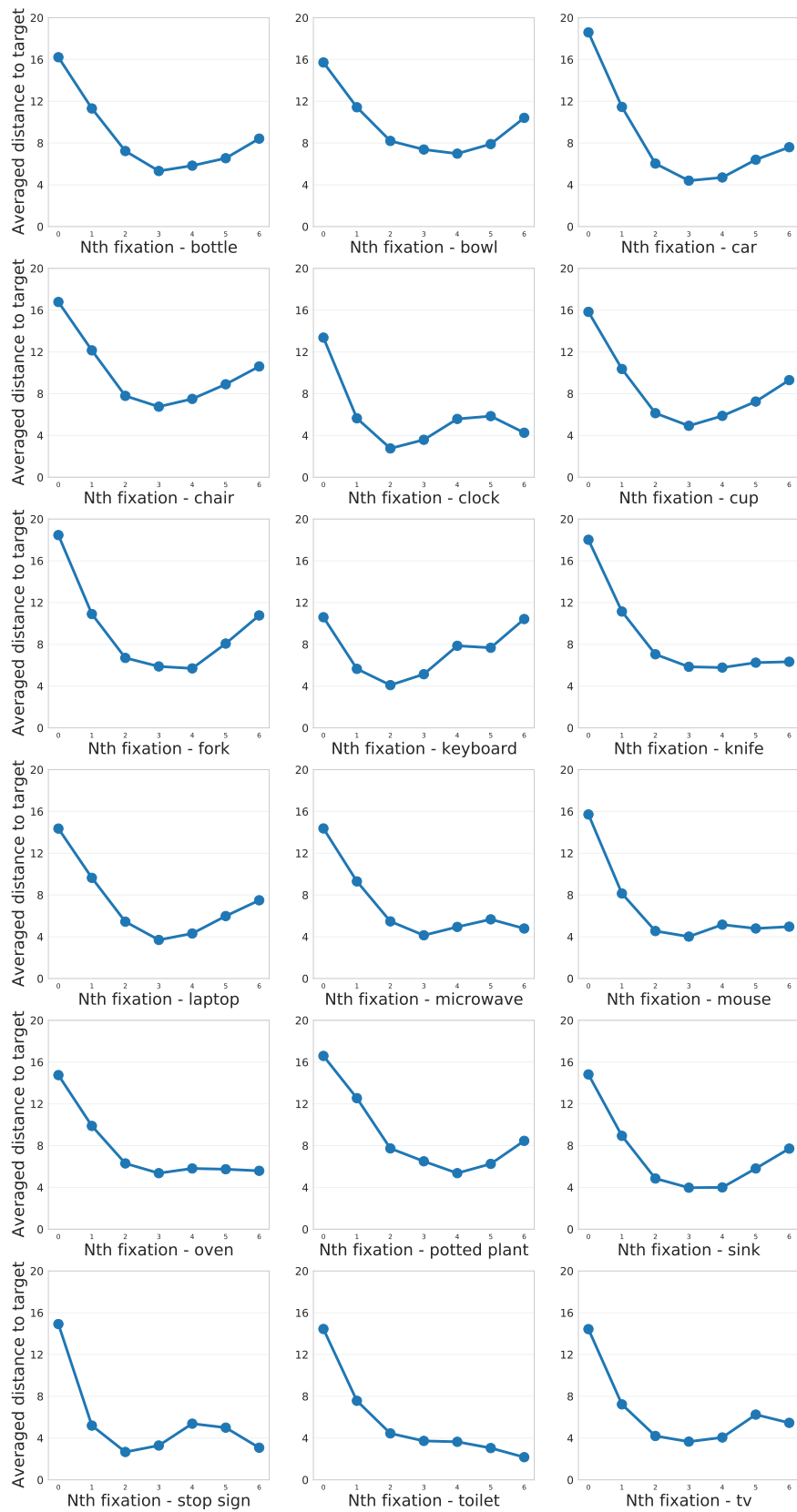
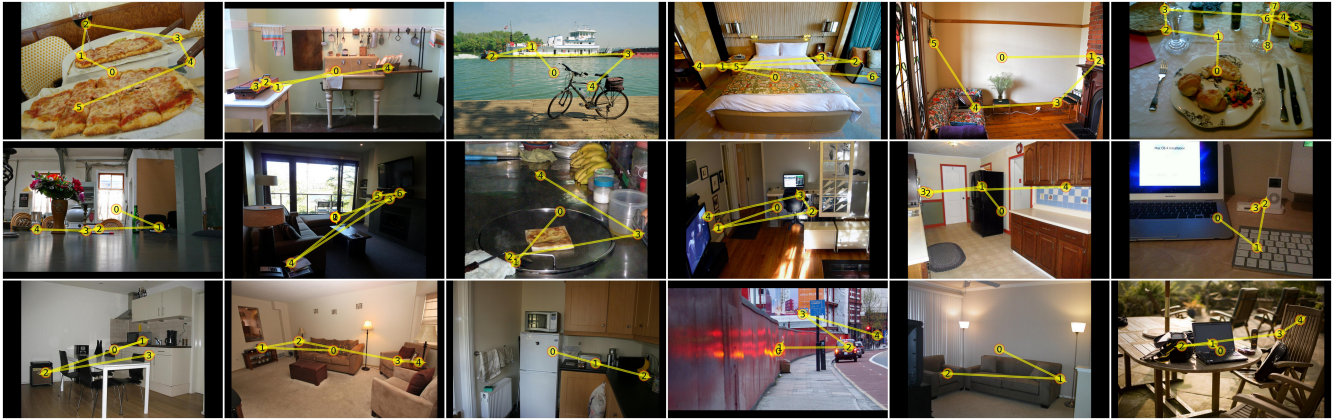


Figure S6. Averaged Euclidean distance (in visual angle) between gaze and the target's center (using COCO bounding-box labels) over the first 6 saccades, grouped by target category.

