

1 **Diagnostic parts are not exclusive in the search template for real-world object categories**

2 Marcel Wurth¹ & Reshanne R. Reeder^{1,2}

3 ¹Magdeburg, Germany

4 ²Center for Behavioral Brain Sciences, Magdeburg, Germany

5

6

7

8

9

10

11

12 Corresponding author:

13 Reshanne R. Reeder

14 Otto-von-Guericke-Universität

15 Institut für Psychologie

16 39106 Magdeburg

17 Germany

18 email: reshanne.reeder@gmail.com

19

20 **1. Introduction**

21 Classical visual search theories predict search will be easy when a target is identifiable by one simple
22 feature that is easily distinguishable from distractors, and difficult when a target is composed of
23 multiple features shared by distractors (Duncan & Humphreys, 1989; Treisman & Gelade, 1980;
24 Wolfe, 1994). But imagine a search in which the target is poorly defined and always changing,
25 identifiable by no single feature, surrounded by many similar distractors, and presented in a cluttered,
26 complex environment. You might think that finding your target in this context would be impossible,
27 but such search tasks are solved easily everyday by billions of people, including toddlers (Hasegawa &
28 Miyashita, 2002), often in a matter of milliseconds and with remarkably low error rates (Thorpe, Fize,
29 & Marlot, 1996; Zelinsky, 2008). Such is the nature of real world search.

30 There is a reason why we succeed in real world search on a regular basis: targets, distractors,
31 and context are all highly trained over a lifetime of experience (Peelen & Kastner, 2014). Nevertheless,
32 the information-richness of an everyday scene is enormous (Malcolm, Groen, & Baker, 2016), while
33 the capacity of the human brain to process this is limited (Desimone & Duncan, 1995). One mechanism
34 that can be used to manage the challenges of visual search is the preparatory search template: a top-
35 down “attentional set” containing information that is relevant to the current search (Duncan &
36 Humphreys, 1989). Activating a template enhances the representation of relevant features in
37 preparation for search (Reeder, Hanke, & Pollmann, 2017), which ultimately aids target detection
38 during the search process (Malcolm & Henderson, 2009, 2010).

39 So what are the necessary contents of a real-world template? Low-level features are not helpful
40 when looking for a member of a broad category – for example, a human body – in an ever-changing
41 natural environment. In line with this, Reeder and Peelen (2013: Experiment 1) found that colors and
42 textures are not part of the search template for category search. In such cases, the template must be
43 flexible to include all relevant members of the target category (see Bravo & Farid, 2012). Previous

44 experimental research indicates that object shape may be an important feature (Reeder & Peelen, 2013;
45 Reeder, van Zoest, & Peelen, 2015). However, because natural objects can appear under varied
46 circumstances, the shapes in the template cannot be rigid (e.g., not view- or orientation-specific). There
47 is evidence that real-world object parts presented in spatially scrambled configurations (e.g., an eye
48 next to a mouth) are processed earlier in the visual hierarchy than spatially intact configurations
49 (Lerner et al., 2001), which suggests that these parts are processed like simple features (Treisman &
50 Gelade, 1980). This led to the hypothesis that a spatially flexible collection of parts dominates the
51 template for real-world category search (see Reeder & Peelen, 2013: Experiment 4). However, an
52 object “part” can be any fragment of the image (e.g., one half of the original image, a single pixel). The
53 parts that are stored in an effective template must still contain category-diagnostic information – so
54 what classifies a part as “diagnostic”?

55 Using computer simulations, Ullman, Vidal-Naquet, & Sali, (2002) found that object features of
56 “intermediate complexity” (median=11%, standard deviation=16% of the original image) were optimal
57 for classification compared to simpler or more specific features. Fragments of images could be
58 classified as belonging to a car or a human face with fewer training images if they contained parts that
59 are commonly shared by the whole class of object (e.g., the eyes of a face), whereas large image
60 fragments (e.g., the eyes, nose, and mouth) are highly specific to the individual image and cannot be
61 used efficiently to classify other objects from the same class, and simpler image fragments (e.g., a
62 dimple) could potentially have features in common with objects from the other class. In this case, as in
63 the current study, “object parts” are not defined by low-level Gestalt principles but rather by the ability
64 to classify image fragments as belonging to one basic category or the other. Applying this to human
65 research, it has been found that removing such diagnostic parts as the eyes, mouth, or limbs from an
66 animal image makes it significantly more difficult to categorize (Delorme, Richard, & Fabre-Thorpe,
67 2010).

68 It is clear from these studies that diagnostic parts are a necessary component of a flexible, real-

69 world, category-level search template. Nevertheless, this does not mean that the template is optimally
70 composed of an exclusive collection of parts. There is evidence from monkey single-cell recordings
71 that some neurons in IT (a region that processes object form) respond selectively to various views of
72 whole objects, and there is behavioral evidence that humans naturally learn to represent the global
73 shape of animate and inanimate object categories that is resistant to variations in viewpoint or changes
74 to local parts (see Logothetis & Sheinberg, 1996). Within the computational literature, Chen et al.
75 (2014) argued that the most accurate and flexible categorical representations should contain various
76 combinations of the “root” (the holistic representation) and its parts. Stemming from this, we
77 hypothesize that the most accurate and flexible search template will contain information about both
78 whole and parts. This, however, has not yet been tested experimentally.

