

What is “where”: Physical reasoning informs object location

Tal Boger^{a,*}, Tomer Ullman^b

^a*Department of Psychology, Yale University, New Haven, CT 06520*

^b*Department of Psychology, Harvard University, Cambridge, MA 02138*

Abstract

A central puzzle the visual system tries to solve is: “what is where?” While a great deal of research attempts to model object recognition (“what”), a comparatively smaller body of work seeks to model object location (“where”), especially in perceiving everyday objects. How do people locate an object, right now, in front of them? In three experiments collecting over 35,000 judgements on stimuli spanning different levels of realism (line drawings, real images, and crude forms), participants clicked “where” an object is, as if pointing to it. We modeled their responses with eight different methods, including both human response-based models (judgements of physical reasoning, spatial memory, free-response “click anywhere” judgements, and judgements of where people would grab the object), and image-based models (uniform distributions over the image, convex hull, saliency map, and medial axis). Physical reasoning was the best predictor of “where,” performing significantly better than even spatial memory and free-response judgements. Our results offer insight into the perception of object locations while also raising interesting questions about the relationship between physical reasoning and visual perception.

Keywords: Object representation, Perception, Physical reasoning

*To whom correspondence should be addressed.

Email addresses: tal.boger@yale.edu (Tal Boger), tomerullman@gmail.com (Tomer Ullman)

1 Introduction

2 When asked “where is the hammer,” with a hammer right in front of you, where
3 would you point? Initially, this question seems trivial; the hammer is “over there.”
4 Yet, one can give many reasonable answers, each highlighting different properties.
5 Perhaps you would point to the center of the hammer (as it appears to you); or to
6 its handle (the part where you would hold it); or to its metal head (the part that
7 performs the action); or to other locations still. The question “where?” points to a
8 subtle problem in our classic definitions of vision.

9 Early definitions of vision distilled the complex process of seeing into a simple
10 question: “what is where?” (Marr, 1982). Such definitions served as touchstones for
11 exploring vision in philosophy, cognitive science, and neuroscience, where researchers
12 discovered an apparent split in the visual system between the ventral stream — which
13 models “what” — and the dorsal stream — which models “where” (Schneider, 1969;
14 though more recent work has significantly complicated this initially neat split, as
15 discussed later).

16 Plenty of research in visual cognition has focused on modeling “what,” and there
17 is an expansive literature about the mechanisms underlying object recognition. While
18 there is also an expansive literature on “where,” by relative comparison it has been
19 less explored than “what,” especially in the perception of everyday objects. Here,
20 we take a step towards exploring the nature of object location by asking: what is
21 “where”?

22 Much of the existing work on modeling “where” analyzes processes different from
23 simply perceiving objects as they appear in front of us. For example, various work
24 explores the nature of object location via spatial memory (Langlois et al., 2021),
25 ambiguous shapes (Huttenlocher et al., 1991), object parts and scenes (Bar and
26 Ullman, 1996), or eye movements (Vishwanath and Kowler, 2003). These all inform
27 our understanding of object localization in the mind and use methods similar to ours,
28 though in a different context. For example, these works only give hints to where we
29 may point at a hammer if it appeared right in front of us – but do not give a well-
30 defined answer. We expand on these works by testing the nature of perceived object
31 location in simple tasks with everyday objects at differing levels, revealing aspects
32 of object location in our daily lives.

33 Here, we present three experiments collecting data from over 35,000 judgements
34 in which participants indicate “where” an object is. The experiments use objects
35 covering a wide range of information and realism (such that they generalize to a
36 range of stimuli). We modeled “where” using methods based on previous work,
37 including both human response-based models and image-based models. Across all

38 levels of realism, a model that relies on physical reasoning (perceived center-of-mass)¹
39 best predicted “where” an object is. Our results provide novel insights into how we
40 model object locations in perception, and point to a surprising relationship between
41 physical reasoning and visual perception.

42 Results

43 Our three experiments span a range of object realism. In Experiment 1, we used
44 line drawings with no depth and color. In Experiment 2, we used images of real
45 objects with depth and color (but no background). Finally, in Experiment 3, we
46 masked and rotated the line drawings, such that they became unidentifiable crude
47 forms.

48 On each trial, participants clicked “where” each object is, as if pointing it out to
49 another person (Figure 1). Each stimulus set consisted of 50 objects, which included
50 a range of everyday entities, both symmetric and asymmetric items, tools, agents,
51 and more.

