

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

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

Department of Psychology, Stony Brook University¹

Department of Computer Science, Stony Brook University²

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^{8,10,23}. In terms of number of fixations, number of target categories, and number of images, COCO-Search18 is far larger. The PET dataset¹⁰ collected search fixations for six animal target categories in 4,135 images selected from the Pascal VOC 2012 dataset⁹, 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 search¹⁶. Gaze fixations were also recorded for only 2 seconds/image, and multiple targets often appeared in each scene. The microwave-clock search dataset (MCS²³) 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 from⁸, 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 (SALICON¹³) or fixations collected using other visual tasks (POET¹⁹), 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 dataset¹⁵ 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^{4,14}. 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.

1017 Applying these exclusion criteria left 32 object categories
1018 from COCO’s original 80. Given that this left still far too
1019 many images for people to practically annotate with search
1020 fixations, we decided to attempt exclusion of images where
1021 targets were highly occluded or otherwise difficult to recog-
1022 nize. We did this out of concern that such images would
1023 largely introduce noise into the search behavior. To do this,
1024 we trained object detectors on cropped views of these 32 cat-
1025 egories, and excluded images if the object bounding boxes
1026 had a classification confidence $< .99$. Specifically, for these
1027 32 categories we created a validation set consisting of images
1028 meeting the selection criteria and a training set consisting of
1029 the images that did not. The bounding box of the object, for
1030 each of the 32 object classes, was then cropped in the image to
1031 obtain the positive training samples. Negative samples were
1032 same-sized image patches that had 25% intersection with the
1033 target (area of intersection divided by area of target), mean-
1034 ing that they were class-specific hard negatives. All cropped
1035 patches (over 1 million) were resized to 224×224 pixels while
1036 maintaining the aspect ratio using padding. The classifier was
1037 a ResNet50 pre-trained on ImageNet, which we fine-tuned
1038 by dilating the last fully-connected layer and re-training on
1039 33 outputs (32+“Negative”). Images were excluded if the
1040 cropped object patch had a classification score of less than
1041 $.99$. This procedure resulted in 18 categories with at least 100
1042 images in each category, totaling 3,131 TP images.

1043 Two final exclusion criteria were implemented by manual
1044 selection. First, for the clock target category we included only
1045 images of analog clocks, meaning that we excluded digital
1046 clocks from being clock targets. We did this because the fea-
1047 tures of analog and digital clocks are highly distinct and very
1048 different, and we were concerned that this would introduce
1049 variability in the search behavior and reduce data quality. Five
1050 images depicting only digital clocks were excluded for this
1051 reason. Lastly, images from all 18 of the target categories
1052 were screened for objectionable content, which we defined
1053 as offensive content or content evoking discomfort or disgust.
1054 The “toilet” category had the most images (17) excluded for
1055 objectionable content, with a total of 25 images excluded
1056 across all target categories. After implementing all exclusion
1057 criteria discussed in this section, we obtained 3,101 TP images
1058 from 18 categories: bottle, bowl, car, chair, (analog) clock,
1059 cup, fork, keyboard, knife, laptop, microwave, (computer)
1060 mouse, oven, potted plant, sink, stop sign, toilet, and tv. See
1061 Figure 2 for the specific number of images in each category.

1062 **Target-absent image selection.** To balance the selection
1063 of the 3,101 TP images, we selected an equal number of TA
1064 images from COCO. To do this, we kept the criteria excluding
1065 images depicting people or animals, extreme width/height
1066 image ratios, and images with objectionable content, all as
1067 described for the TP image selection, but added two more
1068 exclusion criteria that were specific to each of the 18 target-
1069 object categories.

1070 (1) Images were excluded if they depicted an instance of the

target, a prerequisite for a TA image. 1071

(2) Images were excluded if they depicted less than two 1072
instances of the target category’s siblings, a criterion 1073
introduced to discourage searchers from making TA res- 1074
ponses purely on the basis of scene type. For example, a 1075
person might be biased to make a TA response if they are 1076
searching for a toilet target and the image is a street scene. 1077
Because COCO has a hierarchical organization, parent, 1078
child, and sibling relationships can be used for image 1079
selection. For example, COCO defines the siblings of a 1080
microwave to be an oven, toaster, refrigerator, and sink, 1081
all under the parent category of appliance. By requiring 1082
that the TA scenes for a target category have at least two 1083
of that category’s siblings, we impose a sort of scene 1084
constraint that minimizes target-scene inconsistency and 1085
makes a scene appropriate to use as a TA image. A scene 1086
that has an oven and a refrigerator is very likely to be 1087
a kitchen, thereby making it difficult to answer on the 1088
basis of scene type alone whether a microwave target is 1089
present or absent. 1090

1091 These exclusion criteria still left us with many thousands
1092 more TA images than we needed, so we sampled randomly
1093 within each of the 18 target categories to match the 3,101 TP
1094 images.

1095 **Order of target-category presentation**

1096 Collecting the search behavior for 6,202 images required di-
1097 viding each participant’s effort into six days of testing. Each
1098 testing session was conducted on a different day, lasted about
1099 2 hours, and consisted of about 1000 search trials, evenly
1100 divided between TP and TA. Because images from different
1101 categories can overlap (e.g., images depicting a microwave
1102 may also depict an oven), the presentation order of the target-
1103 category blocks was constrained to minimize the repetition
1104 of images in consecutive categories and consecutive sessions.
1105 For example, because 49 images satisfied the selection criteria
1106 for both the sink and microwave target categories, we pre-
1107 vented the microwave and sink categories from appearing in,
1108 not only the same session, but the sessions preceding and fol-
1109 lowing. We did this to minimize possible biases resulting from
1110 seeing the same scene in different search contexts. A heuristic
1111 for maximizing this distance between repeating images
1112 resulted in the following fixed target category presentation
1113 order across the six sessions:

- 1114 (1) tv + sink;
- 1115 (2) fork + chair;
- 1116 (3) car + bowl + potted plant + mouse;
- 1117 (4) knife + keyboard + oven + clock;
- 1118 (5) cup + laptop + toilet;
- 1119 (6) bottle + stop sign + microwave.

1120 Each participant viewed from Session 1 to Session 6, or
1121 from Session 6 to Session 1, with this order counterbalanced
1122 across participants.

1123 **Data-collection procedure**

1124 Participants were 10 Stony Brook University undergraduate
1125 and graduate students, 6 males and 4 females, with ages rang-
1126 ing from 18–30 years. All had normal or corrected to normal
1127 vision, by self report, were naive with respect to task design
1128 and paradigm when recruited, and were compensated with
1129 course credit or money for their participation. Informed con-
1130 sent was obtained from each participant at the beginning of
1131 testing, in accordance with the Institutional Review Board
1132 responsible for overseeing human-subjects research at Stony
1133 Brook University.

1134 The target category was designated to participants at the
1135 start of each block. This was done using the type of display
1136 shown in Figure S2 for the potted-plant and analog clock
1137 categories. The name of the target category was shown in
1138 text at the top, with examples of objects that would, or would
1139 not, qualify as exemplars of the named category. In selecting
1140 exemplars to illustrate as positive target-category members,
1141 we attempted to capture key categorical distinctions at a level
1142 immediately subordinate to the target category. When needed,
1143 we also gave negative examples by placing a red X through
1144 the object. We did this to minimize potential confusions and
1145 to enable the participant to better define the target category’s
1146 boundary.

1147 The procedure (Figure S3) on each trial began with a fixa-
1148 tion dot appearing at the center of the screen. To start a trial,
1149 the participant would press the “X” button on a game-pad con-
1150 troller while carefully looking at the fixation dot. An image
1151 of a scene would then be displayed and the participant’s task
1152 would be to answer, “yes” or “no”, whether an exemplar of the
1153 target category appears in the displayed scene by pressing the
1154 right or left triggers of the game-pad, respectively. The search
1155 scene remained visible until the manual response. Participants
1156 were told that there were an equal number of TP and TA trials,
1157 and that they should make their responses as fast as possible
1158 while maintaining high accuracy. No accuracy or response
1159 time feedback was provided.

1160 The presentation of images during the experiment was con-
1161 trolled by Experiment Builder (SR research Ltd., Ottawa,
1162 Ontario, Canada). Stimuli were presented to participants on
1163 a 22-inch LCD monitor (1680×1050 pixel resolution) at a
1164 viewing distance of 47cm from the monitor, enforced by chin
1165 and head rests. These viewing conditions resulting in hori-
1166 zontal and vertical visual angles of $54^\circ \times 35^\circ$, respectively.
1167 Participants were asked to keep their gaze on the fixation point
1168 at the start of each trial, but were told that they should feel free
1169 to move their eyes as they searched. Eye movements were
1170 recorded throughout the experiment using an EyeLink 1000
1171 eye-tracker in tower-mount configuration (SR research Ltd.,
1172 Ottawa, Ontario, Canada). Eye-tracker calibrations occurred
1173 before every block or whenever necessary, and these 9-point
1174 calibrations were not accepted unless the average calibration
1175 error was $\leq .51^\circ$ and the maximal error was $\leq .94^\circ$. The ex-
1176 periment was conducted in a quiet laboratory room under dim
1177 lighting conditions.

1178 **SM2: Behavioral evaluation of COCO-Search18**

1179 **Effects of set size and target eccentricity**

1180 The visual search literature has done excellent work in identi-
1181 fying many of the factors that increase search difficulty (for
1182 reviews, see:^{6,7,21,22}). Larger set sizes (number of items in
1183 the search display), smaller target size, larger target eccen-
1184 tricity, and greater target-distractor similarity are all known to
1185 make search more difficult. However, most of this work was
1186 done in the context of simple stimuli, and generalization to
1187 realistic images is challenging. For example, what to consider
1188 an object in a scene is often unclear, making it difficult to de-
1189 fine a set size¹⁸. Objects in images also do not usually come
1190 annotated with labels and bounding boxes. These problems of
1191 object segmentation and identification, which largely do not
1192 exist for search studies using object arrays, become significant
1193 obstacles to research when scaled up to images of scenes.