A



B

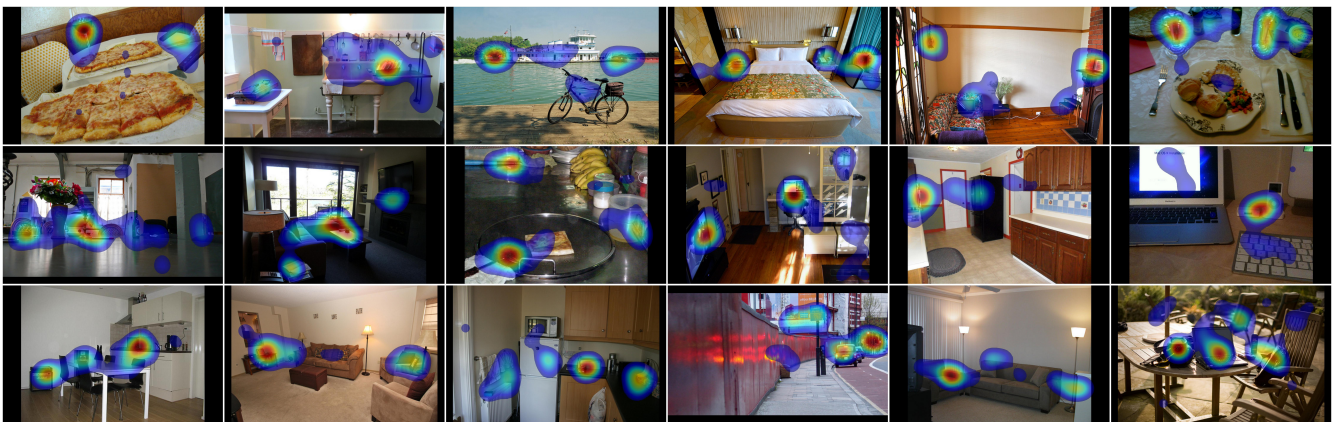


Figure S7. (A). Examples of a target-absent image for each of the 18 target categories. Yellow lines and numbered discs indicate a representative search scanpath from a single participant. From left to right, top to bottom: bottle, bowl, car, chair, (analog) clock, cup, fork, keyboard, knife, laptop, microwave, mouse, oven, potted plant, sink, stop sign, toilet, tv. (B). Examples of fixation density maps for the same target-absent images.

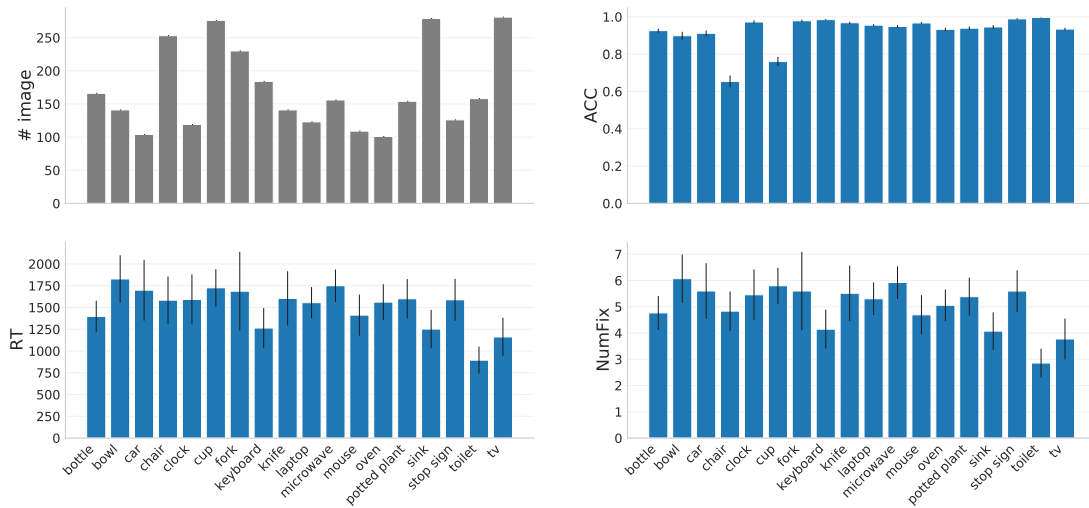


Figure S8. COCO-Search18 analyses for all 18 target categories in target-absent trials. Top: number of images in each category (gray), and response accuracy (ACC). Bottom: reaction time (RT) and number of fixations made before the button press (NumFix). Values are means over 10 participants, and error bars represent standard errors.