52 We modeled participant responses for “where” using eight different models. The
53 first three models were based on human responses from other tasks, collected sepa-
54 rately: (1) center of mass (“click on the object’s center of mass”), (2) spatial memory
55 (“click where the object was” after the object disappeared), and (3) free-response
56 clicks as an attention proxy (“click anywhere on the object”). Participants in each
57 task and experiment were unique and independent. We also considered four image-
58 based models: (4) a uniform distribution across the object, (5) a uniform distribu-
59 tion across the object’s convex hull, (6) the object’s saliency map (as generated by
60 OpenCV fine-grained saliency maps), and (7) medial axis (as generated by scikit-
61 image). The models provide a balance between new proposals specific to this work,
62 and existing models that have been shown to perform well in similar tasks (e.g.,
63 medial axis from Firestone and Scholl, 2014). After we tested these broad models
64 of “where,” we pre-registered and analyzed a final, more specific model: (8) human
65 responses on where they would grasp the object to pick it up.

66 With regards to the center-of-mass model, we emphasize that the true center of
67 mass cannot be accurately recovered, and is also irrelevant even if it could be, as
68 people have no direct access to it. The primary aspect that matters for our analysis

¹Note that, while other relevant dimensions for physical reasoning in humans exist beyond center of mass, we use center of mass as a proxy for physical reasoning. Computing center of mass requires some kind of physical reasoning, which previous work has shown to be quite sensitive, or at least inaccurate in consistent ways which still imply a physical computation (Cholewiak et al., 2013, 2015; Firestone and Keil, 2016), making such a proxy reasonable and well-defined.



Figure 1: Example participant data from Experiment 1 (line drawings). Each blue dot shows a participant’s click in response to the query “where is the [object]?” Plots of “where” data for all our stimuli are accessible on our OSF repository (osf.io/nhj7k). Readers may also try each task for themselves at: tb.perceptionresearch.org/what_is_where.

69 based on this model is people’s subjective judgement of the center of mass, and how
 70 that relates to the perception of “where.”

71 To test the performance of our models, we first fit a Gaussian mixture model
 72 (GMM) to the “where” data provided by participants for each object, such that
 73 we could compare distributions of participant responses (Figure 2). The number of
 74 mixtures was chosen via three-fold cross-validation (for between 1 and 5 mixtures).
 75 We then calculated the mean negative log-likelihood of our models under this GMM
 76 for “where” on each object. Finally, we compared the models using paired Wilcoxon
 77 signed-rank tests on the mean negative log-likelihood scores to ask which model best
 78 predicted the “where” responses. All analysis plans, choice of models, and experi-
 79 mental designs were pre-registered; materials and data are available at osf.io/nhj7k.

80 In all three experiments, “where” responses were best predicted by center-of-
 81 mass judgements, followed by spatial memory judgements and free-response clicks,

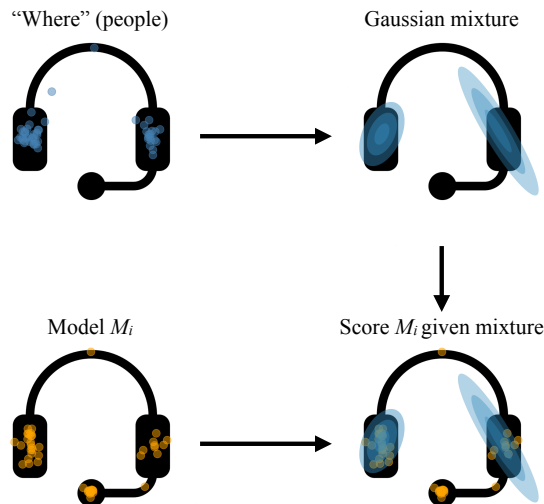


Figure 2: Schematic illustration of our modeling paradigm, using a line drawing of headphones as an example. First, we collect participant judgements for “where” an object is. After fitting a Gaussian mixture to this data, we score a proposed model M_i under this Gaussian mixture, using its negative log-likelihood.

82 respectively (Figure 3). The differences between these models was significant: physical
 83 reasoning was a significantly better predictor of “where” than spatial memory
 84 across experiments (Experiment 1: $p < 0.01$; Experiment 2: $p < 0.001$; Experiment
 85 3: $p < 0.001$). Spatial memory in turn significantly outperformed the free-response
 86 model, though the difference was slightly smaller (Experiment 1: $p = 0.01$; Experiment
 87 2: $p = 0.04$; Experiment 3: $p < 0.01$). The various image-based models, while
 88 based on previous work and reasonable assumptions, performed poorly by compar-
 89 ison.