1194 With COCO-Search18, we can begin to ask how the search
1195 for targets in images is affected by set size and target eccen-
1196 tricity. Set size is determined based on the COCO object and
1197 stuff labels, which collectively map every pixel in an image
1198 to an object or stuff category. Set size is the count of the
1199 number of these labels for a given image. Figure S4 shows
1200 the relationship between the number of fixations made on an
1201 image, averaged over participants, and the set size of that im-
1202 age, grouped by target category. Some target categories, such
1203 as laptop, oven, microwave, and potted-plant, have significant
1204 positive set size effects ($r = .21$ to $.37$, $ps \leq .01$), indicating
1205 a less efficient search with more objects. A similar pattern is
1206 shown in Figure S5 for the relationship between the number of
1207 fixations on a search image and the initial visual eccentricity
1208 of the target (distance between the image center and the target
1209 bounding-box center), where for these same objects there was
1210 a decrease in search efficiency with increasing target eccen-
1211 tricity. For other target object categories, such as: stop sign,
1212 fork, and keyboard, search efficiency was unaffected by either
1213 set size or target eccentricity ($ps > .05$), possibly because
1214 these objects are either highly salient (stop sign) or highly
1215 constrained by scene context (keyboard).

1216 **Distance between search fixations and the target**

1217 How much closer does each search fixation bring gaze to
1218 the target? We analyzed this measure of search efficiency
1219 and report the results in Figure S6. Plotted is the Euclidean
1220 distance between the target location and the locations of the
1221 starting fixation (0) and the fixation locations after the first six
1222 eye movements (1-6). The most salient pattern is the rapid
1223 decrease in fixation-target distance in the first two new fix-
1224 ations, which dovetails perfectly with the steep increase in
1225 the cumulative probability of target fixation over these same
1226 eye movements reported in Figure 4A. From a starting lo-
1227 cation near the center of the image, these eye movements
1228 brought gaze steadily closer to the target. Note that because
1229 this fixation-target distance is averaged over images and partic-
1230 ipants, the roughly 5 degrees of visual angle at the bottom of
1231 these functions should not be misinterpreted as gaze being this
1232 distance from the target on a given trial. More interpretable

1233 are the overall trends, where a steep drop in distance is fol- 1287
1234 lowed by a plateau, or even a smaller increase in distance with 1288
1235 the 5th and 6th new fixations. This small increase is likely an 1289
1236 artifact of these 5 and 6-fixation trials being the most difficult, 1290
1237 with more idiosyncratic search behavior. 1291

1238 **Target-absent search fixations**

1239 In the main text we focused on the TP data, where the guid- 1292
1240 ance signal is clearer and the modeling goals are better defined, 1293
1241 but we conducted largely parallel analyses of the TA data. Fig- 1294
1242 ure S7A shows representative TA images with fixation data 1295
1243 from one participant, and Figure S7B shows FDMs from all 1296
1244 participants for the same images. Comparing these data with 1297
1245 the TP data from Figure 1, it is clear that people made many 1298
1246 more fixations in the absence of a target. This was expected 1299
1247 from the search literature, but it should also be noted that the 1300
1248 FDMs are still much sparser than what would be hypothesized 1301
1249 by an exhaustive search. Paralleling Figure 3, in Figure S8 we 1302
1250 report applicable analyses of the TA search behavior. These 1303
1251 are grouped by manual accuracy and response time, and the 1304
1252 mean number of fixations made before the target-absent but- 1305
1253 ton press terminating a trial. Note that accuracy was high 1306
1254 (low false positive error rate) for all of the target categories 1307
1255 except chairs and cups, with the reason for the former already 1308
1256 discussed in the context of mislabeling and the reason for the 1309
1257 latter likely reflecting an occasionally challenging category 1310
1258 distinction (e.g., some bottles can look like some cups). Also 1311
1259 note that there was an average of only five fixations made 1312
1260 during search, even on the TA search trials. As in Figure 5, 1313
1261 Figure S9 visualizes the agreement and other patterns among 1314
1262 these measures. The rows show ranked performance, with 1315
1263 dark red indicating more difficult (or least efficient) search 1316
1264 and dark blue indicating relatively easy or efficient search. 1317
1265 The columns in Figure S9A group the measures by target 1318
1266 category. Similar to the TP data, there was again good con- 1319
1267 sistency among the measures. Also consistent is the fact that 1320
1268 bottles and cups were among the most difficult target cate- 1321
1269 gories, whereas the toilet category was the easiest. There was 1322
1270 also evidence in the TA data for a speed-accuracy trade-off 1323
1271 for some target categories. For example, microwaves and stop 1324
1272 signs had relatively low error rates, but these categories were 1325
1273 searched with relatively high effort, as measured by ranked 1326
1274 response time and number of fixations. Figure S9B visualizes 1327
1275 the measures by participant instead of category, where we 1328
1276 again found individual differences between participants in 1329
1277 search efficiency. 1330

1278 **Practice effects**

1279 Each of the participants contributing to COCO-Search18 1331
1280 searched more than 6000 images, making it possible to ana- 1332
1281 lyze how their search efficiency improved with practice. Fig- 1333
1282 ure S10 shows practice effects for both response time (top) 1334
1283 and the number of fixations before the button press (bottom), 1335
1284 where we define practice effects as performance on the first 1336
1285 1/3 of the trials compared to performance on the last 1/3 of 1337
1286 the trials for each target category. Practice effects were larger for 1338

TA trials (right) than for TP trials (left), noting the differences 1287
in y-axis scales, and that considerable differences existed 1288
across categories. Some categories, such as bottles, showed 1289
large practice effects, while other categories, such as analog 1290
clocks, showed none at all. We speculate that this difference is 1291
due to some categories requiring more exemplars to fully learn 1292
compared to others. For example, analog clock was perhaps 1293
the most well defined of COCO-Search18’s categories, and 1294
bottle certainly one of the least well defined, creating greater 1295
opportunity to better learn the bottle category with practice 1296
over trials. 1297

1298 **Search fixation durations**

1299 Figures S11 and S12 show density histograms of the search 1299
fixation durations for the TP and TA data, respectively, plot- 1300
ted for each of the target categories. Fixation durations are 1301
plotted across the x-axes with a bin size of 50ms, and y-axes 1302
show the normalized probability density at each fixation. Of 1303
note in the TP data is that the mode initial fixation durations 1304
(blue lines) were a bit longer than the mode duration of the 1305
rest (averaged mode difference = 63ms), consistent with the 1306
very strong guidance observed in the initial eye movements, 1307
and they tended to have more bi-modal distributions. The 1308
main peak was at ~250 ms, with a smaller and very short- 1309
latency peak at ~50 ms that is likely a truncation artifact of 1310
fixation duration being measured relative to the onset of the 1311
search display. In contrast, the distributions of second fixa- 1312
tions (orange lines) were consistently shorter, even relative to 1313
the subsequent fixations. Speculatively, this may be due to 1314
a greater proportion of the first new fixations being “off ob- 1315
ject”²⁴, which are often followed by short-latency corrective 1316
saccades that bring gaze accurately to an object. This inter- 1317
pretation is consistent with the high probability of the target 1318
being fixated by the second eye movement (Figure 4A). As 1319
for the subsequent fixations, they tended to be short (~200ms) 1320
and not highly variable in their durations. The TA fixations 1321
showed similar trends, except for the durations of the second 1322
fixations no longer differing from the rest. 1323

1324 **Saccade amplitudes**

1325 We also analyzed the distribution of saccade amplitudes dur- 1325
ing visual search, defined here as the Euclidean distance be- 1326
tween consecutive fixations in visual angle. Figure S13 and 1327
Figure S14 show the distributions of saccade amplitudes in 1328
the TP and TA data, respectively. In the TP data, saccade 1329
amplitudes were larger in some categories (toilet and stop 1330
sign) than others (bottle and potted plant), likely because eas- 1331
ier target categories could be identified from farther in the 1332
visual periphery. There was also evidence for bimodality in 1333
the amplitude distributions, shown most clearly for clocks, 1334
forks, stop signs, and tvs. We speculate that this bimodal- 1335
ity reflects larger-amplitude exploratory saccades mixed with 1336
smaller-amplitude saccades used in the verification of an ob- 1337
ject category. Mean saccade amplitudes in the TA data were 1338
clearly larger than for the TP data ($t(17) = 11.79, p < .001$), 1339
and this difference was consistent across target categories (all 1340

1341 $ps \leq .001$). We attribute this to the relatively large viewing
1342 angle of the search displays (54×35 degrees of visual angle)
1343 creating a greater need for exploration, but this is also specula-
1344 tion. The distributions of saccade amplitudes were also more
1345 consistent across categories in the TA data, with there being
1346 weaker evidence of bi-modality.

1347 **SM3: Model Methods**

1348 *Training and testing datasets*

1349 Model success depends on the training dataset being an accu-
1350 rate reflection of the test dataset. When the training dataset
1351 includes a behavioral annotation, as does COCO-Search18, it
1352 is therefore important to know that similar patterns exist in
1353 the training and testing search behavior. The analyses shown
1354 in Figure 5A included images from all of COCO-Search18,
1355 which recall were randomly split into 70% for training, 10%
1356 for validation, and 20% for testing. Figure S15 replots the
1357 data from Figure 5A, but divides it into the training/validation
1358 (left) and testing (right) datasets. Note the high agreement
1359 between the testing and train/val datasets across this battery
1360 of behavioral performance measures.