79 In our previous studies (Reeder & Peelen, 2013; Reeder et al., 2015) we used a novel design
80 that combined visual search and “contingent attention capture” (Folk, Remington, & Johnston, 1992)
81 that allowed us to investigate different possible features of the naturalistic search template without
82 changing the search task (which would invariably change the search template). Specifically, subjects
83 were cued to search for object categories (cars, people) on every trial. On the majority of trials, the cue
84 was followed by a search task – however, this task was only included to ensure that subjects activated a
85 category-level search template following the cue. The critical condition occurred on a minority of trials
86 intermixed with the search task – in which subjects were required to indicate whether a briefly
87 presented dot probe appeared on the left or right of fixation. The dot was always preceded by features
88 of search targets (e.g., textures, shapes), but subjects were instructed to ignore these, and to respond as
89 quickly and accurately as possible to the location of the dot. We hypothesized that viewed features
90 presented on dot-probe trials that matched the active template would capture attention involuntarily (as
91 indicated by faster responses to a subsequent dot-probe on the same side of fixation as template-
92 matching features), even if they are task-irrelevant and presented in search-irrelevant locations (Seidl-
93 Rathkopf, Turk-Browne, & Kastner, 2015; Wyble, Folk, & Potter, 2013). In other words, any attention

94 capture by such irrelevant items must be due to their matching the template rather than being used
95 explicitly to improve search performance.

96 In the current study, we present three experiments that test the efficacy of parts versus wholes in
97 capturing attention during real-world category search (see Figures 1 and 2). Experiment 1 is a direct
98 replication of Experiment 4 of Reeder & Peelen (2013), only with a larger number of subjects (25
99 compared to 13). We hypothesized that higher power associated with a larger N (as calculated by a
100 power analysis) could potentially bring out small differences in efficacy between parts and wholes to
101 capture attention. In Experiment 2, we presented capture stimuli in search-irrelevant locations to
102 determine whether capture differences between parts and wholes could be due to subjects activating an
103 inflexible, spatially rigid search template. Finally, Experiment 3 was conducted to determine whether
104 capture differences could be attributed to some parts being less diagnostic than others; to control for
105 this, we presented collections of four object parts and compared them to whole objects as capture
106 stimuli. In all experiments, we found capture effects for wholes but not parts, suggesting parts alone are
107 not a sufficient search template for real-world object categories.

108

109 **2. Materials and Methods**

110 *2.1 Subjects*

111 90 students and faculty of Otto-von-Guericke University, Magdeburg, were recruited for this study (29
112 for Experiment 1, 30 for Experiment 2, and 31 for Experiment 3). Nine subjects took part in more than
113 one experiment. All subjects had normal or corrected to normal vision, received psychology credits or
114 money as reimbursement for their participation, and provided written informed consent prior to
115 experimentation. These measures conformed to the Declaration of Helsinki and were approved by the
116 research ethics committee of Otto-von-Guericke University.

117 Outlier exclusion was strict to control for possible confounds in subject performance. Subjects
118 were excluded if their mean visual search accuracy was below 75%, as in the previous study. This was

119 to ensure that subjects could activate an appropriate category-level search template. Subjects were also
120 excluded if their mean dot-probe detection accuracy was below 90%. In the dot-probe task (see 2.3
121 *Procedure*ed, and accuracy
122 lower than 90% would suggest a lack of attention or inability to understand the task rather than normal
123 human error. Furthermore, high dot-probe accuracy was necessary for sufficient power to analyze
124 reaction times (RT; only reported for accurate trials) due to a low number of dot-probe trials per run
125 (16, as in the previous study). Following these exclusions, subjects were lastly excluded if their mean
126 dot-probe RT was slower than two standard deviations from the group mean. This was done to ensure
127 that subjects did not consciously deliberate about their responses based on the capture stimuli that
128 appeared prior to the dot.

129 Because of the strict exclusion criteria and to ensure an adequate number of subjects per
130 experiment, we continued to run subjects in each experiment until a final number of 25 suitable
131 subjects per experiment was reached. We determined that 25 subjects would provide adequate power in
132 each experiment following an analysis run in GPower (Faul et al., 2009). From the results of
133 Experiment 4 of Reeder & Peelen (2013), a power analysis (power (1- β) set at 0.95 and $\alpha = 0.05$, two-
134 tailed) indicated that to find a reliable difference in RT between consistent and inconsistent dot-probe
135 trials on which silhouettes of object parts appeared, we would need to run 22 subjects. Although a
136 sample size of 12 was found to be adequate based on a power analysis run on the capture effects of
137 whole silhouettes, we used a sample size that could reveal any true capture effect by parts on their own
138 because we specifically sought to target potential RT differences between parts and wholes.

139 Four subjects were excluded from Experiment 1 (two due to low search accuracy, and one each
140 due to low dot-probe accuracy and slow dot-probe RT); five subjects were excluded from Experiment 2
141 (three due to low search accuracy, one due to low dot-probe accuracy, and one due to slow dot-probe
142 RT); and six subjects were excluded from Experiment 3 (three due to low search accuracy, two due to
143 both low search accuracy and low dot-probe accuracy, and one due to slow dot-probe RT). 75 subjects

144 contributed to the final results; 25 in Experiment 1 (age range = 18-30 years, mean age = 21.64 years, 5
145 men, 2 left-handed), 25 in Experiment 2 (age range = 19-30 years, mean age = 22.35 years, 4 men, 1
146 left-handed), and 25 in Experiment 3 (age range = 19-27 years, mean age = 22.05 years, 6 men, 1 left-
147 handed).

148

149 *2.2 Stimuli*

150 All stimuli were presented on a 24-inch Samsung with a 1920 x 1080 screen resolution and a 60 Hz
151 refresh frequency (Samsung Electronics Co., Ltd., Suwon, South Korea). Stimuli were created in
152 Python using PsychoPy functions (Peirce, 2008). The fixation cross and the letter cues “A” and “M”
153 were text stimuli with a height of 31 pixels (0.8 cm) for the fixation and 70 pixels (1.8 cm) for the letter
154 cues, presented in the center of the screen in Times New Roman font. Subjects viewed all stimuli from
155 a free-viewing distance of approximately 57 cm.