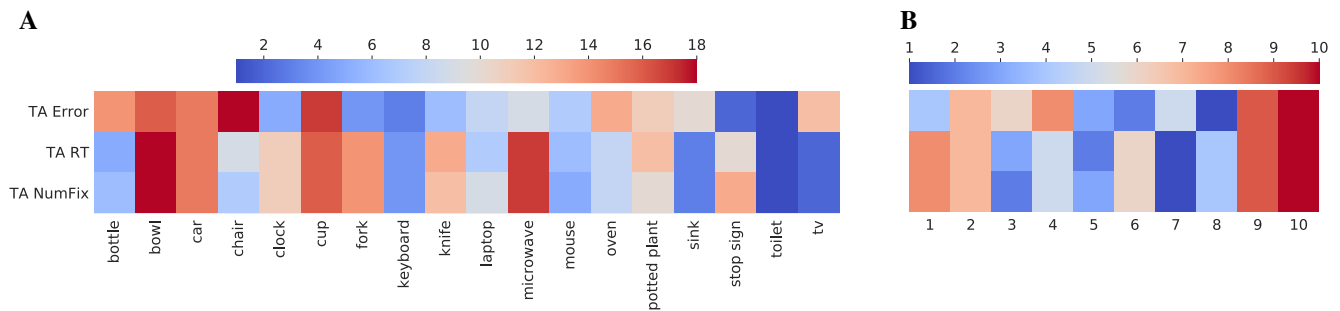


Figure S9. (A). Target-absent data, ranked [1-18] by target category (columns) and averaged over participants, shown for multiple performance measures (rows). These include: response error, reaction time (RT), and number of fixations (NumFix). Redder color indicates higher rank and harder search targets, bluer color indicates lower rank and easier search. (B) Target-absent data, now ranked by participant [1-10] and averaged over target category (columns). Performance measures and color coding are the same as in (A).

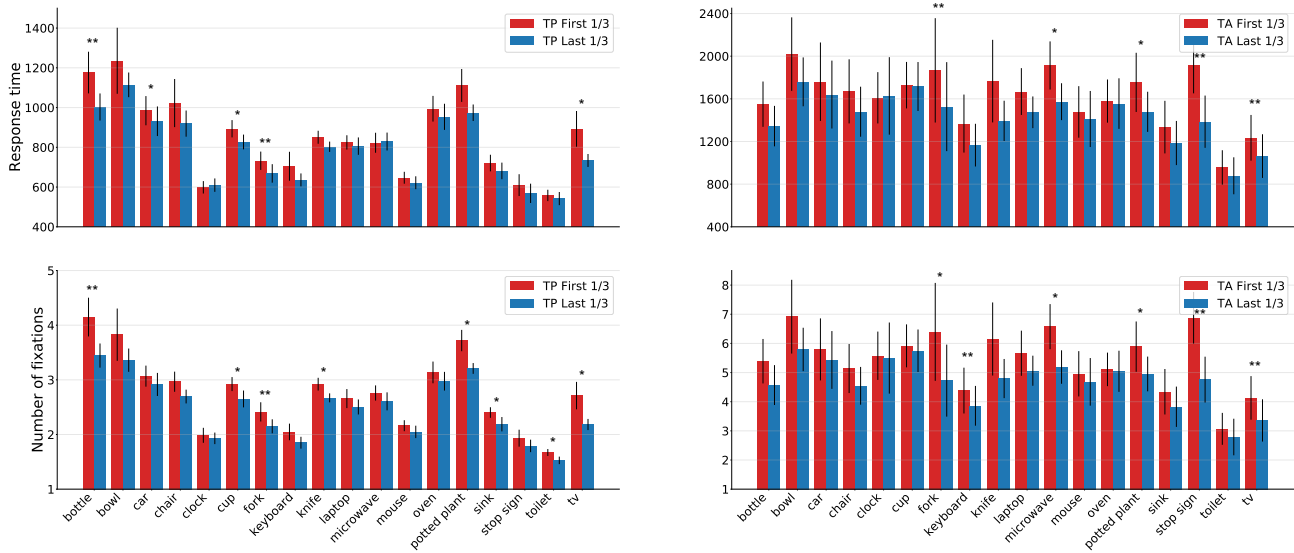


Figure S10. Practice effects, visualized as the difference in search performance between the red (first 1/3 of the trials) and the blue (last 1/3 of the trials) bars, grouped by the 18 target categories. The top row shows response time, and the bottom row shows the number of fixations before the button press. Target-present data are shown on the left, target-absent data are shown on the right. Only correct trials were included. *: $p < .05$, **: $p < .01$