1361 *Inverse Reinforcement Learning*

1362 The specific inverse-reinforcement learning (IRL) method
1363 that we used was generative adversarial imitation learning
1364 (GAIL¹²) with proximal policy optimization (PPO)²⁰. The
1365 model policy is a generator that aims to create state-action
1366 pairs that are similar to human behavior. The reward function
1367 (the logarithm of the discriminator output) maps a state-action
1368 pair to a numeric value. The generator and discriminator are
1369 trained within an adversarial optimization framework to obtain
1370 the policy and reward functions. The discriminator’s task is
1371 to distinguish whether a state-action pair was generated by
1372 a person (real) or by the generator (fake), with the generator
1373 aiming to fool the discriminator by maximizing the similarity
1374 between its state-action pairs and those from people. The
1375 reward function and policy that are learned from the fixation-
1376 annotated images during training are then used to predict new
1377 search fixations in the unseen test images.

1378 **SM4: Performance metrics and model evaluation**

1379 *Metrics for comparing search efficiency and scanpaths*

1380 We considered five metrics for quantifying search efficiency
1381 and comparing search scanpaths (Table 1). Two metrics for
1382 quantifying search efficiency follow directly from the group
1383 target-fixation probability (TFP) function shown in Figure 4.
1384 The first of these computes the area under the TFP curve, a
1385 metric we refer to as TFP-auc. Second, and relatedly, we
1386 compute the sum of the absolute differences between the hu-
1387 man and model target-fixation-probabilities in a metric that
1388 we refer to as Probability Mismatch. A third metric for quan-
1389 tifying overt search efficiency is Scanpath Ratio. It is the
1390 Euclidean distance between the initial fixation location and
1391 the target divided by the summed Euclidean distances between
1392 the fixation locations in the search scanpath¹¹. It is an effi-
1393 ciency metric because an initial saccade that lands directly

on the target would give a Scanpath Ratio of 1, meaning that
the distance between starting fixation and the target would
be the same as the summed saccade distance. These three
metrics emphasize target-fixation efficiency by penalizing ei-
ther the number of fixations or the saccade-distance traveled
to achieve the target goal. The final two metrics focus on
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 algorithm¹⁷ 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^{1,5} to measure the scanpath similarity at the pixel
level. MultiMatch measures five aspects of scanpath simi-
larity: 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 S3 for the results of statistical
tests comparing predictions from each pair of models.

1415 *Comparing predicted and behavioral fixation-density 1416 maps (FDMs)*

1417 Search has a temporal dynamic, making a metric for capturing
1418 the spatio-temporal sequence of fixations preferred over ones
1419 that compare only FDMs, where this temporal component is
1420 disregarded. However, the prediction of FDMs is common
1421 for free-viewing tasks, and because there is no technical rea-
1422 son why FDM metrics cannot be applied to search we do so
1423 here in the hope that the visual saliency literature finds this
1424 comparison useful. Models generated scanpaths having a max-
1425 imum length of 6 new fixations, but FDMs were constructed
1426 only from those fixations leading up to the first fixation on
1427 the target, just as FDMs were constructed from the behav-
1428 ioral fixations. We used three widely accepted metrics for
1429 comparing predicted against observed FDMs. Area Under
1430 the Receiver Operating Characteristic Curve (AUC) uses a
1431 predicted priority map as a binary classifier to discriminate
1432 behavioral fixation locations from non-fixated locations. Nor-
1433 malized Scanpath Saliency (NSS) finds the model predictions
1434 at each of the behavioral fixation locations, then averages and
1435 normalizes these values. Lastly we computed a Pearson’s
1436 Correlation Coefficient (CC) between the predicted and be-
1437 havioral FDMs, although this metric reflects only the degree
1438 of linear relationship between predicted and behavioral FDMs
1439 (for additional discussion, see: Borji & Itti²; Bylinskii et al.³).
1440 Table S2 reports the results of an evaluation comparing model
1441 predictions of search FDMs to behavioral search FDMs using
1442 each of these metrics. The findings that we report in the main
1443 text in the context of scanpath prediction also hold in the case
1444 of FDM prediction. Specifically, the IRL-Hi-Low-C model
1445 outperformed the others, and did so for all three metrics. Ad-
1446 ditionally, the Detector-Hi model also performed relatively
1447 well in all the metrics, supporting our conclusion that a simple
1448 detector does a relatively good job in predicting fixations in

1449 visual search.

1450 References

- 1451 1. Nicola C Anderson, Fraser Anderson, Alan Kingstone,
1452 and Walter F Bischof. A comparison of scanpath compar-
1453 ison methods. *Behavior research methods*, 47(4):1377–
1454 1392, 2015.
- 1455 2. Ali Borji and Laurent Itti. State-of-the-art in visual atten-
1456 tion modeling. *PAMI*, 35(1):185–207, 2012.
- 1457 3. Zoya Bylinskii, Tilke Judd, Aude Oliva, Antonio Torralba,
1458 and Frédo Durand. What do different evaluation metrics
1459 tell us about saliency models? *IEEE transactions on*
1460 *pattern analysis and machine intelligence*, 41(3):740–
1461 757, 2018.
- 1462 4. Moran Cerf, Jonathan Harel, Wolfgang Einhäuser, and
1463 Christof Koch. Predicting human gaze using low-level
1464 saliency combined with face detection. In *Advances in*
1465 *neural information processing systems*, pages 241–248,
1466 2008.
- 1467 5. Richard Dewhurst, Marcus Nyström, Halszka Jaro-
1468 dzka, Tom Foulsham, Roger Johansson, and Kenneth
1469 Holmqvist. It depends on how you look at it: Scan-
1470 path comparison in multiple dimensions with multimatch,
1471 a vector-based approach. *Behavior research methods*,
1472 44(4):1079–1100, 2012.
- 1473 6. John Duncan and Glyn W Humphreys. Visual search
1474 and stimulus similarity. *Psychological review*, 96(3):433,
1475 1989.
- 1476 7. Miguel P Eckstein. Visual search: A retrospective. *Jour-
1477 nal of vision*, 11(5):14–14, 2011.
- 1478 8. Krista A Ehinger, Barbara Hidalgo-Sotelo, Antonio Tor-
1479 ralba, and Aude Oliva. Modelling search for people in
1480 900 scenes: A combined source model of eye guidance.
1481 *Visual cognition*, 17(6-7):945–978, 2009.
- 1482 9. M. Everingham, L. Van Gool, C. K. I. Williams,
1483 J. Winn, and A. Zisserman. The PASCAL Visual
1484 Object Classes Challenge 2012 (VOC2012) Results.
1485 <http://host.robots.ox.ac.uk/pascal/VOC/index.html>.
- 1486 10. Syed Omer Gilani, Ramanathan Subramanian, Yan Yan,
1487 David Melcher, Nicu Sebe, and Stefan Winkler. Pet:
1488 An eye-tracking dataset for animal-centric pascal object
1489 classes. In *2015 IEEE International Conference on Mul-
1490 timedia and Expo (ICME)*, pages 1–6. IEEE, 2015.
- 1491 11. John M Henderson, Phillip A Weeks Jr, and Andrew
1492 Hollingworth. The effects of semantic consistency on
1493 eye movements during complex scene viewing. *Jour-
1494 nal of experimental psychology: Human perception and*
1495 *performance*, 25(1):210, 1999.
- 1496 12. Jonathan Ho and Stefano Ermon. Generative adversarial
1497 imitation learning. In *Advances in Neural Information*
1498 *Processing Systems*, pages 4565–4573, 2016.
- 1499 13. Ming Jiang, Shengsheng Huang, Juanyong Duan, and
1500 Qi Zhao. Salicon: Saliency in context. In *The IEEE*
1501 *Conference on Computer Vision and Pattern Recognition*
1502 *(CVPR)*, June 2015.
- 1503 14. Tilke Judd, Krista Ehinger, Frédo Durand, and Antonio
1504 Torralba. Learning to predict where humans look. In
1505 *2009 IEEE 12th international conference on computer*
1506 *vision*, pages 2106–2113. IEEE, 2009.
- 1507 15. Tsung-Yi Lin, Michael Maire, Serge Belongie, James
1508 Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and
1509 C Lawrence Zitnick. Microsoft coco: Common objects
1510 in context. In *European conference on computer vision*,
1511 pages 740–755. Springer, 2014.
- 1512 16. Justin T Maxfield and Gregory J Zelinsky. Searching
1513 through the hierarchy: How level of target categorization
1514 affects visual search. *Visual Cognition*, 20(10):1153–
1515 1163, 2012.
- 1516 17. Saul B Needleman and Christian D Wunsch. A general
1517 method applicable to the search for similarities in the
1518 amino acid sequence of two proteins. *Journal of molecu-
1519 lar biology*, 48(3):443–453, 1970.
- 1520 18. Mark B Neider and Gregory J Zelinsky. Exploring set
1521 size effects in scenes: Identifying the objects of search.
1522 *Visual Cognition*, 16(1):1–10, 2008.
- 1523 19. Dim P Papadopoulos, Alasdair DF Clarke, Frank Keller,
1524 and Vittorio Ferrari. Training object class detectors from
1525 eye tracking data. In *European conference on computer*
1526 *vision*, pages 361–376. Springer, 2014.
- 1527 20. John Schulman, Filip Wolski, Prafulla Dhariwal, Alec
1528 Radford, and Oleg Klimov. Proximal policy optimization
1529 algorithms. *arXiv:1707.06347*, 2017.
- 1530 21. Anne Treisman and Stephen Gormican. Feature analy-
1531 sis in early vision: evidence from search asymmetries.
1532 *Psychological Review*, 95(1):15, 1988.
- 1533 22. Gregory Zelinsky. A theory of eye movements during tar-
1534 get acquisition. *Psychological review*, 115(4):787, 2008.
- 1535 23. Gregory Zelinsky, Zhibo Yang, Lihan Huang, Yupei Chen,
1536 Seoyoung Ahn, Zijun Wei, Hossein Adeli, Dimitris Sama-
1537 ras, and Minh Hoai. Benchmarking gaze prediction for
1538 categorical visual search. In *Proceedings of the IEEE*
1539 *Conference on Computer Vision and Pattern Recognition*
1540 *Workshops*, pages 0–0, 2019.
- 1541 24. Gregory J Zelinsky. Tam: Explaining off-object fixations
1542 and central fixation tendencies as effects of population
1543 averaging during search. *Visual Cognition*, 20(4-5):515–
1544 545, 2012.