156

157 *2.2.1 Search Stimuli*

158 Search stimuli were natural scene photographs used in Reeder & Peelen (2013; see Figure 1a). 960
159 total scenes were used, 240 containing cars but not people, 240 containing people but not cars, 240
160 containing both cars and people, and 240 containing neither cars nor people. No scene was viewed
161 twice. Scenes contained objects in various natural viewpoints, levels of occlusion, distances from the
162 observer, colors, sizes, genders (for people), and makes and models (for cars). Furthermore, one or
163 multiple targets could appear in the same scene. Our goal with the variety of our images was to
164 simulate natural everyday views as well as possible, and to encourage subjects to activate realistic
165 templates for real-world search. Each scene had a size of 548 x 411 pixels, corresponding to 13.9 x
166 10.4 cm. Two scenes were presented to the left and right of fixation with 76 pixels (1.9 cm) separating
167 the two images in Experiment 1, or above and below fixation with 261 pixels (6.6 cm) separating the
168 two images in Experiments 2 and 3. The larger spatial separation in the latter experiments was to

169 ensure that the scenes did not spatially overlap with the capture stimuli.

170

171 2.2.2 Capture stimuli

172

173 Capture stimuli were whole silhouettes of cars and people, and silhouettes of parts of cars and bodies,
174 used in Reeder & Peelen (2013; see Figure 1b). People were presented without heads to remain
175 consistent with the previous study. 160 different images were used, each with 80 cars and 80 people in
176 different perspectives and positions. Silhouette parts were created by cutting out sections of the total
177 area of each of the whole silhouette images (range = 4.61%-24.19%, mean = 14.66% for cars, and
178 range = 7.68%-23%, mean = 14.27% for people). A test designed to ensure that the part pictures were
179 clearly discriminable as belonging to a car or person showed that subjects ($N = 8$) could correctly
180 discriminate the object categories 97.1% of the time (Reeder & Peelen, 2013). Parts were scaled in a
181 way that allowed them to be presented in the same sizes and locations as whole silhouettes. In
182 Experiments 1 and 2, silhouettes were presented in one of three sizes, with the longest dimension of the
183 image at 100, 180, or 200 pixels (2.5, 4.6, or 5.1 cm), determined randomly for each stimulus on a trial-
184 by-trial basis. These were presented at one of three distances from the center of the screen to the center
185 of the image at 130, 160 or 250 pixels (3.3, 4.1, or 6.4 cm), also determined randomly and
186 independently for each image. This is in accordance with the stimulus parameters detailed in Reeder &
187 Peelen (2013).

188 Part collections in Experiment 3 were 160 newly generated images (80 per image category)
189 composed of four silhouette parts arranged in a 2 x 2 imaginary bounding box (see Figure 1c). Images
190 were randomly selected from the set of 80 single part images per category. Because there were only 80
191 single part images per category, parts were repeated across different part collections. Each part was
192 randomly assigned one of the four locations in the collection, and no precise combination of parts and
193 part locations was repeated. The total area covered by part collections was determined based on
194 visibility of the individual parts, with each image randomly assigned to a size of 180 x 180, 200 x 200

195 or 280 x 280 pixels (4.6 x 4.6, 5.1 x 5.1, or 7.1 x 7.1 cm) on a trial-by-trial basis. Distance between the
196 center of fixation and the center of the image could be 160, 200, or 250 pixels (4.1, 5.1, or 6.4 cm),
197 randomly assigned to each stimulus separately. This kept part collections within the field of vision
198 while remaining identifiable.

199

200 Figure 1. All stimuli are taken from Reeder & Peelen (2013). a.) examples of search stimuli. b.)

201 examples of whole silhouette capture stimuli. c.) examples of silhouette part capture stimuli. In

202 Experiment 2, only one part from each category was shown in isolation on either side of fixation. In

203 Experiment 3, four parts from each category appeared equally spaced in a 2 x 2 grid as depicted here.

204 Images are not shown to scale.

a.) Natural scenes



b.) Whole silhouettes



c.) Silhouette parts



205

206 *2.3 Procedure*

207 Subjects performed two different intermixed tasks in the experiment (see Figure 2). In the search task
208 (75% of trials, 48 trials per block), subjects were required to indicate, with the directional arrow keys,
209 which of two scenes presented to the left and right of fixation (Experiment 1) or above and below
210 fixation (Experiments 2 and 3) contained a cued category—a car or a person. Cars and people were
211 chosen as the search categories for three reasons: first, they were used in our previous study and we
212 wanted to follow the methods as closely as possible; second, both are highly familiar and are viewed
213 every day by most people; and finally, they are both highly variable categories that can contain myriad
214 different features (colors, sizes, shapes) – and yet both are quickly and easily identified. This latter
215 point highlights the importance of discovering how such information-rich and variable objects are
216 represented at a categorical level in the search template. The letters “A” (for the German word *Auto*,
217 which means car) or “M” (for the German word *Mensch*, which means human), presented prior to the
218 scenes and intermixed an equal number of times within a block, cued subjects as to which category to
219 look for. Each scene pair either contained cars in one image and people in the other, or cars and people
220 in one image and neither in the other. This made it so that both the relevant and irrelevant category
221 appeared on every trial. Each scene type (car, person, both, neither) appeared an equal number of times
222 on either side of fixation. A complete search trial consisted of a 500 ms fixation, followed by a 500 ms
223 letter cue, 1000 ms fixation, 67 ms search scenes followed by a 350 ms mask, and a final 1660 ms of
224 fixation. Participants were instructed to use the arrow keys (right and left in Experiment 1, up and
225 down in Experiments 2 and 3) to indicate in which scene the target category appeared.