90 Furthermore, our eighth model (where people would grab the object), which we
 91 pre-registered and explored after testing our initial seven, performed significantly
 92 worse than our initial three human response-based models. In all three experiments,
 93 the grasping model performed the worse than the center of mass, spatial memory,
 94 and “click anywhere” models ($p < 0.001$ when compared to the center-of-mass model
 95 in each experiment).

96 We estimate ceiling performance as the log likelihood of people’s “where” judge-
 97 ments under its own GMM – another model should not predict the “where” data
 98 better than the “where” data itself. In Experiments 1 and 2, the “where” data was
 99 significantly more likely under the GMM than the center of mass data (Experiment

100 1: $p < 0.01$; Experiment 2: $p < 0.01$). However, in Experiment 3, we observe
 101 near-ceiling performance for the physical reasoning model, as its likelihood is not
 102 distinguishable from the likelihood of the “where” data itself ($p = 0.43$).

103 Though participants tend to click near the center of objects — perhaps leading
 104 to a bias towards the center-of-mass model — this does not explain away our results.
 105 First, this bias of clicking near the center of the objects would apply to *all* human-
 106 response models, not just the center-of-mass model. Further, had this been the case,
 107 then simply predicting responses for “where” by distance from the image centroid
 108 would be the best model. However, this was not the case, suggesting that this
 109 preference for center of mass goes beyond mere clicking biases.

110 Discussion

111 What is “where”? Our three experiments explore how we judge the location
 112 of objects and find that, across a range of object realism, a judgement rooted in
 113 physical reasoning (center of mass) is the strongest predictor of perceived object
 114 location. We suggest that judgements of object location rely on physical properties.
 115 This idea echoes other work in visual perception, neuroscience, and developmental
 116 psychology.

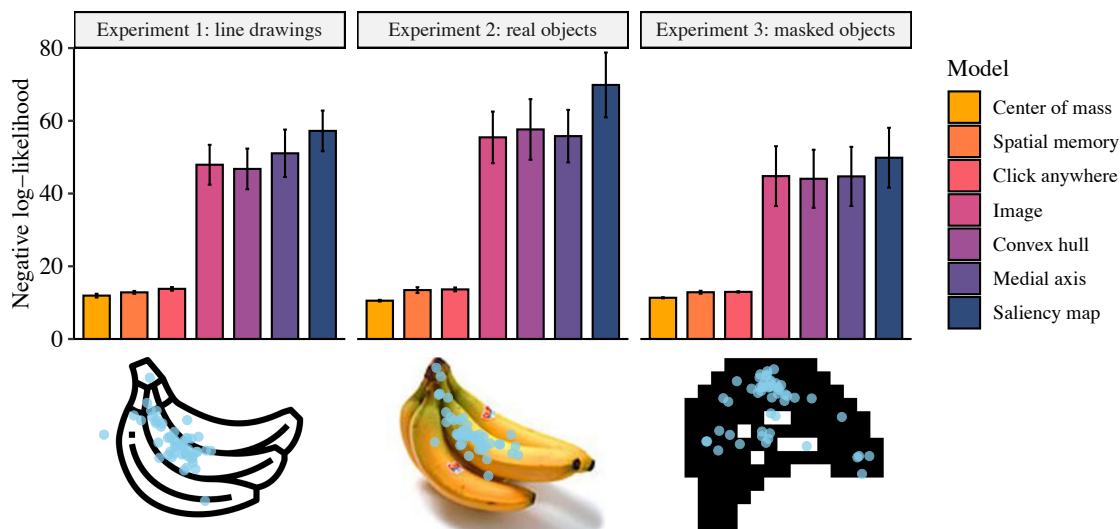


Figure 3: Mean negative log-likelihood of our initial seven models in each experiment. In all three experiments, the center-of-mass model performed significantly better than the other models. Under each experiment’s results is a depiction of a banana in that experiment’s form of stimuli, overlaid with “where” data. All analyses and stimuli are available at (osf.io/nhj7k).