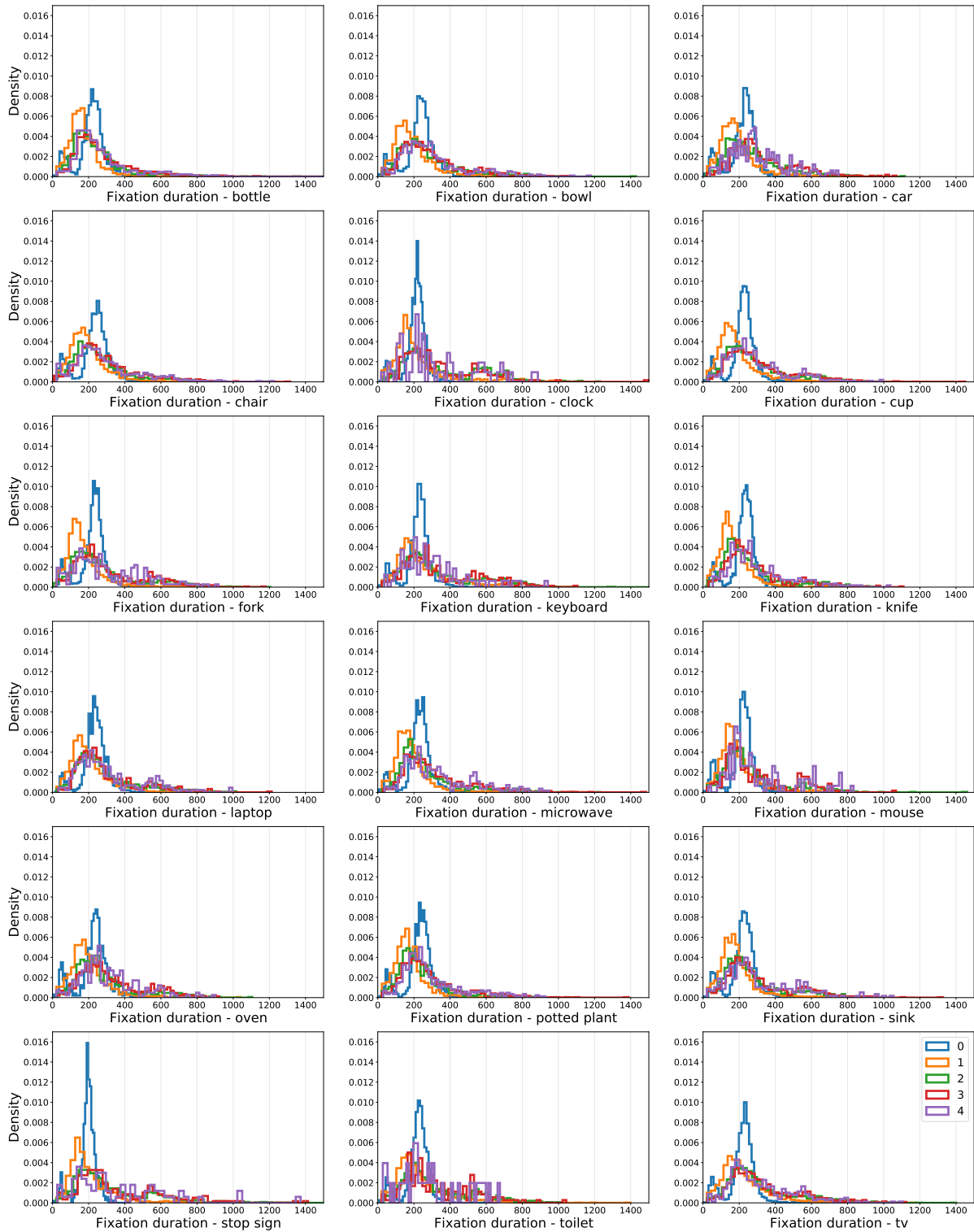


Figure S11. Density distributions of target-present fixation durations, plotted for each of the target categories (bin size = 50ms). The color lines refer to the initial fixation durations (0, blue), followed by the first four new fixations (1-4).