226 In the dot-probe task (25% of trials, 16 trials per block), subjects were instructed only to
227 respond to the location of a black dot that appeared on the left or right of fixation, using the left and
228 right arrow keys. Capture stimuli (the silhouettes) appeared prior to the dot, but subjects were told to
229 ignore these to the best of their ability. One car and one body silhouette appeared on each capture trial
230 and appeared an equal number of times on the left and right of fixation. The location of the target-
231 matching silhouette was tied to the same side of fixation as the dot-probe on half of trials (consistent

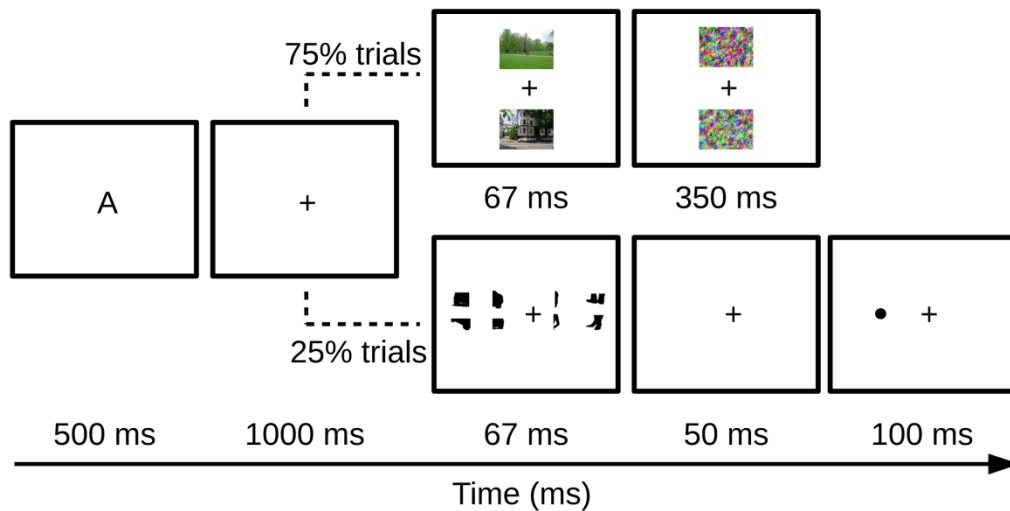
232 trials). The target-matching silhouette appeared on the opposite side of fixation from the dot-probe on
233 the other half of trials (inconsistent trials). In Experiments 1 and 2, capture stimuli were either both
234 whole silhouettes or both silhouette parts, with these trials randomly mixed within a block. In
235 Experiment 3, capture stimuli were either both whole silhouettes or both collections of silhouette parts.
236 A capture trial started the same as a search trial, with a 500 ms fixation, followed by a 500 ms letter
237 cue, and 1000 ms fixation. Because trials from the search task and dot-probe task were mixed, subjects
238 did not know what type of stimuli would appear following the cue. This ensured that subjects would
239 form a search template on every trial. Following the appearance of the two capture stimuli for 67 ms,
240 there was a fixation for 50 ms, the dot-probe for 100 ms, and a final fixation of 1660 ms.

241 Trials were organized in ten 64-trial blocks, with the first block being a practice block that did
242 not contribute to the analysis. Before each experiment, a short, slowed practice block was performed
243 under the supervision of the experimenter, to ensure that each subject understood the task correctly.

244

245 Figure 2. A schematic of the search task (75% of trials) and dot-probe task (25% of trials). Here is an
246 example of the search scenes appearing above and below fixation (Experiments 2 and 3) in the search
247 task and collections of silhouette parts (Experiment 3) appearing in the dot-probe task. In Experiment
248 1, search scenes appeared to the left and right of fixation. In Experiments 1 and 2, single silhouette
249 parts were shown on half of dot-probe trials.

250



251

252

253 2.4 Analysis

254 Our main goal was to find out if subjects performed better on consistent dot-probe trials than on
 255 inconsistent dot-probe trials. A consistent trial is defined as a trial on which the cue-matching silhouette
 256 (e.g., the car silhouette following the “A” cue) appears on the same side of fixation as the dot-probe.
 257 An inconsistent trial is a trial on which the cue-matching silhouette appears on the other side of fixation
 258 from the dot-probe. We calculated RT from the onset of the dot-probe. Only correct trials were input
 259 into the RT analysis. Attention capture was defined as faster RT on consistent trials compared to
 260 inconsistent trials.

261

262 3. Results

263 3.1 Experiment 1

264 Experiment 1 was a replication of Experiment 4 in Reeder & Peelen (2013). Search scenes, capture
 265 stimuli, and dot-probes were all presented to the left and right of fixation. Capture stimuli were either
 266 whole silhouettes or single silhouette parts. Our aim was to see if the earlier results could be replicated
 267 with a larger sample, and if the resulting increase in power could reveal a more detailed understanding

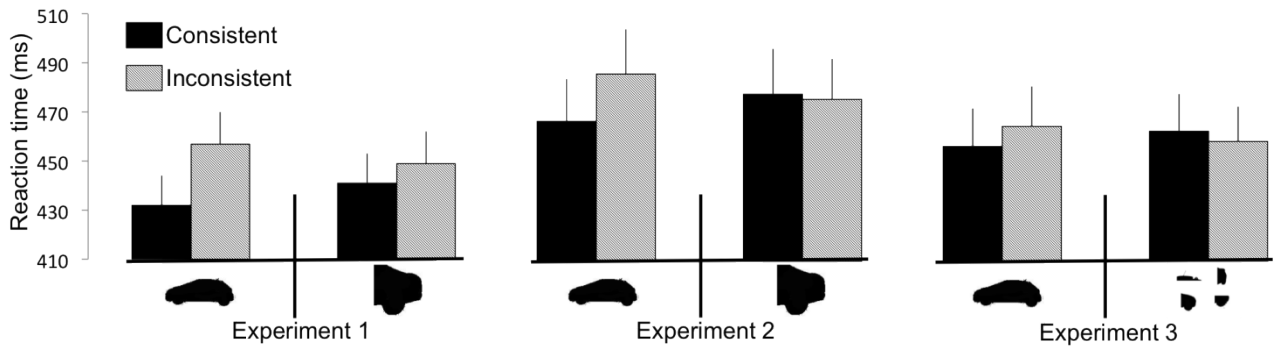
268 of the role that diagnostic object parts play in the search template.