117 Recently, many researchers hypothesized that a mental “simulation engine” un-
118 derlies much of our intuitive physical reasoning (Battaglia et al., 2013; Fischer et al.,
119 2016). The first step in these mental simulation models assumes we de-render a visual
120 image into a physical scene representation, but this process is not yet formally solved.
121 With a few exceptions (e.g., Little and Firestone, 2021; Boger and Firestone, 2022),
122 the question of how physical reasoning fits into vision remains relatively unexplored.

123 Our results speak to a different direction in this process; we suggest that physical
124 reasoning and vision rely on each other, rather than one process exclusively relying
125 on the other. This relationship is strong enough that physical reasoning predicts
126 perceived object locations better than seemingly closer processes such as spatial
127 memory and attention. The instructions for the spatial memory judgements (“where
128 was the object”), free-response judgements (“click anywhere on the object”), and
129 “where” judgements (“where is the object”) are all semantically similar, compared
130 to the physical reasoning judgements (“click on the object’s center of mass”). This
131 makes our empirical findings all the more striking.

132 Beyond intuitive physical reasoning, our work adds a new angle to existing liter-
133 ature on multiple object tracking (MOT) and object-based attention. For example,
134 foundational work in MOT suggests that attention is deployed to objects rather than
135 features, implying that there is a sort of physical “objectness” crucial to vision (Scholl
136 et al., 2001). This tracking ability persists dynamically through some physical events
137 such as occlusion, but not others (such as deletion) (Scholl and Pylyshyn, 1999). Such
138 object-based attention even exists over object representations implicitly created by
139 perceptual completion (Moore et al., 1998). However, much of these results exist
140 over either dynamic objects (i.e., in the case of MOT) or abstract objects. Here, we
141 show the impact of physical “objectness” on the judgements of object location in a
142 simple, static paradigm involving everyday objects.

143 Our results also have analogs in neuroscience, further aligning with our proposed
144 relationship between physical reasoning and vision. For example, the MT complex
145 — thought to be responsible for motion perception — has been shown to mediate
146 attentive tracking (Culham et al., 1998). A large body of work has taken this classic
147 split between the “what” and “where” streams to include “how,” which modulates
148 how we interact with objects, noting that people can direct accurate motor move-
149 ments at objects they fail to localize (Goodale et al., 1991; Goodale and Milner, 1992;
150 Kravitz et al., 2011). This matches our results that the “where” of objects may be
151 constrained by their physical behaviors.

152 At an even more basic level of representation, foundational work in infant cogni-
153 tion shows that physical reasoning may supersede representations of object identity.
154 Even when infants forgot the features of a set of objects, they still expected them

155 to remain physically consistent. For example, infants who fail to notice changes in
156 object shape are surprised to see the object disappear entirely (Kibbe and Leslie,
157 2011; Zosh and Feigenson, 2012). This provides an even richer demonstration of
158 earlier work showing that infants’ representations of objects conform to more foun-
159 dational physical properties such as continuity and rigidity (Baillargeon et al., 1985;
160 Spelke and Van de Walle, 1993; Spelke et al., 1992). We suggest that such repre-
161 sentations persist across development (as in Kibbe, 2015) in even simpler ways than
162 object tracking or planning, and rather influence our judgements of object locations
163 in simple, static settings.

164 Though our experiments covered a range of realism, they were still limited to
165 static objects with no background. Physical reasoning in the real world deals with
166 complex scenes and moving objects. Future work may explore how dynamics, re-
167 lations, and interactions affect judgements of “where.” In everyday life, we do not
168 recognize and locate a hammer as an isolated object, but as a hammer *next to* a
169 cup, *behind* a book, *on* a tiger, and so on (Hafri and Firestone, 2021). These re-
170 lations require extracting rich physical and visual information that may affect our
171 perception of object location. For example, participating in a causal event such as a
172 collision creates a reliable illusion in the spatial relations of two objects (Scholl and
173 Nakayama, 2004). Future work may explore how “where” changes if, rather than
174 seeing a hammer with no background, we see a hammer supported by a table from
175 below, or supported by a string from above.

176 Beyond the insights our results provide about physical reasoning and object lo-
177 cation, they raise intriguing questions about how “where” relates to “what.” For
178 example, when judging the location of a pineapple, people consistently clicked on its
179 body, as if ignoring the stem on top. What does this mean about the nature of how
180 we perceive the pineapple? Judgements of “where” may reveal a unique way to an-
181 alyze the perceived essence of an object (Gelman et al., 2003). Because participants
182 must choose only a single point on the object, they may ask themselves which part of
183 the object most represents its essence. Perhaps, in line with work in infant cognition
184 and MOT, we represent the essence of objects not only by their visual features, but
185 also by their physical properties, in some cases even relying more on the latter than
186 the former.