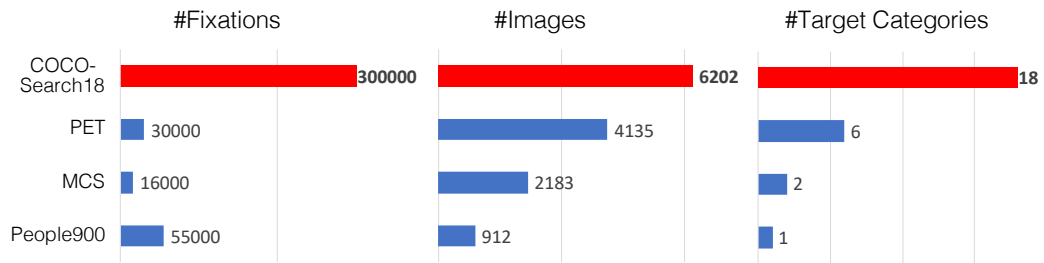


Figure S1. Comparisons between COCO-Search18 and other large-scale datasets of search behavior. COCO-Search18 is the largest in terms of number of fixations (~300,000), number of target categories (18), and number of images (6,202).

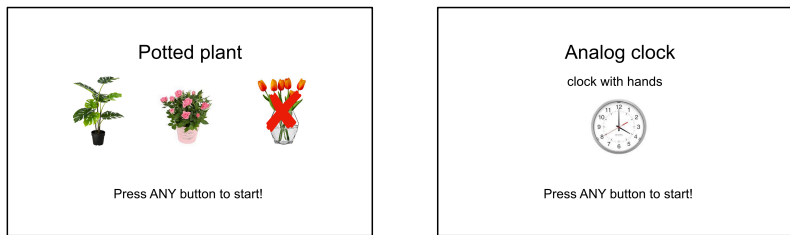


Figure S2. Examples of target-designation displays, shown for the potted-plant and analog clock targets, that preceded the block of trials for a given target category.

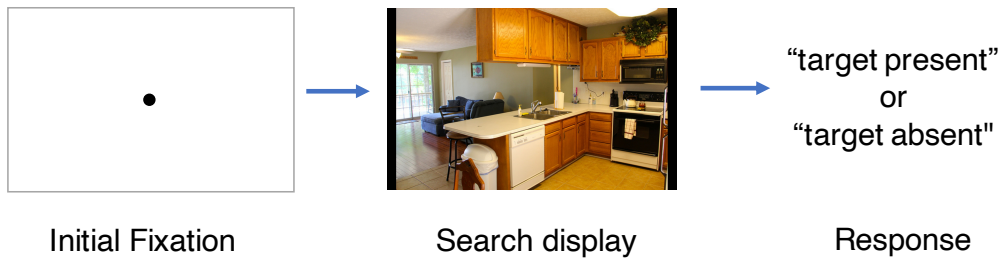


Figure S3. Example of the search procedure. Each trial began with a fixation dot appearing at the center of the screen. Participants would start a trial by pressing a button on a game-pad controller while carefully looking at the fixation dot. An image of a scene would then be displayed and the participant’s 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.

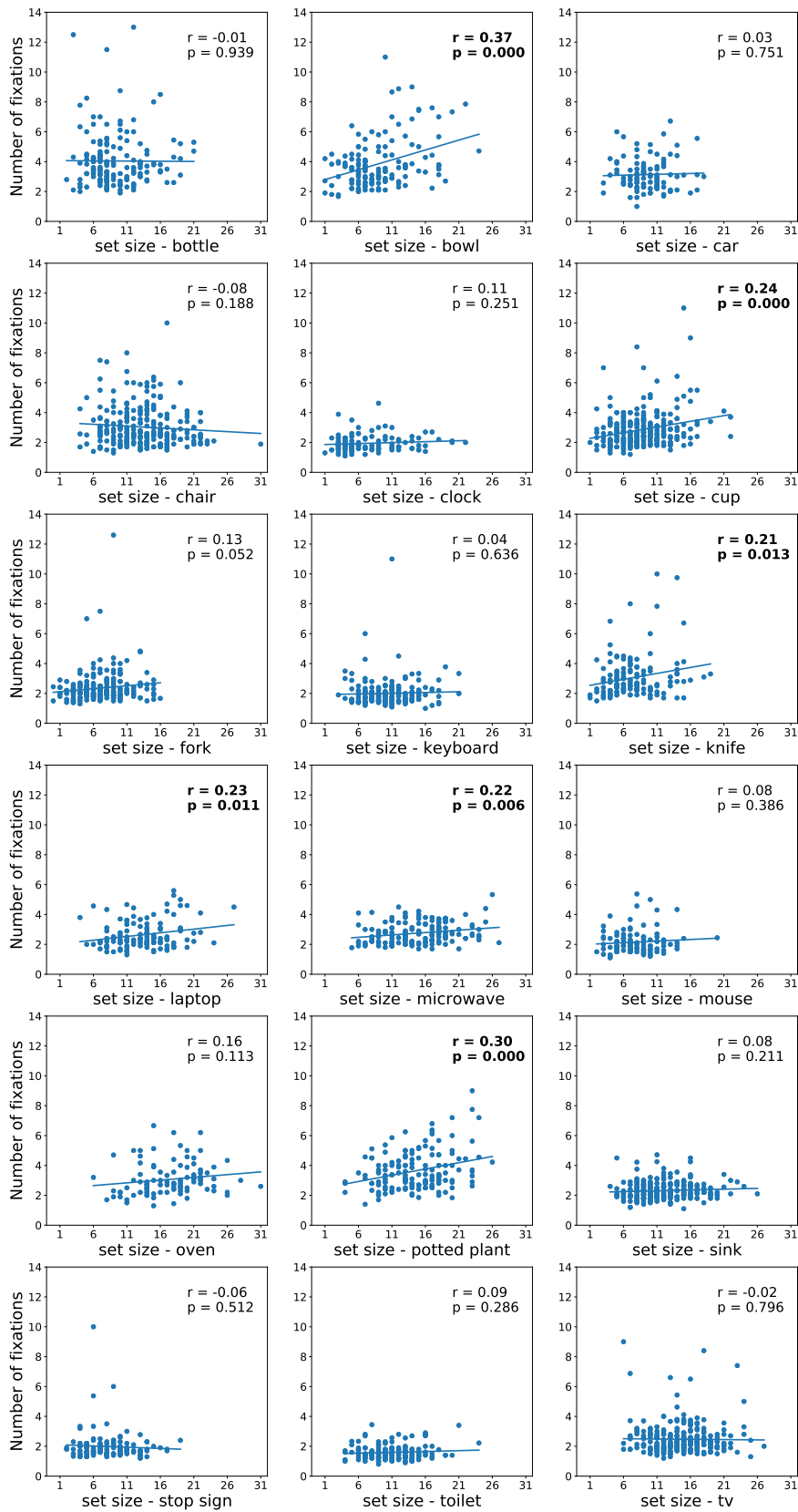


Figure S4. Number of fixations made on the target-present images plotted as a function of the set sizes of those images (using COCO object and stuff labels), averaged over participants and grouped by target category.