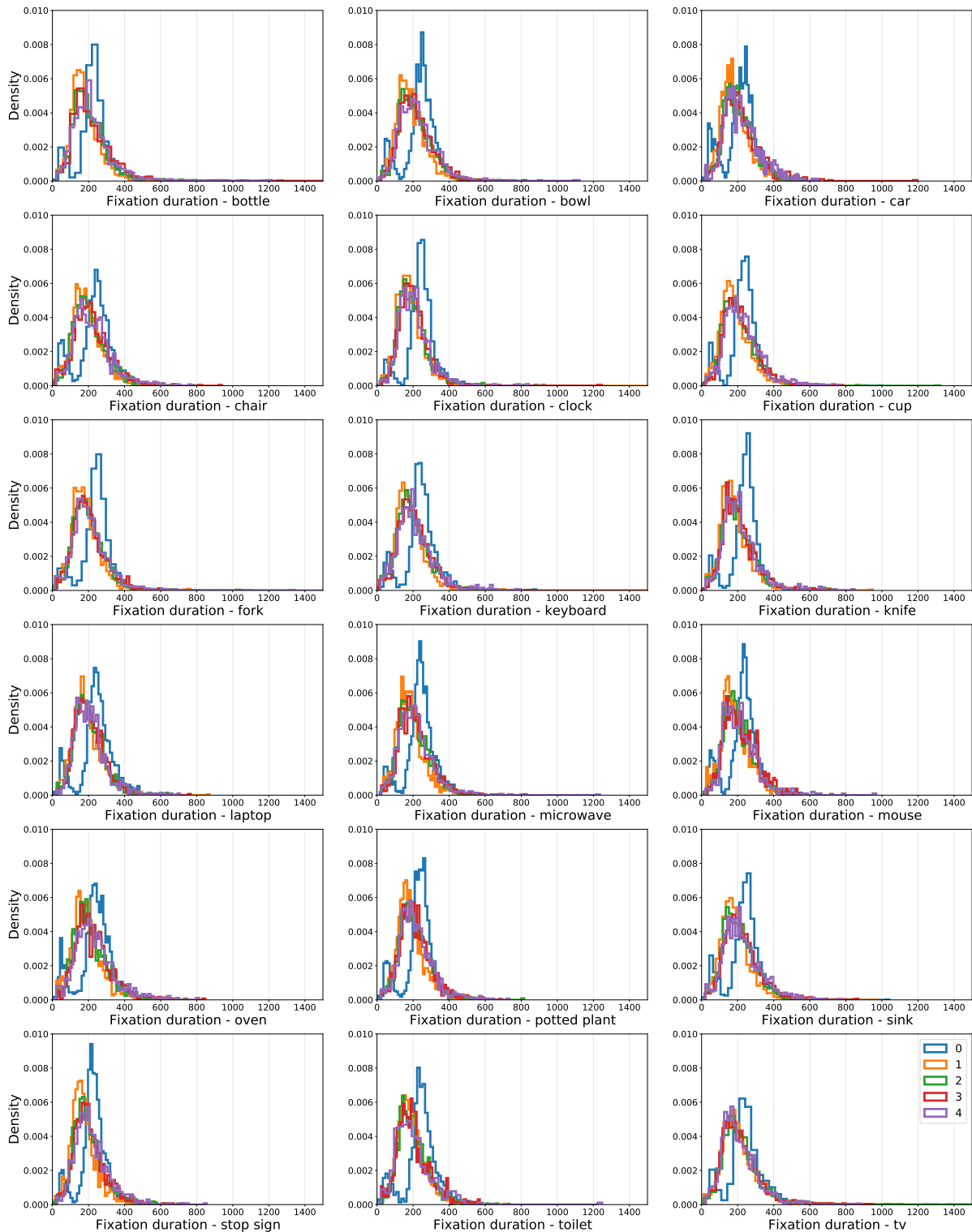


Figure S12. Density distributions of target-absent fixation durations, plotted for each of the target categories (bin size = 50ms). The color lines refer to the initial fixation durations (0, blue), followed by the first four new fixations (1-4).

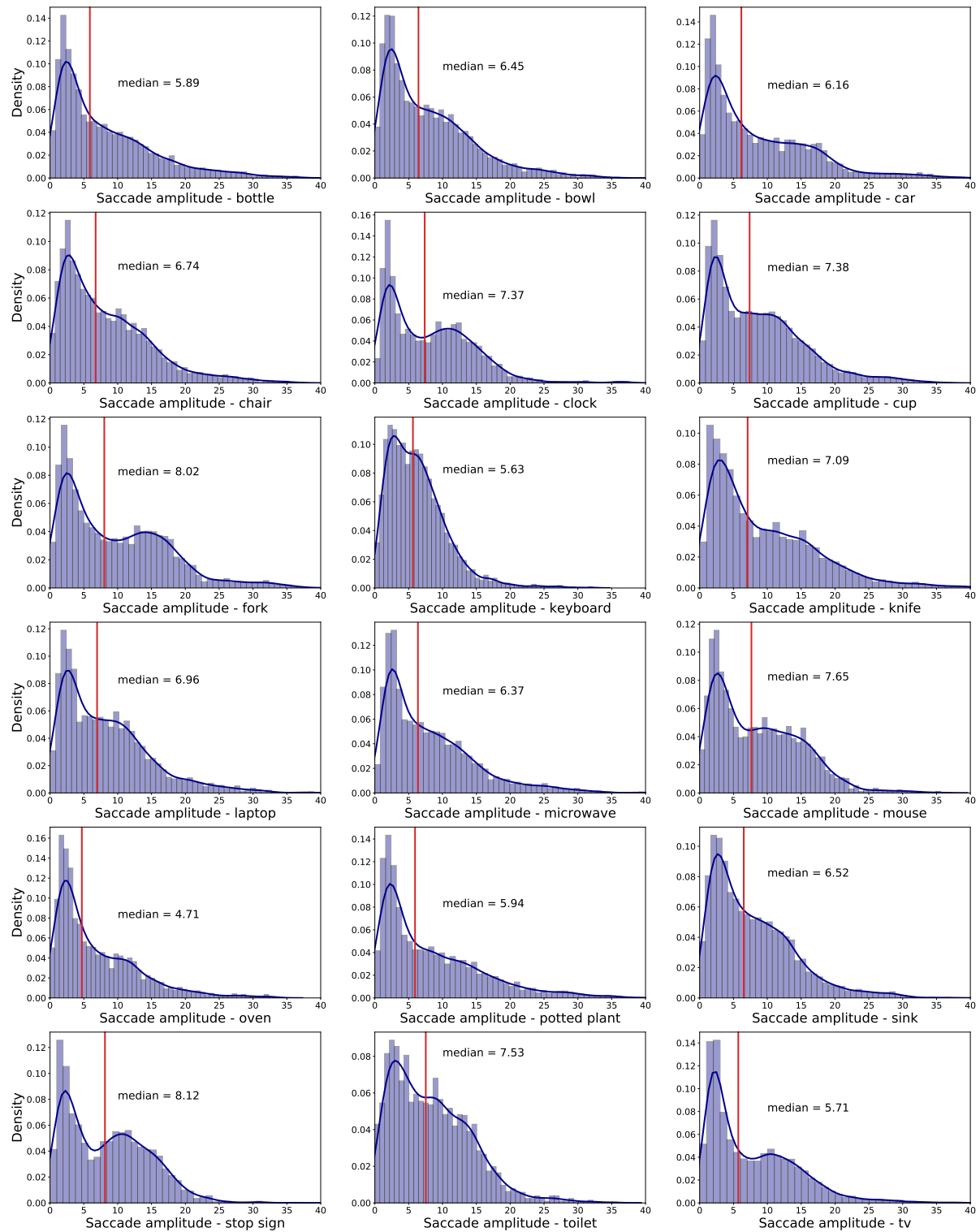


Figure S13. Density distributions of target-present saccade amplitudes (in visual angle), plotted by target category. Red vertical lines indicate median amplitudes. Dark blue lines represent Gaussian kernel density estimates.