187 While our experiments ask each participant for a single judgement, our percep-
188 tion of “where” likely depends on more than just a single point. This representation
189 may instead resemble a “point cloud.” However, we can treat each single-point esti-
190 mate as a sample drawn from such a cloud distribution, in line with other proposals
191 on sampling-based cognition (Vul et al., 2014). By aggregating judgements across
192 participants, we generate well-powered cumulative distributions of object location,

193 in much the same way that cumulative distributions reveal mathematically defined
194 shape skeletons (Firestone and Scholl, 2014). In this sense, our distributions recap-
195 ture the potential point-cloud distributions, which turn out to be neither uniform,
196 nor skeletal.²

197 By modeling “where” in a simple and direct setting, we take a step towards
198 understanding how we represent object location. Our results reveal a surprising
199 bidirectional relationship between physical reasoning and visual perception. More
200 broadly, we believe this work suggests a novel avenue for future work on modeling
201 “where” in vision. While it does not fully resolve the question of what is “where,”
202 it suggests where to look.

203 **Frequently asked questions**

204 We thought it would be useful to directly address a few common questions and
205 comments we have received regarding this work. We hope this helps lead to open
206 conversation with readers, and serves as a “theoretical supplement.”

207 *Surely people represent locations as more than a single dot, something more like an*
208 *area-cloud?*

209 We agree participant representations of “where” may involve multiple locations
210 on the object, rather than the single location we ask each participant to produce.
211 But, we believe these single point estimates form well-powered point clouds, which
212 together reflect the “where” distribution.

213 In many ways, giving participants the option to click on the object is the best
214 solution to such a problem. First, such issues exist in other single-choice clicking tasks
215 (e.g., Firestone and Scholl, 2014) which also find mathematically strong distributions.
216 Second, clicking tasks give participants maximum flexibility to represent this single
217 location as best they can, whereas, for example, a forced-choice task adds ambiguity
218 to this point-cloud across participants.

²An additional way to address this point cloud hypothesis would be with a series of object localization tasks that do not rely on clicking, such as a vernier acuity task (for review, see McKee and Westhe, 1978). Relying on two-alternative forced-choice responses for object positions eliminates aspects of the fine-grained modeling approach we present here. However, it also presents a higher-level, coarser interpretation of object location; future research may seek to use these types of paradigms to further our understanding of object localization. In this work, we stick to clicking-based tasks given their simplicity and prevalence in related work, such as in Firestone and Scholl, 2014.

219 *Wouldn't any dot on the image be a valid answer? If I point to any part of an image*
220 *and ask "is this the [object]" the answer should be "yes."*

221 In principle, valid "where" responses would be any location on the surface or edges
222 of the object (though in our experiments people also point to empty areas, such as
223 the middle of a bicycle). In practice, this is not what people do when generating
224 responses. Under the hypothesis that any image part is a valid answer, people would
225 conflate "where" with "any non-background pixel," and responses should then either
226 form a uniform distribution across the object (a uniform point cloud), or the single-
227 best error-minimizing point sample from that cloud (the image center). We don't
228 observe either of these. Rather, we see that center of mass is highly predictive of
229 "where" responses.

230 *Is this about vision? Isn't this actually about social things, such as communication?*

231 We do not know for sure that our results contain no social component, and it
232 would be interesting if they did. However, we cannot think of a theoretical account at
233 the moment for how social features explain our results, and dictate people's responses
234 in a way that a hypothetical "social-free" version would not. Put as a question, why
235 would "point an object out to someone" cause participants to produce clicks that
236 match judgements about center of mass (a non-social judgement), but simply locating
237 an object for yourself result in different judgements?

238 Also, such social or communication components exist (via task demands) in other
239 studies that are taken to be about vision, and cannot be fully removed. For example,
240 (Firestone and Scholl, 2014) — who ran similar tasks to explore shape skeletons in
241 the visual system — ask participants to tap anywhere on a shape. Though there is no
242 language about "pointing it out to someone," the experiments do require participants
243 to tap the shape on an iPad held by an experimenter, requiring some form of pointing
244 it out to, and communicating with the experimenter.

245 More broadly, the question about whether this is *actually* about vision in turn
246 raises the question of what vision is, and a classic answer has been "vision is about
247 what is where," bringing us full circle.