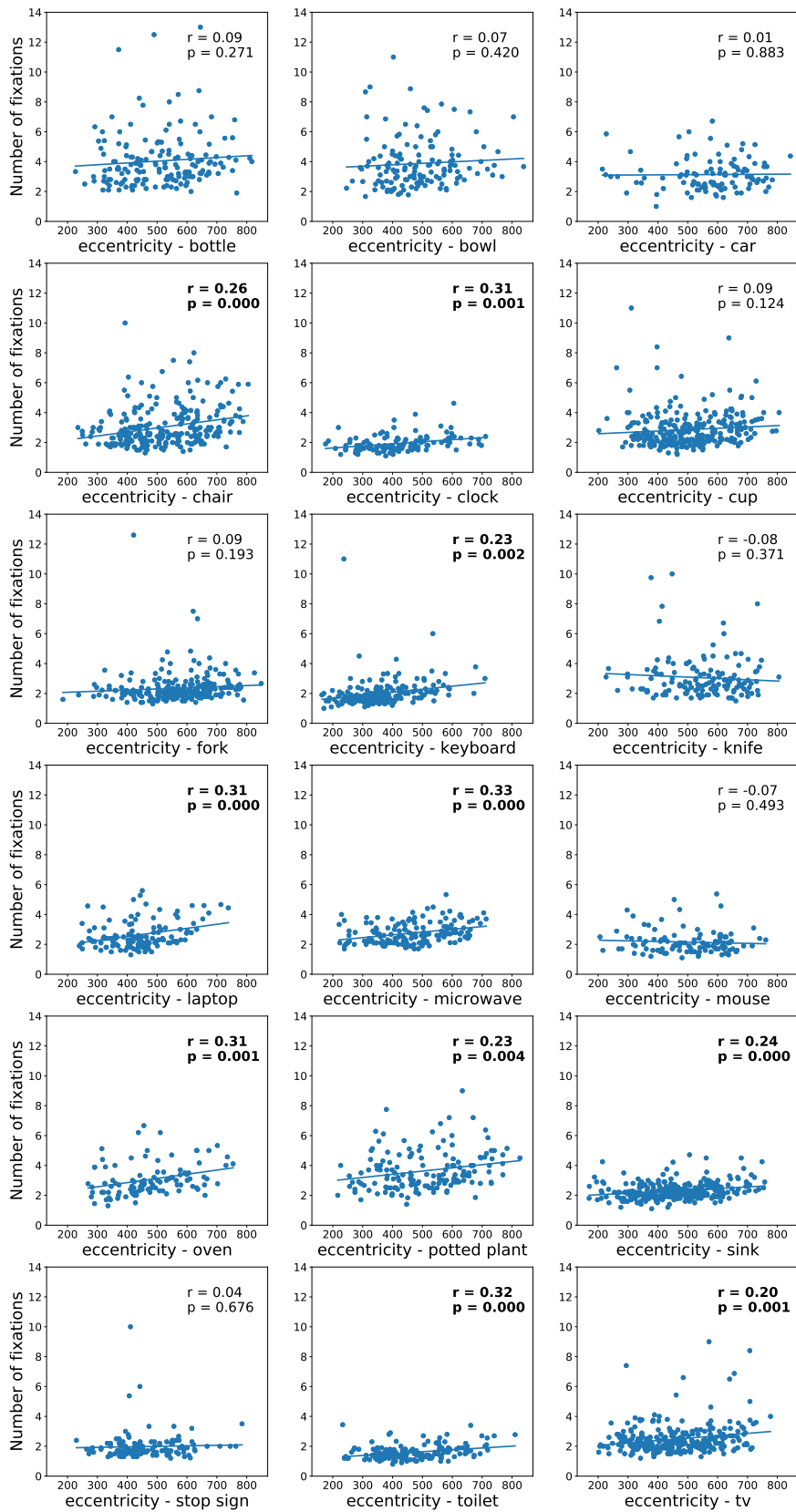


Figure S5. Number of fixations made on the target-present images plotted as a function of initial target eccentricity (using the center of the COCO bounding-box), averaged over participants and grouped by target category.

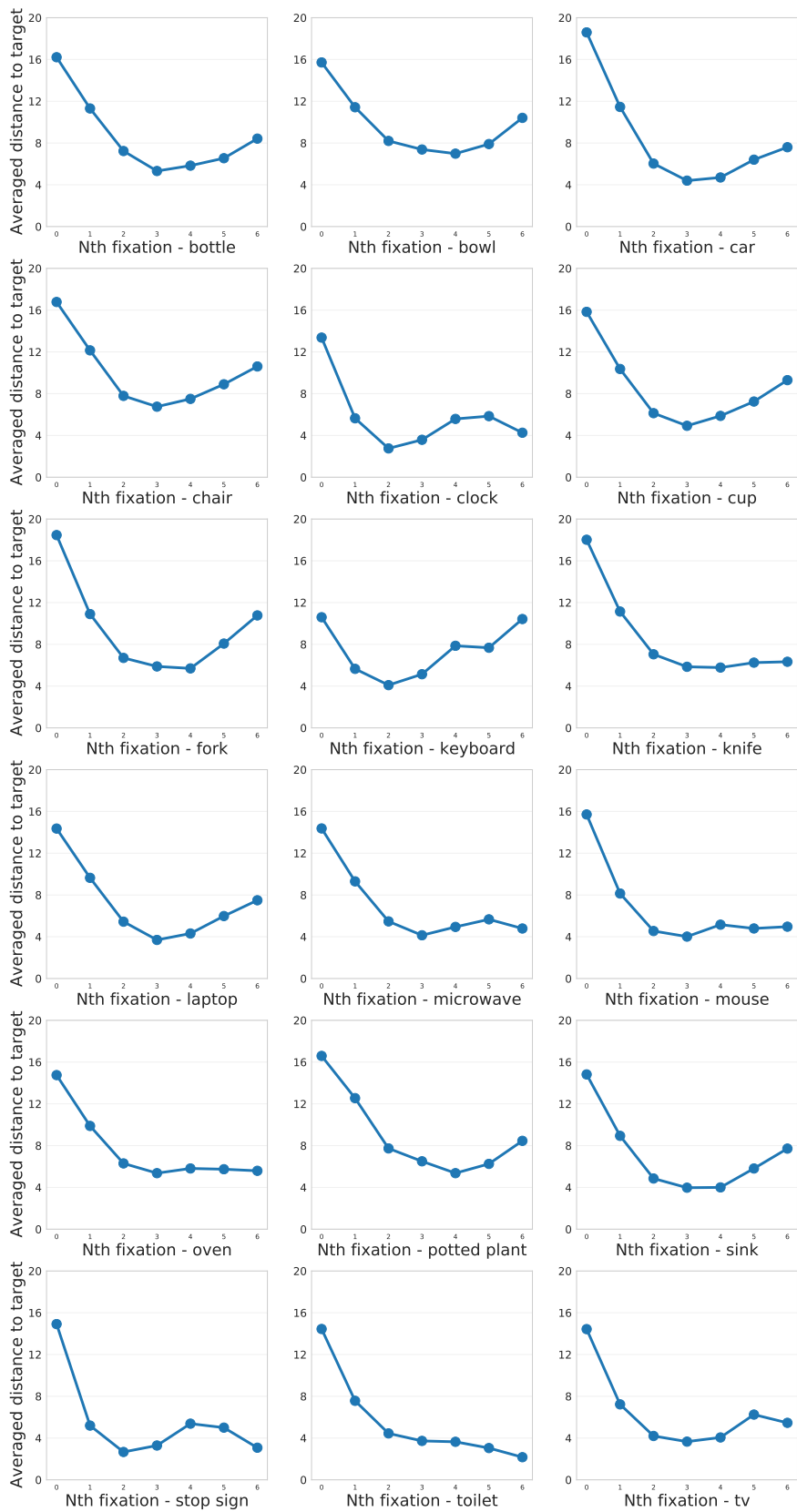
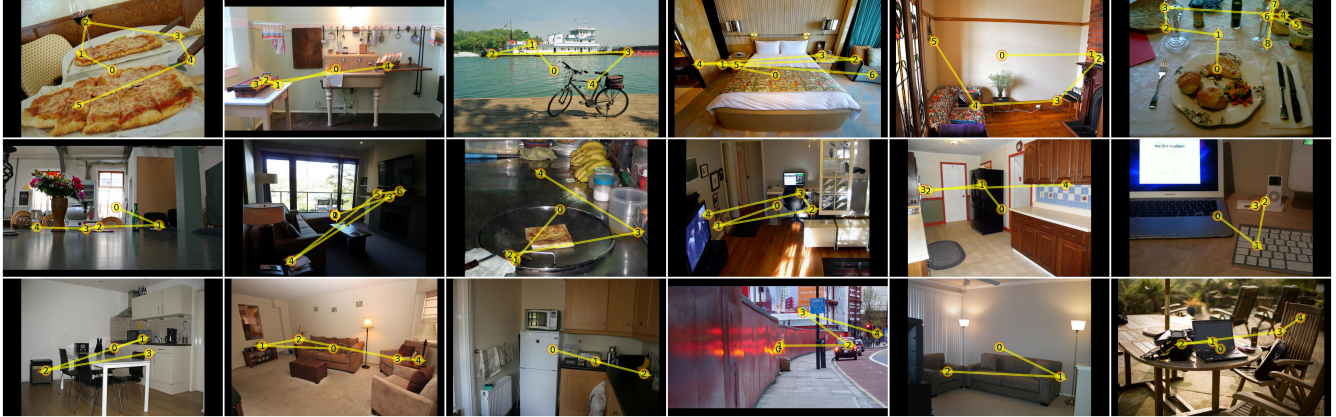


Figure S6. Averaged Euclidean distance (in visual angle) between gaze and the target's center (using COCO bounding-box labels) over the first 6 saccades, grouped by target category.

A



B

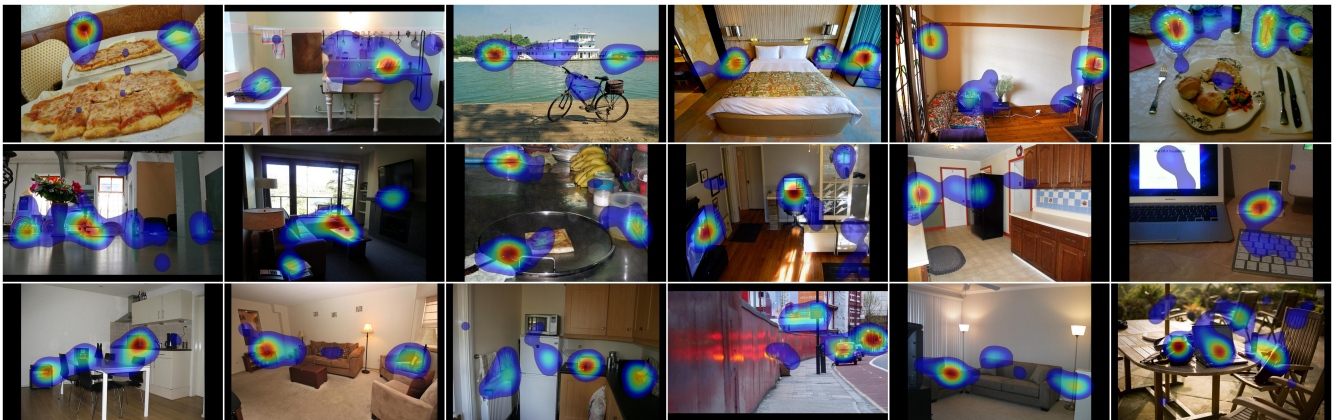


Figure S7. (A). Examples of a target-absent image for each of the 18 target categories. Yellow lines and numbered discs indicate a representative search scanpath from a single participant. From left to right, top to bottom: bottle, bowl, car, chair, (analog) clock, cup, fork, keyboard, knife, laptop,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