269 RTs were submitted to a 2 x 2 repeated-measures ANOVA, with consistency (consistent,
270 inconsistent) and silhouette type (whole object, object part) as factors (see Figure 3). There was a main
271 effect of consistency ($F(1,24) = 20.918, p < 0.001, \eta_p^2 = 0.466$), indicating generally faster responses on
272 consistent trials (mean = 436 ms, standard error (SE) = 11 ms) compared to inconsistent trials (mean =
273 453 ms, SE = 13 ms). There was no main effect of silhouette type ($F(1,24) = 0.050, p = 0.825, \eta_p^2 =$
274 0.002), but a significant interaction between consistency and silhouette type ($F(1,24) = 9.322, p =$
275 0.005, $\eta_p^2 = 0.280$). Paired-samples t-tests revealed that the consistency effect for whole silhouettes was
276 significant ($t(24) = -4.984, p < 0.001$; consistent RT mean = 432 ms, SE = 12 ms; inconsistent RT
277 mean = 457 ms, SE = 13 ms), whereas the consistency effect for silhouette parts was not ($t(24) = -$
278 1.848, $p_{\text{stent}} \text{ RT mean} = 441 \text{ ms, SE}$
279 = 12 ms; inconsistent RT mean = 449 ms, SE = 13 ms).

280

281 Figure 3. Error bars show the standard error of the mean. RT results for consistent and inconsistent dot-
282 probe trials are labeled by experiment. The silhouette shown below the bars represents whether the
283 capture stimuli for those data were whole silhouettes (represented by a whole car silhouette here),
284 single silhouette parts (represented by a single car part silhouette here), or collections of parts
285 (represented by a collection of car parts here).

286



287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305 3.2 Experiment 2

306

The previous finding that diagnostic object parts capture attention in a comparable way to

whole objects (Reeder & Peelen, 2013: Experiment 4) could not be replicated in a larger sample size.

However, we were still able to find a reliable capture effect by whole objects, suggesting that whole

object shape may be a better representation of the search template. Nevertheless, it is possible that

subjects were not using a flexible template for this task for whatever reason (perhaps German students

rely on more rigid templates than Italian students). In our previous study, we found that the category-

level template was not only flexible in terms of changes in viewpoint, but also in terms of location

around the display. The representation of features without a particular spatial configuration (such as the

representation of object parts) relies on rapid, spatially non-specific processing (e.g., Quinlan, 2003;

Singh & Hoffmann, 2001), so it is possible that the representation of parts in the template is also

associated with a spatially global representation. Therefore, in Experiment 2, we presented capture

stimuli and search stimuli in separate, non-overlapping locations. If the capture effect for whole

silhouettes is impaired in the next experiment, it would suggest a reliance on a spatially focused search

template. Alternatively, a replication of the pattern of results seen in Experiment 1 would provide

evidence that even a flexible, spatially non-specific template is not exclusively composed of parts of

objects.

Experiment 2 was the same as Experiment 1 with the exception that search scenes now appeared above

307 and below fixation, with capture stimuli and dot-probes still appearing to the left and right. RT results
308 were submitted to a 2 x 2 repeated-measures ANOVA with consistency and silhouette type as factors
309 (see Figure 3). This revealed a main effect of consistency ($F(1,24) = 20.096, p < 0.001, \eta_p^2 = 0.456$),
310 again indicating faster RT on consistent trials (mean = 471 ms, SE = 17 ms) compared to inconsistent
311 trials (480 ms, SE = 17 ms). There was again no main effect of silhouette type ($F(1,24) = 0.195, p =$
312 $0.663, \eta_p^2 = 0.008$), but a significant interaction between consistency and silhouette type ($F(1,24) =$
313 $13.638, p = 0.001, \eta_p^2 = 0.362$). Paired-samples t-tests revealed the same pattern as before, with whole
314 silhouettes showing a significant consistency effect ($t(24) = -5.353, p < 0.001$; consistent RT mean =
315 466 ms, SE = 17 ms; inconsistent RT mean = 485 ms, SE = 18 ms), and silhouette parts showing none
316 ($t(24) = 0.536, p = 0.475$,
317 SE = 16 ms). This shows that the search template for object categories is spatially global, and optimally
318 composed of information about both diagnostic parts and their configuration around the whole.

319 The argument arises that perhaps the reason for a lack of a capture effect with object parts is
320 because, despite prior evidence of high general diagnosticity, it is possible that some object part images
321 are less diagnostic than others. Furthermore, a whole object contains multiple diagnostic parts that
322 could all contribute to the capture effect. This could lead to a natural disadvantage for the silhouette
323 part trials compared to the whole silhouette trials. To control for this, we conducted Experiment 3,
324 which presented collections of four diagnostic parts instead of one on silhouette part capture trials. If
325 some parts are less diagnostic than others, or if stronger capture by whole objects is due to a summation
326 of diagnostic part information, then part collections should prove to be a more equatable condition to
327 the whole silhouettes.