248 *Have you considered [this other model] instead?*

249 In this work, we analyzed eight different models, which is straining a short paper.
250 We chose these models to form an encompassing package, while trying to not be
251 overbearing, and not claiming to be exhaustive. In the process, it's quite possible we
252 left out other reasonable models.

253 We're happy to explore new models, or additions to the current models. We also
254 encourage proposals for why our existing models work or don't work. We believe

255 part of the appeal of this work is in spurring new directions. However, we have two
256 suggestions for any new models or proposals.

257 First, new models or proposals should match the data already at hand, at a basic
258 level. For example, several proposals beyond our set turn out to be equivalent to
259 a uniform distribution or center-of-image model, which does not match the existing
260 data. We considered this above, in the interpretation that “where” ambiguously
261 leads to “any non-background pixel,” and we’ve also come across proposals that
262 people might “minimize the error of a mis-click,” which turn out to be similar.

263 Second, new models or proposals should be able to generalize in a way that can
264 capture both our broad stimulus set, and visual representations more generally. This
265 is perhaps a main pitfall of our eighth model, that asks participants where they would
266 grab the object. Many objects in our stimulus set (and the world) are not graspable,
267 especially the crude forms we use in Experiment 3.

268 We invite interested readers to test new models and proposals; all our data and
269 experimental code are available on our OSF repository (osf.io/nhj7k).

270 *How do we know participants are calculating the center of mass accurately? Why not*
271 *calculate the true center of mass, instead of relying on people’s judgements?*

272 We would stress that a “true” center of mass cannot be calculated from our
273 images, given that the weight of each object part is unknown. So, there is no way
274 to know such calculations are capturing a ground truth. (Though previous work has
275 shown that people accurately judge an object’s center of mass ([Cholewiak et al., 2015](#))).
276 More importantly, even if we *could* calculate the true center of mass, it
277 would be irrelevant for judgements of “where.” People do not have access to the
278 ground-truth center of mass beyond the mental calculations they perform in the task
279 asking them to estimate the center of mass, which is what we asked them to do.

280 *How do we know participants are not merely clicking on the center of the object for*
281 *“where,” and that’s why center of mass is the best model?*

282 As with the above question of additional models, we believe that this concern
283 would need to be first validated by the data. Before performing any analysis, we can
284 see that people are not merely clicking on the object’s center, and rather that the
285 clicks possess a unique distribution which seems to have some structure.

286 However, this concern can also be tested empirically; if the main reason the center-
287 of-mass model predicts the “where” data well is because of a bias to click towards
288 the center, then a model predicting “where” clicks using the image centroid should
289 perform the strongest. However, this is not the case, as it performs significantly
290 worse than all human response-based model.

291 Finally, if such a center bias existed in the “where” clicks, it would likely extend
292 to other models. It is especially hard to explain why such a bias would not extend
293 to the “click anywhere” model (and why that model is not the strongest) under this
294 explanation; the instructions for the “click anywhere” and “where” tasks are almost
295 identical, such that a center bias in one *should* extend to the other if it existed.
296 However, this is not what we observe, so we believe this concern is unsubstantiated
297 by our data.

298 **Materials and methods**

299 *Participants*

300 Each of the three experiments recruited 50 unique participants for each of the
301 tasks and each of the three forms of stimuli (“click where the object is,” “click on
302 the object’s center of mass,” “click where the object was,” “click anywhere on the
303 object,” and “click where you would grab the object to pick it up”; total N=750).
304 All participants were recruited from the online platform Prolific (for a discussion of
305 the reliability of this subject pool, see [Peer et al., 2017](#)). Unique participants were
306 used for each condition and experiment such that no participant appeared in more
307 than one model or in both the dependent and independent variables. Participants
308 were excluded if they did not contribute a complete dataset or if they clicked the
309 same location in five consecutive trials.

310 *Stimuli*

311 Line drawings for Experiment 1 were taken from The Noun Project. Object
312 images for Experiment 2 were taken from a variety of online sources. The kinds of
313 objects in Experiment 2 were the same as those in Experiment 1 (i.e., if Experiment
314 1 included a line drawing of a gorilla, Experiment 2 included a real image of a
315 gorilla). The masked objects for Experiment 3 were created by applying a random
316 mask to the line drawings, then vertically flipping them to remove any identifying
317 information. The images were randomly padded both vertically and horizontally such
318 that responding in the center of the screen each time would not produce reasonable
319 data. All images were 500x500 pixels large in the participant’s web browser.