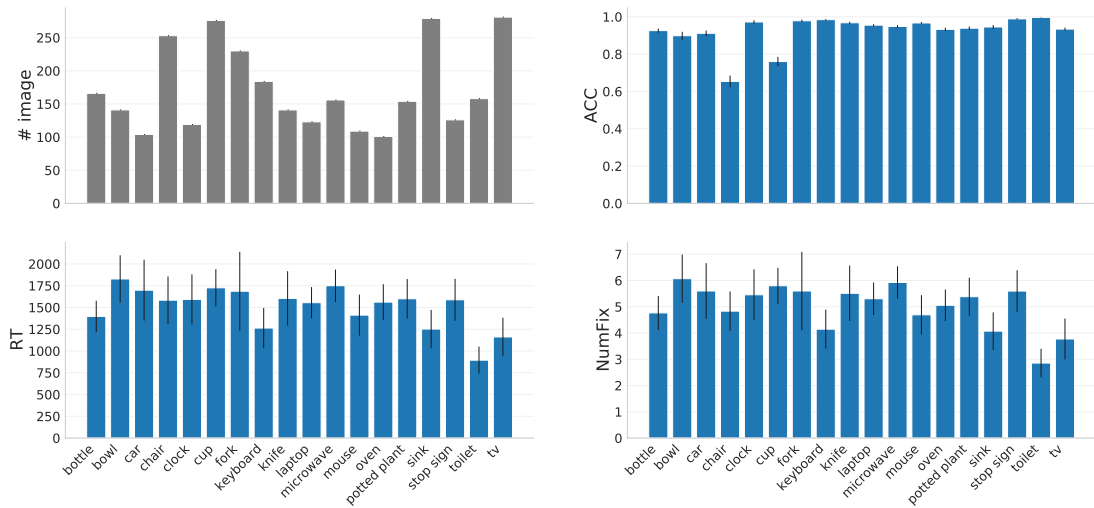


Figure S8. COCO-Search18 analyses for all 18 target categories in target-absent trials. Top: number of images in each category (gray), and response accuracy (ACC). Bottom: reaction time (RT) and number of fixations made before the button press (NumFix). Values are means over 10 participants, and error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

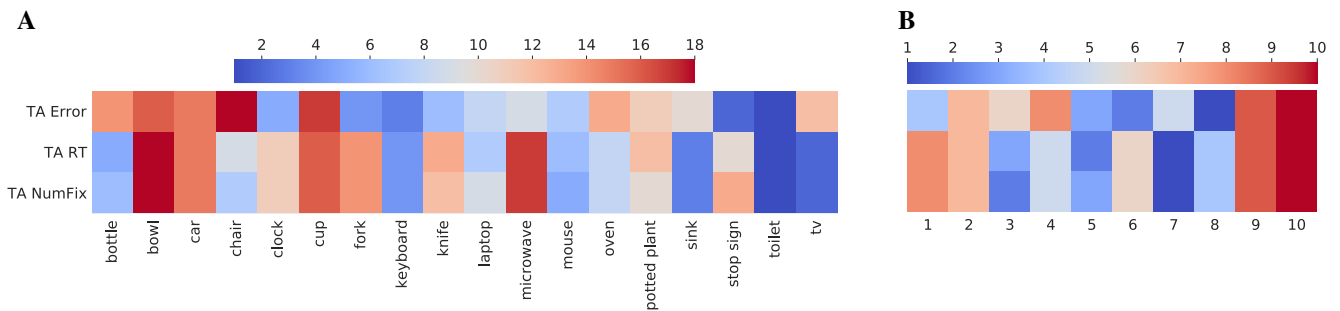


Figure S9. (A). Target-absent data, ranked [1-18] by target category (columns) and averaged over participants, shown for multiple performance measures (rows). These include: response error, reaction time (RT), and number of fixations (NumFix). Redder color indicates higher rank and harder search targets, 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).

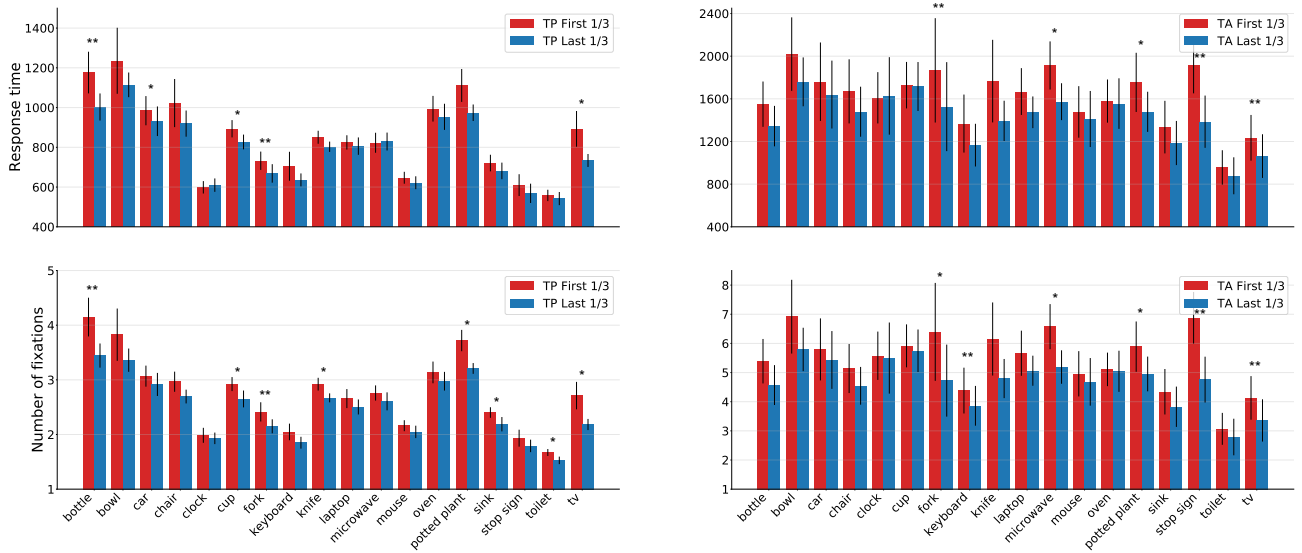


Figure S10. Practice effects, visualized as the difference in search performance between the red (first 1/3 of the trials) and the blue (last 1/3 of the trials) bars, grouped by the 18 target categories. The top row shows response time, and the bottom row shows the number of fixations before the button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

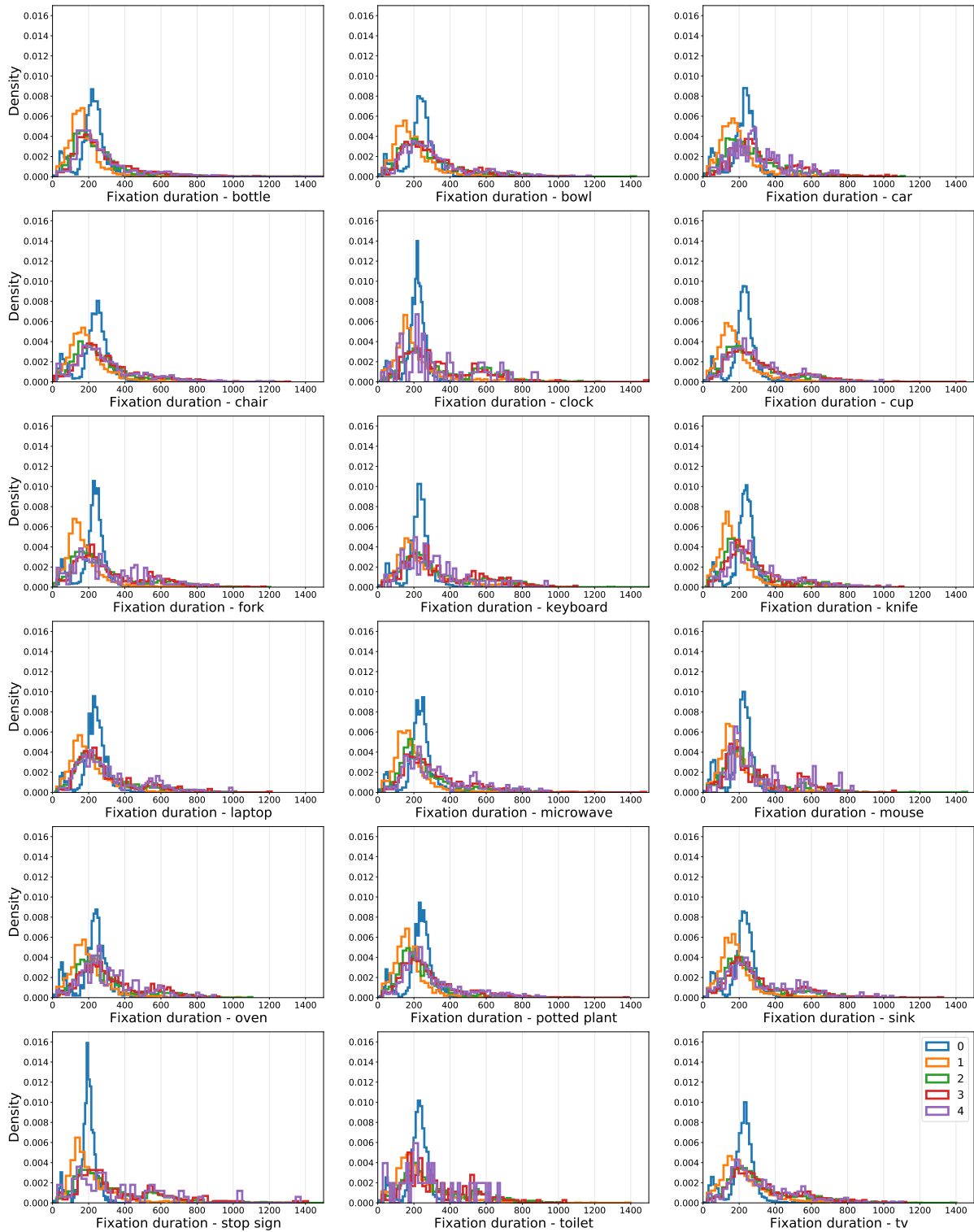


Figure S11. Density distributions of target-present fixation durations, plotted for each of the target categories (bin size = 50ms). The color lines refer to the initial fixation durations (0, blue), followed by the first four new fixations (1-4).