328

329 3.3 Experiment 3

330 Experiment 3 was the same as Experiment 2 except now a collection of four parts instead of one

331 appeared on either side of fixation on silhouette part capture trials. A 2 x 2 repeated-measures ANOVA
332 was conducted with consistency and silhouette type (whole object, part collection) as factors (see
333 Figure 3). This revealed no main effects (consistency: $F(1,24) = 0.557, p = 0.463, \eta_p^2 = 0.023$;
334 silhouette type: $F(1,24) < .001, p = 0.994, \eta_p^2 < 0.001$), but a significant interaction between
335 consistency and silhouette type ($F(1,24) = 8.472, p = 0.008, \eta_p^2 = 0.261$). Paired-samples t-tests
336 revealed a significant consistency effect for whole silhouettes ($t(24) = -2.069, p = 0.049$), with
337 consistent trials showing faster RT (mean = 456 ms, SE = 15 ms) than inconsistent trials (464 ms, SE =
338 16 ms). Conversely, part collections showed no consistency effect ($t(24) = 1.268, p = 0.217$), with
339 consistent trials (mean = 462 ms, SE = 15 ms) showing no performance advantage compared to
340 inconsistent trials (mean = 458 ms, SE = 14 ms).

341 These results indicate that the null effect for parts as seen in the previous experiments is not
342 likely due to a lack of diagnostic information being present in the capture stimuli, but rather due to
343 parts being a less accurate match to the template than whole silhouettes.

344

345 *3.4 Cars versus People*

346 It could be argued that the lack of a capture effect for parts reported in these experiments is driven by a
347 strong, natural holistic representation of conspecifics (Reed, Stone, Bozova, & Tanaka, 2003)
348 compared to other objects. In other words, perhaps diagnostic parts are the favored template for cars
349 (and generally other objects), but a holistic and configural representation of human bodies steers the
350 results of the current study. This is unlikely, since the bodies in our study were headless, and research
351 suggests that faces on bodies is what necessitates such effects (Brandman & Yovel, 2010; Yovel, Pelc,
352 & Lubetzky, 2010). Nevertheless, to test for this, we analyzed the consistency effects in our
353 experiments separately for cars and people. Because separating the results on this added dimension
354 leads to few trials per condition (2 per experiment run), we analyzed all 75 subjects together to increase

355 power and decrease noise.

356 We conducted a 2 x 2 x 2 repeated-measures ANOVA with category (cars, people), consistency
357 (consistent, inconsistent), and silhouette type (whole objects, object parts) as factors. There was no
358 main effect of category ($F(1,74) < 0.001$, $p = 0.996$, $\eta_p^2 < 0.001$), and no interaction between category
359 or any other factor (all $F_s < 1.388$, all $p_s > 0.242$, all $\eta_p^2 < 0.019$). These results suggest that the reported
360 capture effects cannot be attributed to a dominant holistic representation of people compared to cars.

361

362 **4. General Discussion**

363 The results of the three experiments presented here provide a critical examination of the contribution of
364 diagnostic object parts and whole object shape to the search template for real-world category search.

365 Whole object silhouettes, presented as task-irrelevant distractors, captured attention across all
366 experimental variations both spatially specific and global. This suggests that whole object shape is
367 represented efficiently and flexibly in preparation to search for object categories that appear in
368 naturalistic contexts. The results are far less promising for diagnostic object parts — when presented at
369 task-irrelevant locations they do not capture attention at all, and if presented in task-relevant locations,
370 only a weak trend emerges. The current study therefore favors a template composed of both diagnostic
371 parts and the whole.

372 In a direct replication of Experiment 4 from Reeder & Peelen (2013), the capture effect for
373 object parts was far weaker than for whole objects, and did not reach statistical significance. One basic
374 axiom of statistical testing might lend an explanation for the difference in results between the current
375 Experiment 1 and its counterpart from 2013: smaller samples tend to produce more extreme results
376 (Huber, 2011). It is possible that with a smaller sample size of 13, we recruited a disproportionate
377 number of subjects who happened to have a high disposition toward an exclusively part-based template
378 (also see Reeder, 2017). Contreras Rubio, Peña, & Santacreu (2010) remarked that people tend to fall

379 on a continuum between relying more on whole objects and relying more on segments for visual
380 search. The likelihood to recruit a sample that is skewed to either side of that spectrum is bigger for a
381 sample size of 13 than it is for a sample size of 25. Another possibility is that parts can capture
382 attention to some extent because they are still a necessary part of the search template, so the difference
383 in the capture effects between silhouette parts and wholes is not as robust as, say, the difference
384 between texture patches and wholes. Larger samples, experiment replication, and task variation can
385 therefore all contribute to a more accurate assessment of the contribution of diagnostic parts to the
386 contents of the category-level template.

387 The results of Experiments 2 and 3 provide evidence that, when presented at search-irrelevant
388 locations, object parts do not even show a trend toward a capture effect. This supports the idea that
389 parts alone do not make an efficient template, and are perhaps easier to inhibit in the latter two
390 experiments. Our results are in contrast to Ullman et al. (2002), who showed computationally that parts
391 of moderate complexity should actually be an ideal candidate to identify objects at the category level.
392 However, the algorithms used in that study did not operate under the same limitations in time and
393 capacity as our very human subjects did. The futility of adding more object parts in Experiment 3 and
394 the trend towards capture effects for parts in Experiment 1 indicate further that it is not a lack of
395 diagnostic information, but rather inefficiency of template matching, that prevents object parts from
396 capturing attention.