320 Note that unique participants are assigned to each condition, where they then
321 see all the stimuli in the given form and answer the given question. In other words, a
322 participant in the Experiment 1 “center of mass” conditions will see 50 line drawings
323 and click on their center of mass; they will not see any images of other types or be
324 told to click according to different instructions. The same set of 50 images are used
325 across all conditions in a given experiment.

326 *Design and procedure*

327 Participants saw 50 images in each experiment. The order of the images was
328 randomized. When gathering judgements for “where,” we instructed participants
329 as follows: “Your friend asks: ‘where is the [object]?’. Click on where you would
330 point to.” In the “center of mass” condition, participants were told to “Click on the
331 center of mass of the [object].” Participants in this condition were provided with
332 an additional instruction of what center of mass means (“average position of all the
333 mass in the object”). Though we cannot calculate the “accuracy” of these responses
334 (given that we cannot calculate a true center of mass from images), the responses
335 appear consistent and reasonable (and previous work shows such judgements are
336 fairly consistent; Cholewiak et al., 2015). In the spatial memory condition, the
337 object appeared for 1000ms, during which time the participant’s mouse was hidden
338 and immovable. The object then disappeared and participants were instructed as
339 follows: “Your friend asks: ‘where was the [object]?’. Click on where you would point
340 to.” In the “click anywhere” condition, participants were told to “Click anywhere
341 you want on the [object].” Finally, the “grasp” condition instructed participants to
342 “Click where you would grab the [object] to pick it up.”

343 *Data availability*

344 All data, code, materials, and pre-registrations are available at osf.io/nhj7k.
345 Readers can also do the tasks for themselves at tb.perceptionresearch.org/what_is_where.

346 **Supportive Information**

347 *Acknowledgements*

348 For helpful discussion and comments on previous drafts, we thank Sami Yousif,
349 Chaz Firestone, members of the Harvard CoCoDev lab, and members of the Com-
350 putation and Language Lab at UC Berkeley. Funding: TU is supported by NSF
351 Science Technology Center Award CCF-1231216, the DARPA Machine Common
352 Sense program, and the Jacobs Foundation.

353 *Author contributions*

354 TB carried out the experiments, analyzed the data, and wrote the first draft of
355 the manuscript. TB and TU jointly edited the paper and designed the research

356 **References**

- 357 Baillargeon, R., Spelke, E. S., and Wasserman, S. (1985). Object permanence in
358 five-month-old infants. *Cognition*, 20(3):191–208.
- 359 Bar, M. and Ullman, S. (1996). Spatial context in recognition. *Perception*, 25(3):343–
360 352.
- 361 Battaglia, P. W., Hamrick, J. B., and Tenenbaum, J. B. (2013). Simulation as an
362 engine of physical scene understanding. *Proceedings of the National Academy of
363 Sciences*, 110(45):18327–18332.
- 364 Boger, T. and Firestone, C. (2022). Automatic simulation of unseen physical events.
365 *Journal of Vision*, 22(14):3637–3637.
- 366 Cholewiak, S. A., Fleming, R. W., and Singh, M. (2013). Visual perception of
367 the physical stability of asymmetric three-dimensional objects. *Journal of vision*,
368 13(4):12–12.
- 369 Cholewiak, S. A., Fleming, R. W., and Singh, M. (2015). Perception of physical
370 stability and center of mass of 3-d objects. *Journal of vision*, 15(2):13–13.
- 371 Culham, J. C., Brandt, S. A., Cavanagh, P., Kanwisher, N. G., Dale, A. M., and
372 Tootell, R. B. (1998). Cortical fmri activation produced by attentive tracking of
373 moving targets. *Journal of neurophysiology*.
- 374 Firestone, C. and Keil, F. C. (2016). Seeing the tipping point: Balance perception
375 and visual shape. *Journal of Experimental Psychology: General*, 145(7):872.
- 376 Firestone, C. and Scholl, B. J. (2014). “please tap the shape, anywhere you like”
377 shape skeletons in human vision revealed by an exceedingly simple measure. *Psy-
378 chological science*, 25(2):377–386.
- 379 Fischer, J., Mikhael, J. G., Tenenbaum, J. B., and Kanwisher, N. (2016). Functional
380 neuroanatomy of intuitive physical inference. *Proceedings of the national academy
381 of sciences*, 113(34):E5072–E5081.
- 382 Gelman, S. A. et al. (2003). *The essential child: Origins of essentialism in everyday
383 thought*. Oxford Cognitive Development.
- 384 Goodale, M. A. and Milner, A. D. (1992). Separate visual pathways for perception
385 and action. *Trends in neurosciences*, 15(1):20–25.