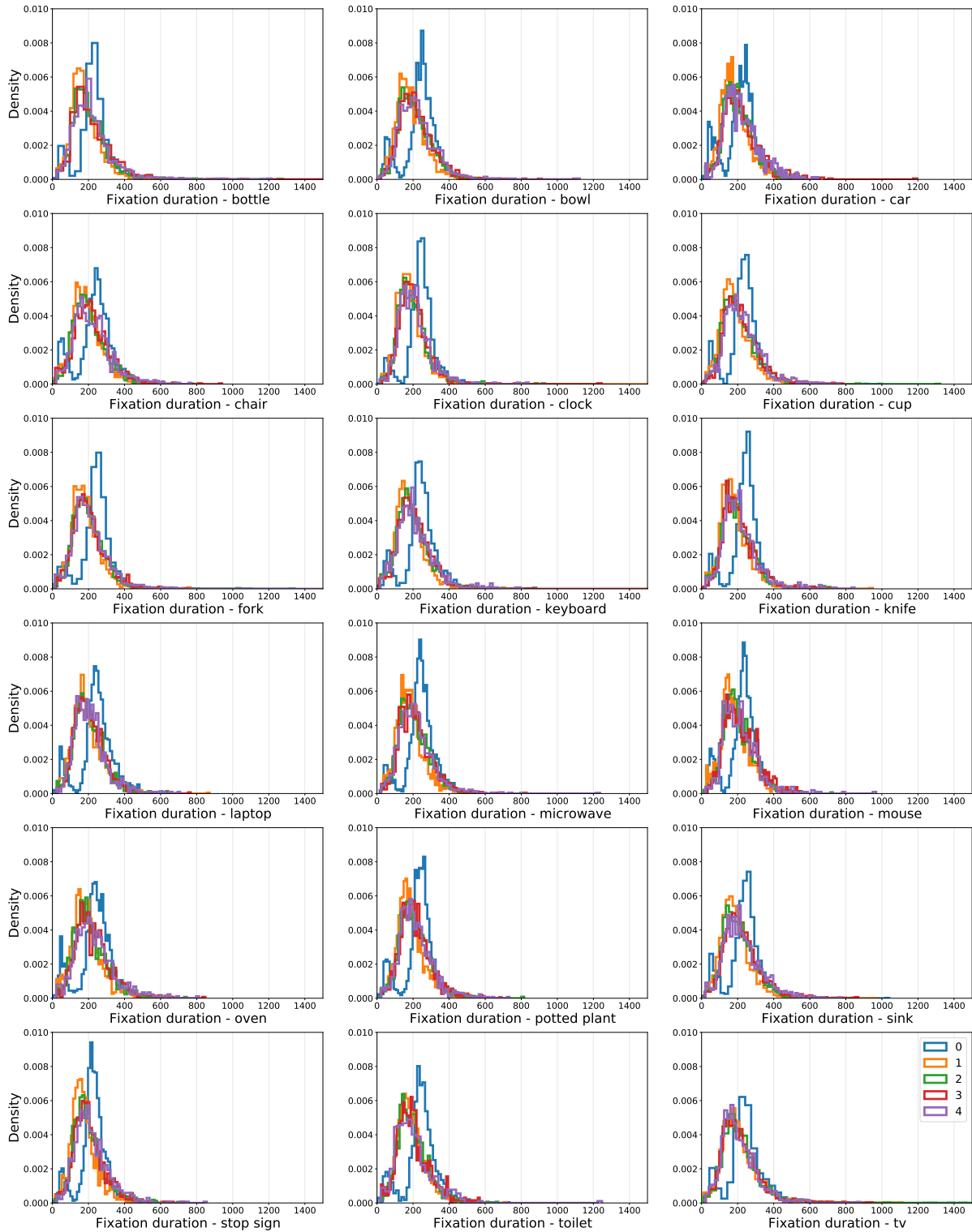


Figure S12. Density distributions of target-absent fixation durations, plotted for each of the target categories (bin size = 50ms). The color lines refer to the initial fixation durations (0, blue), followed by the first four new fixations (1-4).

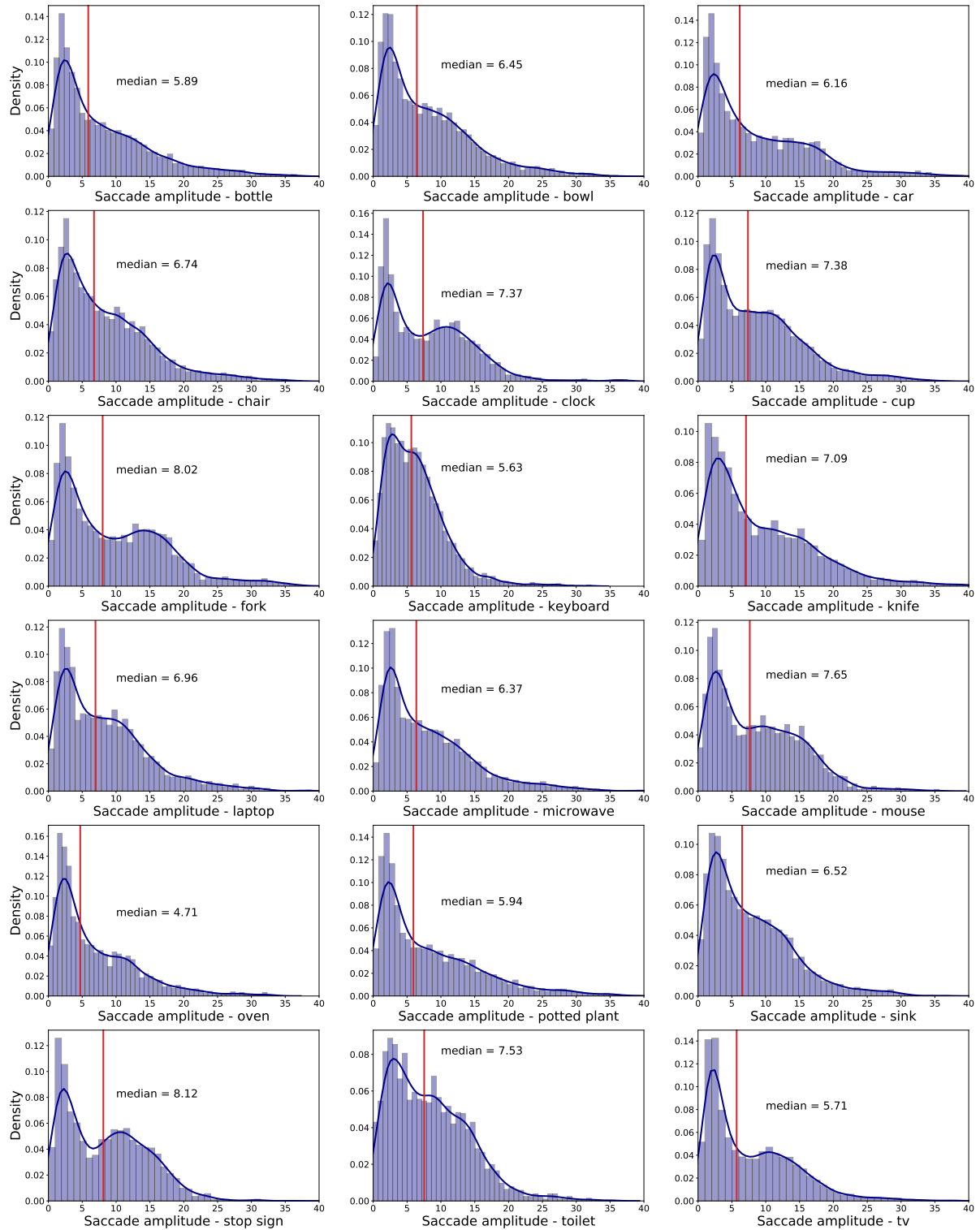


Figure S13. Density distributions of target-present saccade amplitudes (in visual angle), plotted by target category. Red vertical lines indicate median amplitudes. Dark blue lines represent Gaussian kernel density estimates.