397 One reason we initially hypothesized that the category-level template for real-world objects is
398 composed of parts rather than wholes, is because we found capture effects for whole object silhouettes
399 that were not only upright, but also inverted and rotated (Reeder & Peelen, 2013: Experiment 2 and 3).
400 Because a canonical viewpoint of object features is not a necessary aspect of the template, we proposed
401 that it would be easier to activate a spatially non-specific collection of parts rather than many variations
402 of whole objects in the template. However, there is another possibility: in the real world, we may
403 encounter whole object shapes that are inverted or rotated: just think of a car that has been turned

404 upside down in an accident, or a person performing a handstand or leaning against a wall. We can still
405 easily recognize these objects; however, the parts are always presented in a sensible relation to one
406 another – even if a car is upside down, the wheels are attached to the lower body, with the doors
407 appearing between the hood and the wheels. The chance of beholding a car with the doors on the
408 ground and the wheels above the hood is almost non-existent. It is more likely, as an adaptation to real
409 life probabilities, that proper spatial relations between parts (regardless of viewpoint) is a necessary
410 part of the template. Therefore, our results point toward a viewpoint- and location-flexible template that
411 nevertheless requires configural information to remain intact. One line of future research would be to
412 present collections of parts in meaningful layouts (e.g., by removing the torso from person images), and
413 to evaluate whether configural information without the body can benefit a capture effect by object
414 parts.

415 The fact that parts do not show a significant consistency effect in this study does not mean that
416 the search template can never be composed of diagnostic parts in isolation. Specific attributes of the
417 scenes we used might have enhanced the value of configural representations as a template. The search
418 was broad (cars and people appeared in all sorts of sizes, colors, positions), swift (search scenes were
419 presented for 67 ms), and included various levels of target occlusion (there are no deliberate
420 occlusions, and many scenes showed entire objects); these are all properties that might favor a template
421 composed of a coarse representation of the whole object shape, rather than the exclusive representation
422 of diagnostic parts. The latter might be a good template for a specific search that requires the
423 representation of more detailed part information (e.g., detecting sports cars). Parts in isolation would
424 also be a more appropriate template if a search required subjects to detect predominantly incomplete or
425 occluded targets. Furthermore, longer presentation times for search scenes would have allowed subjects
426 to attend smaller or richer details of objects, which could in turn bias the template more heavily toward
427 diagnostic parts. Finally, the results should not be generalized to other stimuli too readily—both people
428 and cars are special with regards to their commonality and the distinctiveness in the shapes of their

429 parts. More exotic or uncommon stimuli would likely require a more specific shape template for
430 detection. There are also many object categories for which shape itself is not diagnostic – one example
431 being the proverbial apples and oranges (for which color-based templates may be more diagnostic).

432

433 **5. Conclusions**

434

435 The results of this study indicate that the search template for object categories presented in natural
436 scenes is composed of whole object shape. Shapes of diagnostic object parts alone are not an adequate
437 template, contrary to previous findings.

438

439 **Author Contributions**

440 RRR developed and designed the experiment. MW and RRR both contributed to data collection, data
441 analysis and interpretation, and writing the manuscript. Both authors approved the final version of the
442 manuscript.

443

444 **Acknowledgements**

445 This research was funded by an Open Research Area grant [DFG PO 548/16-1].

446

447

448

449

450

451

452

453

454

455

456

457

458 **References**

459 Brandman, T., & Yovel, G. (2010). The body inversion effect is mediated by face-selective, not body-
460 selective, mechanisms. *Journal of Neuroscience*, *30*(31), 10534-10540.

461 <https://doi.org/10.1523/JNEUROSCI.0911-10.2010>

462 Bravo, M. J., & Farid, H. (2012). Task demands determine the specificity of the search template.

463 *Attention, perception & psychophysics*, *74*(1), 124–131. <https://doi.org/10.3758/s13414-011->

464 [0224-5](https://doi.org/10.3758/s13414-011-0224-5)

465 Chen, X., Mottaghi, R., Liu, X., Fidler, S., Urtasun, R., & Yuille, A. (2014). Detect what you can:

466 Detecting and representing objects using holistic models and body parts. In *Proceedings of the*
467 *IEEE Conference on Computer Vision and Pattern Recognition* (pp. 1971-1978).

468 <https://doi.org/10.1109/CVPR.2014.254>

469 Contreras, M. J., Rubio, V. J., Peña, D., & Santacreu, J. (2010). On the Robustness of Solution Strategy

470 Classifications. *Journal of Individual Differences*, *31*(2), 68–73. <https://doi.org/10.1027/1614->

471 [0001/a000012](https://doi.org/10.1027/1614-0001/a000012)

472 Delorme, A., Richard, G., & Fabre-Thorpe, M. (2010). Key visual features for rapid categorization of
473 animals in natural scenes. *Frontiers in psychology*, *1*, 21.

474 <https://doi.org/10.3389/fpsyg.2010.00021>

475 Desimone, R., & Duncan, J. (1995). Neural mechanisms of selective visual attention. *Annual review of*

476 *neuroscience*, *18*, 193–222. <https://doi.org/10.1146/annurev.ne.18.030195.001205>

477 Duncan, J., & Humphreys, G. W. (1989). Visual search and stimulus similarity. *Psychological Reviews*,

478 96(3), 433-458. <https://doi.org/10.1037/0033-295X.96.3.433>

479 Faul, F., Erdfelder, E., Buchner, A., & Lang, A.-G. (2009). Statistical power analyses using G*Power
480 3.1: Tests for correlation and regression analyses. *Behavior Research Methods*, 41, 1149-1160.
481 <https://doi.org/10.3758/BRM.41.4.1149>

482 Folk, C. L., Remington, R. W., & Johnston, J. C. (1992). Involuntary covert orienting is contingent on
483 attentional control settings. *Journal of Experimental Psychology: Human perception and*
484 *performance*, 18(4), 1030-1044. <https://doi.org/10.1037/0096-1523.18.4.1030>

485 Hasegawa, I., & Miyashita, Y. (2002). Categorizing the world: expert neurons look into key features.
486 *Nature neuroscience*, 5(2), 90–91. <https://doi.org/10.1038/nn0202-90>