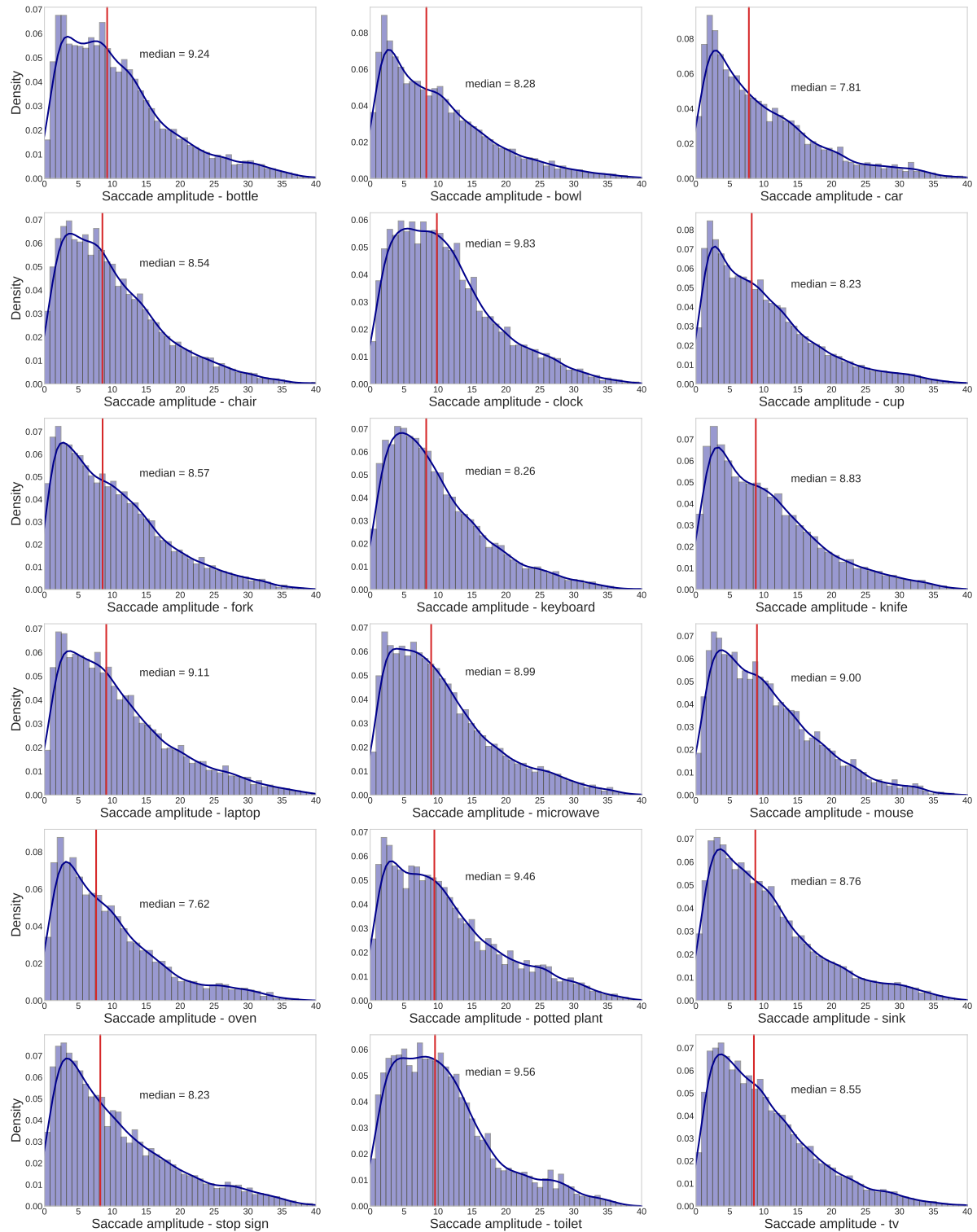


Figure S14. Density distributions of target-absent saccade amplitudes (in visual angle), plotted by target category. Red vertical lines indicate median amplitudes. Dark blue lines represent Gaussian kernel density estimates.