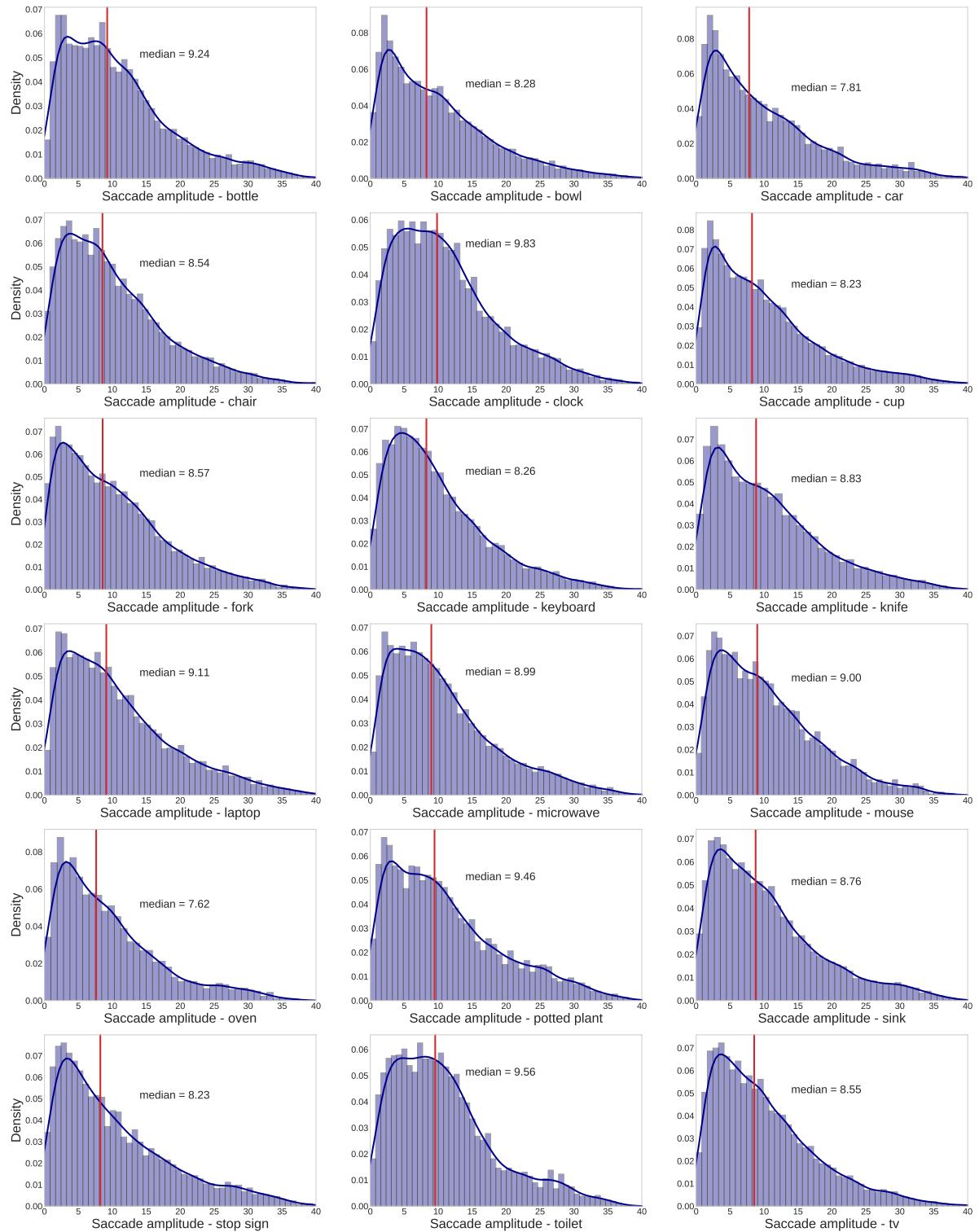


Figure S14. Density distributions of target-absent saccade amplitudes (in visual angle), plotted by target category. Red vertical lines indicate median amplitudes. Dark blue lines represent Gaussian kernel density estimates.

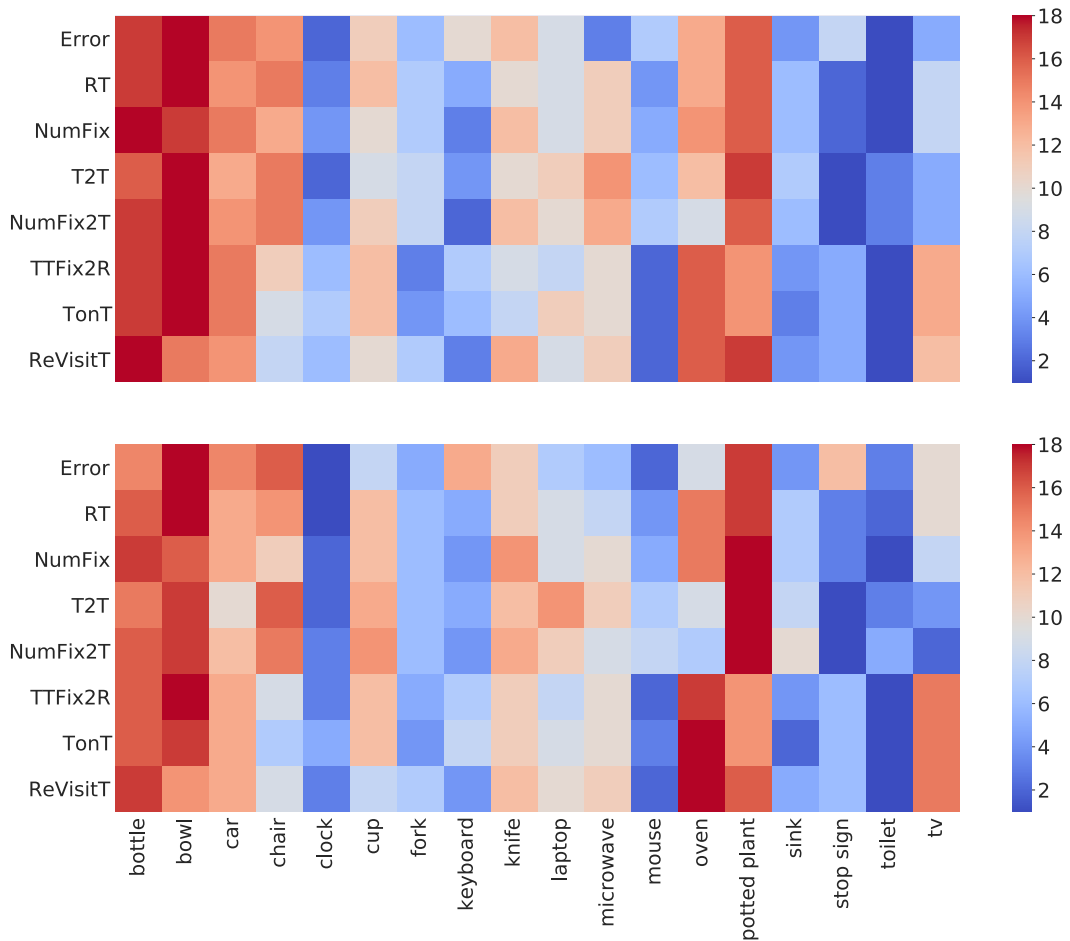


Figure S15. Target-present data, ranked by target category (1-18, columns) and shown for multiple performance measures (rows) in the trainval (top) and test (bottom) COCO-Search18 datasets. Redder color indicates higher rank and harder search targets, bluer color indicates lower rank and easier search. 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).

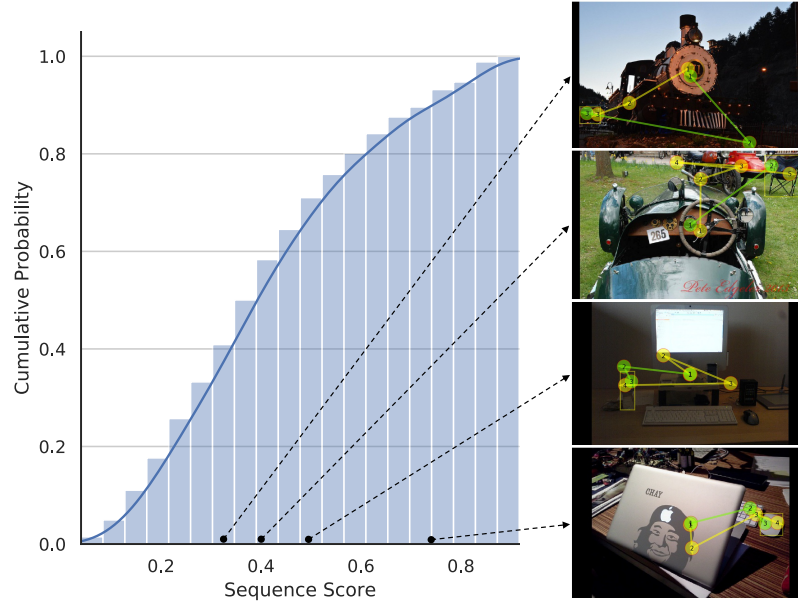


Figure S16. Left: cumulative distribution of average sequence scores computed between each scanpath generated by the IRL model and each behavioral scanpath for the test images of COCO-Search18. Right: Examples illustrating the scanpaths producing four different sequence scores. Behavioral scanpaths are 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.

A

Participants	Error	RT (ms)	NumFix	T2T (ms)	NumFix2T	TTFix2R (ms)	TonT (ms)	ReVisitT
1	0.06	993.91	2.92	372.81	1.73	769.90	717.79	1.05
2	0.09	878.03	2.81	355.50	1.88	625.61	584.05	0.75
3	0.09	780.30	2.41	349.16	1.76	651.39	622.66	0.58
4	0.10	783.61	2.49	329.89	1.96	489.23	441.90	0.45
5	0.10	761.47	2.63	308.09	1.80	544.66	525.32	0.74
6	0.07	811.03	2.96	352.92	2.12	490.26	460.37	0.71
7	0.08	633.56	2.15	310.47	1.65	429.51	415.33	0.44
8	0.07	713.69	2.44	331.08	1.79	494.28	465.60	0.60
9	0.08	1027.95	2.97	404.65	2.08	564.27	495.76	0.61
10	0.05	825.27	2.37	391.44	1.79	528.57	504.68	0.49
Mean	0.08	820.88	2.61	350.60	1.86	558.77	523.35	0.64

B

Participants	TA Error	TA RT (ms)	TA NumFix
1	0.07	1834.48	5.83
2	0.08	1384.30	4.91
3	0.07	961.80	3.07
4	0.08	1119.09	3.80
5	0.07	954.97	3.21
6	0.06	1336.61	4.92
7	0.07	897.78	2.99
8	0.05	1016.36	3.48
9	0.09	2652.84	8.02
10	0.10	2919.68	9.98
Mean	0.07	1507.79	5.02

Table S1. (A): Detailed behavioral data for 10 participants on 8 measures in target-present (TP) images. (B): Detailed behavioral data for 10 participants on 3 measures in target-absent (TA) images.

	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 S2. 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 S3. *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.