- 386 Goodale, M. A., Milner, A. D., Jakobson, L., and Carey, D. (1991). A neurological
387 dissociation between perceiving objects and grasping them. *Nature*, 349(6305):154–
388 156.
- 389 Hafri, A. and Firestone, C. (2021). The perception of relations. *Trends in Cognitive*
390 *Sciences*.
- 391 Huttenlocher, J., Hedges, L. V., and Duncan, S. (1991). Categories and particulars:
392 prototype effects in estimating spatial location. *Psychological review*, 98(3):352.
- 393 Kibbe, M. M. (2015). Varieties of visual working memory representation in infancy
394 and beyond. *Current Directions in Psychological Science*, 24(6):433–439.
- 395 Kibbe, M. M. and Leslie, A. M. (2011). What do infants remember when they forget?
396 location and identity in 6-month-olds’ memory for objects. *Psychological science*,
397 22(12):1500–1505.
- 398 Kravitz, D. J., Saleem, K. S., Baker, C. I., and Mishkin, M. (2011). A new neural
399 framework for visuospatial processing. *Nature Reviews Neuroscience*, 12(4):217–
400 230.
- 401 Langlois, T. A., Jacoby, N., Suchow, J. W., and Griffiths, T. L. (2021). Serial
402 reproduction reveals the geometry of visuospatial representations. *Proceedings of*
403 *the National Academy of Sciences*, 118(13).
- 404 Little, P. C. and Firestone, C. (2021). Physically implied surfaces. *Psychological*
405 *Science*, 32(5):799–808.
- 406 Marr, D. (1982). *Vision: A Computational Investigation into the Human Represen-*
407 *tation and Processing of Visual Information*. Henry Holt and Co., Inc., USA.
- 408 McKee, S. P. and Westhe, G. (1978). Improvement in vernier acuity with practice.
409 *Perception & psychophysics*, 24(3):258–262.
- 410 Moore, C. M., Yantis, S., and Vaughan, B. (1998). Object-based visual selection:
411 Evidence from perceptual completion. *Psychological science*, 9(2):104–110.
- 412 Peer, E., Brandimarte, L., Samat, S., and Acquisti, A. (2017). Beyond the turk:
413 Alternative platforms for crowdsourcing behavioral research. *Journal of Experi-*
414 *mental Social Psychology*, 70:153–163.

- 415 Schneider, G. E. (1969). Two visual systems: Brain mechanisms for localiza-
416 tion and discrimination are dissociated by tectal and cortical lesions. *Science*,
417 163(3870):895–902.
- 418 Scholl, B. J. and Nakayama, K. (2004). Illusory causal crescents: Misperceived
419 spatial relations due to perceived causality. *Perception*, 33(4):455–469.
- 420 Scholl, B. J. and Pylyshyn, Z. W. (1999). Tracking multiple items through occlusion:
421 Clues to visual objecthood. *Cognitive psychology*, 38(2):259–290.
- 422 Scholl, B. J., Pylyshyn, Z. W., and Feldman, J. (2001). What is a visual object?
423 evidence from target merging in multiple object tracking. *Cognition*, 80(1-2):159–
424 177.
- 425 Spelke, E. S., Breinlinger, K., Macomber, J., and Jacobson, K. (1992). Origins of
426 knowledge. *Psychological review*, 99(4):605.
- 427 Spelke, E. S. and Van de Walle, G. (1993). Perceiving and reasoning about ob-
428 jects: Insights from infants. *Spatial representation: Problems in philosophy and*
429 *psychology*, pages 132–161.
- 430 Vishwanath, D. and Kowler, E. (2003). Localization of shapes: Eye movements and
431 perception compared. *Vision Research*, 43(15):1637–1653.
- 432 Vul, E., Goodman, N., Griffiths, T. L., and Tenenbaum, J. B. (2014). One and done?
433 optimal decisions from very few samples. *Cognitive science*, 38(4):599–637.
- 434 Zosh, J. M. and Feigenson, L. (2012). Memory load affects object individuation in
435 18-month-old infants. *Journal of Experimental Child Psychology*, 113(3):322–336.