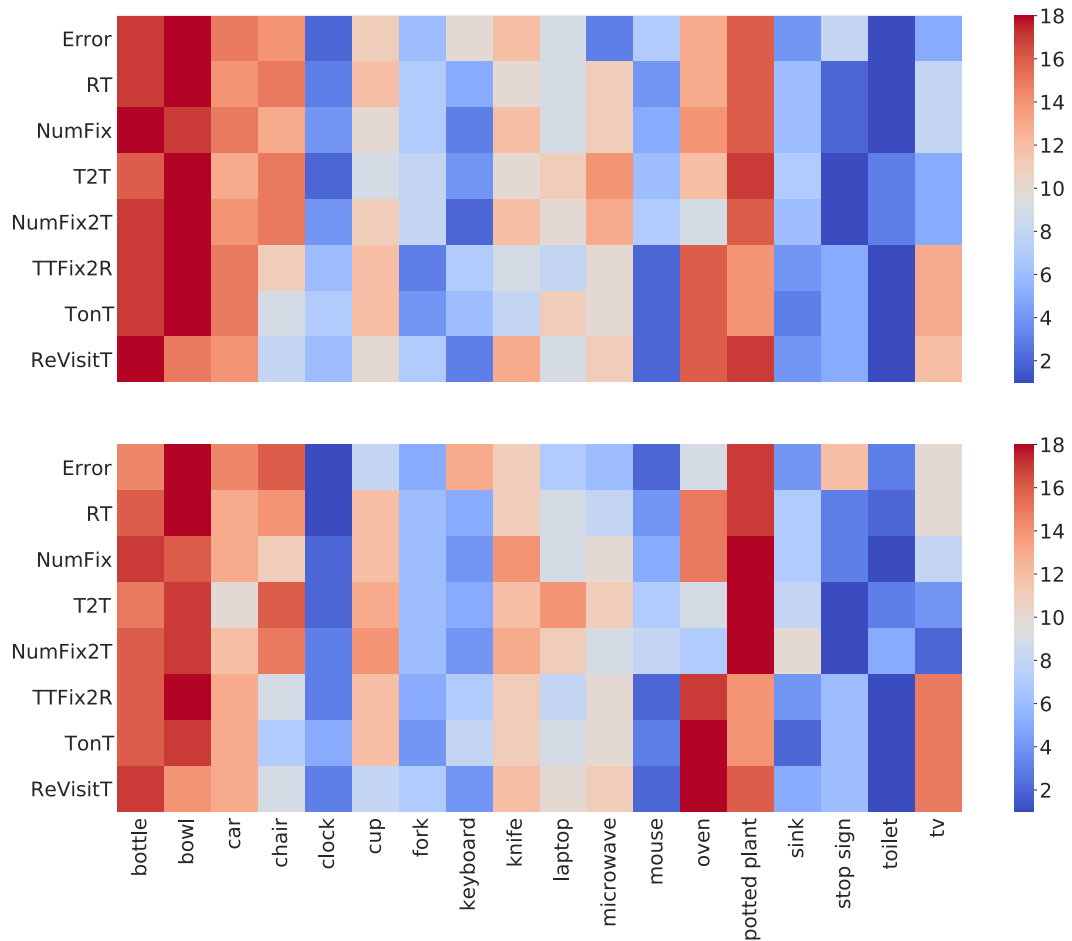


Figure S15. Target-present data, ranked by target category (1-18, columns) and shown for multiple performance measures (rows) in the trainval (top) and test (bottom) COCO-Search18 datasets. Redder color indicates higher rank and harder search targets, bluer color indicates lower rank and easier search. Measures include: response error, reaction time (RT), number of fixations (NumFix), time to target (T2T), number of fixations to target (NumFix2T), time from first target fixation until response (TTFix2R), time spent fixating the target (TonT), and the number of target re-fixations (ReVisitT).

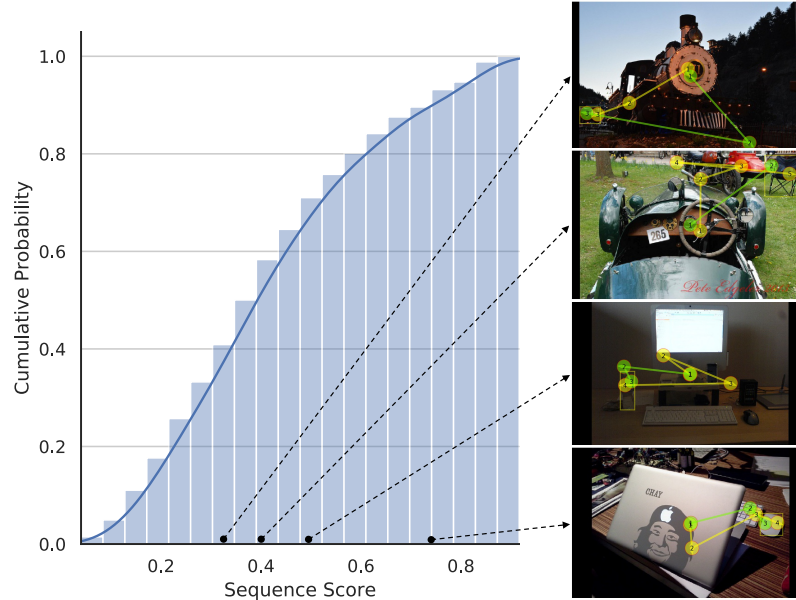


Figure S16. Left: cumulative distribution of average sequence scores computed between each scanpath generated by the IRL model and each behavioral scanpath for the test images of COCO-Search18. Right: Examples illustrating the scanpaths producing four different sequence scores. Behavioral scanpaths are colored in yellow, and the IRL-generated scanpaths are in green. Sequence scores for the four illustrated examples are 0.33, 0.40, 0.50, and 0.75, from top to bottom. Note that these results are from a slightly different version of the IRL model than the one reported here.

	AUC ↑	NSS ↑	CC ↑
Human	0.675	3.396	0.356
Random	0.531	0.280	0.039
Detector-Hi	0.605	1.210	0.163
Detector-Hi-Low	0.575	0.792	0.105
Deep Search-Hi	0.620	1.122	0.153
Deep Search-Hi-Low	0.598	0.864	0.118
IRL-ReT-C	0.595	1.601	0.214
IRL-Hi-Low-C	0.628	1.806	0.246
IRL-Hi-Low	0.621	1.728	0.235

Table S1. Results from models (rows) predicting behavioral fixation-density maps (FDMs) using three spatial comparison metrics (columns), applied to the COCO-Search18 test images. “Human” refers to an oracle method whereby the FDM from half of the searchers was used to predict the FDM from the other half of the searchers. See the supplemental text for additional details about the spatial fixation comparison metrics.