487 Huber, P. J. (2011). Robust Statistics. In M. Lovric (Ed.), *International Encyclopedia of Statistical*
488 *Science* (pp.1248–1251). Berlin, Heidelberg: Springer Berlin Heidelberg.
489 https://doi.org/10.1007/978-3-642-04898-2_594

490 Lerner, Y., Hendler, T., Ben-Bashat, D., Harel, M., & Malach, R. (2001). A hierarchical axis of object
491 processing stages in the human visual cortex. *Cerebral cortex*, 11(4), 287-297.
492 <https://doi.org/10.1093/cercor/11.4.287>

493 Logothetis, N. K., & Sheinberg, D. L. (1996). Visual object recognition. *Annual review of*
494 *neuroscience*, 19(1), 577-621. <https://doi.org/10.1146/annurev.ne.19.030196.003045>

495 Malcolm, G. L., Groen, I. I., & Baker, C. I. (2016). Making sense of real-world scenes. *Trends in*
496 *cognitive sciences*, 20(11), 843-856. <https://doi.org/10.1016/j.tics.2016.09.003>

497 Malcolm, G. L., & Henderson, J. M. (2009). The effects of target template specificity on visual search
498 in real-world scenes: evidence from eye movements. *Journal of vision*, 9(11), 8.1-13.
499 <https://doi.org/10.1167/9.11.8>.

500 Malcolm, G. L., & Henderson, J. M. (2010). Combining top-down processes to guide eye movements

501 during real-world scene search. *Journal of vision*, 10(2), 4.1-11. <https://doi.org/10.1167/10.2.4>

502 Peelen, M. V., & Kastner, S. (2014). Attention in the real world: toward understanding its neural basis.
503 *Trends in cognitive sciences*, 18(5), 242–250. <https://doi.org/10.1016/j.tics.2014.02.004>

504 Peirce, J. W. (2008). Generating Stimuli for Neuroscience Using PsychoPy. *Frontiers in*
505 *neuroinformatics*, 2, 10. <https://doi.org/10.3389/neuro.11.010.2008>

506 Quinlan, P. T. (2003). Visual feature integration theory: Past, present, and future. *Psychological*
507 *bulletin*, 129(5), 643-673. <https://doi.org/10.1037/0033-2909.129.5.643>

508 Reed, C. L., Stone, V. E., Bozova, S., & Tanaka, J. (2003). The body-inversion effect. *Psychological*
509 *science*, 14(4), 302-308. <https://doi.org/10.1111/1467-9280.14431>

510 Reeder, R. R. (2017). Individual differences shape the content of visual representations. *Vision*
511 *research*, 141, 266-281. <https://doi.org/10.1016/j.visres.2016.08.008>

512 Reeder, R. R., Hanke, M., & Pollmann, S. (2017). Task relevance modulates the cortical representation
513 of feature conjunctions in the target template. *Scientific reports*, 7(1), 4514.
514 <https://doi.org/10.1038/s41598-017-04123-8>

515 Reeder, R. R., & Peelen, M. V. (2013). The contents of the search template for category-level search in
516 natural scenes. *Journal of vision*, 13(3). <https://doi.org/10.1167/13.3.13>

517 Reeder, R. R., van Zoest, W., & Peelen, M. V. (2015). Involuntary attentional capture by task-
518 irrelevant objects that match the search template for category detection in natural scenes. *Attention,*
519 *Perception, & Psychophysics*, 77(4), 1070–1080. <https://doi.org/10.3758/s13414-015-0867-8>

520 Seidl-Rathkopf, K. N., Turk-Browne, N. B., & Kastner, S. (2015). Automatic guidance of attention
521 during real-world visual search. *Attention, Perception, & Psychophysics*, 77(6), 1881-1895.
522 <https://doi.org/10.3758/s13414-015-0903-8>

523 Singh, M., & Hoffman, D. D. (2001). Part-based representations of visual shape and implications for

524 visual cognition. In *Advances in psychology* (Vol. 130, pp. 401-459). North-Holland.
525 [https://doi.org/10.1016/S0166-4115\(01\)80033-9](https://doi.org/10.1016/S0166-4115(01)80033-9)

526 Thorpe, S., Fize, D., & Marlot, C. (1996). Speed of processing in the human visual system. *Nature*,
527 *381*(6582), 520-522. <https://doi.org/10.1038/381520a0>

528 Treisman, A., & Gelade, G. (1980). A feature integration theory of attention. *Cognitive Psychology*, *12*,
529 97–136. [https://doi.org/10.1016/0010-0285\(80\)90005-5](https://doi.org/10.1016/0010-0285(80)90005-5)

530 Ullman, S., Vidal-Naquet, M., & Sali, E. (2002). Visual features of intermediate complexity and their
531 use in classification. *Nature neuroscience*, *5*(7), 682–687. <https://doi.org/10.1038/nn870>

532 Wolfe, J. M. (1994). Guided Search 2.0 A revised model of visual search. *Psychonomic Bulletin &*
533 *Review*, *1*(2), 202–238. <https://doi.org/10.3758/BF03200774>

534 Wyble, B., Folk, C., & Potter, M. C. (2013). Contingent attentional capture by conceptually relevant
535 images. *Journal of Experimental Psychology: Human Perception and Performance*, *39*(3), 861-
536 871. <https://doi.org/10.1037/a0030517>

537 Yovel, G., Pelc, T., & Lubetzky, I. (2010). It's all in your head: why is the body inversion effect
538 abolished for headless bodies?. *Journal of Experimental Psychology: Human Perception and*
539 *Performance*, *36*(3), 759-767. <https://doi.org/10.1037/a0017451>

540 Zelinsky, G. J. (2008). A theory of eye movements during target acquisition. *Psychological review*,
541 *115*(4), 787–835. <https://doi.org/10.1037/a0013118>