Compared Models	TFP-AUC	Probability Mismatch	Scanpath Ratio	Sequence Score	MultiMatch			
					shape	direction	length	position
IRL-ReT-C vs. IRL-Hi-Low-C	<i>n.s.</i>	<i>n.s.</i>	<i>n.s.</i>	<i>n.s.</i>	<i>n.s.</i>	<i>n.s.</i>	<i>n.s.</i>	<i>n.s.</i>
IRL-ReT-C vs. IRL-Hi-Low	<i>n.s.</i>	<i>n.s.</i>	<i>n.s.</i>	<i>n.s.</i>	<i>n.s.</i>	<i>n.s.</i>	<i>n.s.</i>	<i>n.s.</i>
IRL-ReT-C vs. Detector-Hi	<i>n.s.</i>	<i>n.s.</i>	<i>n.s.</i>	<i>n.s.</i>	<i>n.s.</i>	<i>n.s.</i>	<i>n.s.</i>	<i>n.s.</i>
IRL-ReT-C vs. Detector-Hi-Low	<i>.0017</i>	<i><.001</i>	<i><.001</i>	<i>n.s.</i>	<i>.005</i>	<i>.0686</i>	<i><.001</i>	<i>.0039</i>
IRL-ReT-C vs. Deep Search-Hi	<i><.001</i>	<i><.001</i>	<i><.001</i>	<i>n.s.</i>	<i>n.s.</i>	<i><.001</i>	<i>n.s.</i>	<i>n.s.</i>
IRL-ReT-C vs. Deep Search-Hi-Low	<i><.001</i>	<i><.001</i>	<i><.001</i>	<i>.0587</i>	<i>n.s.</i>	<i><.001</i>	<i>n.s.</i>	<i>n.s.</i>
IRL-Hi-Low-C vs. IRL-Hi-Low	<i>n.s.</i>	<i>n.s.</i>	<i>n.s.</i>	<i>n.s.</i>	<i>n.s.</i>	<i>n.s.</i>	<i>n.s.</i>	<i>n.s.</i>
IRL-Hi-Low-C vs. Detector-Hi	<i>n.s.</i>	<i>n.s.</i>	<i>.0653</i>	<i>n.s.</i>	<i>n.s.</i>	<i>n.s.</i>	<i>.0235</i>	<i>n.s.</i>
IRL-Hi-Low-C vs. Detector-Hi-Low	<i><.001</i>	<i><.001</i>	<i><.001</i>	<i>n.s.</i>	<i><.001</i>	<i>.0515</i>	<i><.001</i>	<i><.001</i>
IRL-Hi-Low-C vs. Deep Search-Hi	<i><.001</i>	<i><.001</i>	<i><.001</i>	<i>n.s.</i>	<i>n.s.</i>	<i><.001</i>	<i>n.s.</i>	<i>n.s.</i>
IRL-Hi-Low-C vs. Deep Search-Hi-Low	<i><.001</i>	<i><.001</i>	<i><.001</i>	<i>.0559</i>	<i>.0298</i>	<i><.001</i>	<i>n.s.</i>	<i>.0110</i>
IRL-Hi-Low vs. Detector-Hi	<i>n.s.</i>	<i>n.s.</i>	<i>.0151</i>	<i>n.s.</i>	<i>n.s.</i>	<i>n.s.</i>	<i>.0206</i>	<i>n.s.</i>
IRL-Hi-Low vs. Detector-Hi-Low	<i><.001</i>	<i><.001</i>	<i><.001</i>	<i>n.s.</i>	<i><.001</i>	<i>.0539</i>	<i><.001</i>	<i><.001</i>
IRL-Hi-Low vs. Deep Search-Hi	<i><.001</i>	<i><.001</i>	<i><.001</i>	<i>n.s.</i>	<i>n.s.</i>	<i><.001</i>	<i>n.s.</i>	<i>n.s.</i>
IRL-Hi-Low vs. Deep Search-Hi-Low	<i><.001</i>	<i><.001</i>	<i><.001</i>	<i>.0506</i>	<i>n.s.</i>	<i><.001</i>	<i>n.s.</i>	<i>.0029</i>
Detector-Hi vs. Detector-Hi-Low	<i>.0019</i>	<i><.001</i>	<i>.0086</i>	<i>n.s.</i>	<i>n.s.</i>	<i>n.s.</i>	<i>n.s.</i>	<i>.0150</i>
Detector-Hi vs. Deep Search-Hi	<i><.001</i>	<i><.001</i>	<i><.001</i>	<i>n.s.</i>	<i>n.s.</i>	<i>.0013</i>	<i><.001</i>	<i>n.s.</i>
Detector-Hi vs. Deep Search-Hi-Low	<i><.001</i>	<i><.001</i>	<i><.001</i>	<i>.0755</i>	<i>n.s.</i>	<i><.001</i>	<i><.001</i>	<i>n.s.</i>
Detector-Hi-Low vs. Deep Search-Hi	<i>n.s.</i>	<i>n.s.</i>	<i>n.s.</i>	<i>n.s.</i>	<i><.001</i>	<i>n.s.</i>	<i><.001</i>	<i><.001</i>
Detector-Hi-Low vs. Deep Search-Hi-Low	<i>n.s.</i>	<i>.0275</i>	<i>n.s.</i>	<i>n.s.</i>	<i>.0446</i>	<i>n.s.</i>	<i><.001</i>	<i>.0511</i>
Deep Search-Hi vs. Deep Search-Hi-Low	<i>n.s.</i>	<i>n.s.</i>	<i>n.s.</i>	<i>n.s.</i>	<i>n.s.</i>	<i>n.s.</i>	<i>n.s.</i>	<i>.0778</i>

Table S2. *P* values from post-hoc t-tests (Bonferroni corrected) comparing predictive models (rows), averaged across the 18 target categories, for multiple scanpath metrics (columns). All *dfs* = 34. For decisively significant comparisons, the more predictive model is indicated in boldface.