

# 1 Title

2 The Impact of Scene Context on Visual Object Recognition: Comparing Humans,  
3 Monkeys, and Computational Models

## 4 **Abbreviated title**

5 Probing primate visual context processing

## 6 **Authors**

7 Sara Djambazovska<sup>1,2</sup>, Anaa Zafer<sup>1,‡</sup>, Hamidreza Ramezanpour<sup>1,‡</sup>, Gabriel Kreiman<sup>2</sup>,  
8 and Kohitij Kar<sup>1</sup>

## 9 **Affiliation**

10

11 1. York University, Department of Biology and Centre for Vision Research,  
12 Toronto, Canada

13 2. Children's Hospital, Harvard Medical School, MA, USA

14

15

16 \* Correspondence should be addressed to Kohitij Kar

17 ‡ Contributed equally to this work

18

19 E-mail: [k0h1t1j@yorku.ca](mailto:k0h1t1j@yorku.ca)

## 20 **Conflict of interests**

21 The authors declare no competing financial interests.

## 22 **Acknowledgments**

23

24 SD was supported by the Bertarelli Foundation Fellowship grant. HR is funded by CIHR  
25 Postdoctoral Fellowship. GK is supported by NIH R01EY026025. KK has been supported by  
26 funds from the Canada Foundation for Innovation (CFI), the Canada Research Chair Program,  
27 the Simons Foundation Autism Research Initiative (SFARI, 967073), the Canada First Research  
28 Excellence Funds (VISTA Program), and a Google Research Award.

## 29 Abstract

30  
31  
32  
33

34 During natural vision, we rarely see objects in isolation but rather embedded in rich and complex  
35 contexts. Understanding how the brain recognizes objects in natural scenes by integrating  
36 contextual information remains a key challenge. To elucidate neural mechanisms compatible with  
37 human visual processing, we need an animal model that behaves similarly to humans, so that  
38 inferred neural mechanisms can provide hypotheses relevant to the human brain. Here we  
39 assessed whether rhesus macaques could model human context-driven object recognition by  
40 quantifying visual object identification abilities across variations in the amount, quality, and  
41 congruency of contextual cues. Behavioral metrics revealed strikingly similar context-dependent  
42 patterns between humans and monkeys. However, neural responses in the inferior temporal (IT)  
43 cortex of monkeys that were never explicitly trained to discriminate objects in context, as well as  
44 current artificial neural network models, could only partially explain this cross-species  
45 correspondence. The shared behavioral variance unexplained by context-naïve neural data or  
46 computational models highlights fundamental knowledge gaps. Our findings demonstrate an  
47 intriguing alignment of human and monkey visual object processing that defies full explanation by  
48 either brain activity in a key visual region or state-of-the-art models.

## 49 Introduction

50

51 The field of visual neuroscience has long been fascinated by the computationally remarkable  
52 process of object recognition<sup>1-3</sup>, a cornerstone of primate visual perception. However,  
53 understanding an image transcends the ability to identify specific and isolated objects<sup>4-6</sup>.  
54 Interpreting an image requires knowledge about object correlations (e.g., bananas tend to co-  
55 occur with trees), relative object sizes (e.g., bananas are often smaller than trees), and relative  
56 object positions (e.g., bananas tend to be near the top part of a tree). Contextual information can  
57 dramatically alter how object information is interpreted<sup>7,8</sup>. There has been a long-standing interest  
58 in the statistics of natural images, and there are foundational behavioral studies of the role of  
59 context in vision<sup>9-14</sup>. The mechanisms behind incorporating contextual cues at the computational  
60 and neurophysiological levels remain poorly understood. Multiple prior studies focused on the role  
61 of context in relatively "low-level" visual phenomena such as extra-classical receptive fields and  
62 surround suppression<sup>15-19</sup>. However, little is known about how the brain represents prior high-  
63 level knowledge and integrates it with incoming inputs to modulate visual cognition.

64

65 Over the last decades, the field has made much progress in identifying the primate ventral visual  
66 pathway as crucial for housing neural circuits essential to object recognition<sup>4,20,5</sup>. A critical factor  
67 that led to progress in this domain has been the availability of rhesus macaques as an animal  
68 model that can mimic human object recognition behavior<sup>21,22</sup>. Given the ability to invasively probe  
69 finer-grain neural mechanisms in macaques<sup>23,24</sup>, studies have shown that a linear combination of  
70 image-driven population activity distributed across the macaque inferior temporal (IT) cortex (at  
71 the apex of the macaque ventral visual pathway) can sufficiently predict human object recognition  
72 behavioral error patterns on a battery of tasks<sup>5,25</sup>. Remarkably, these responses are typically  
73 recorded in monkeys who passively view the images without actively engaging in (or learning) the  
74 task -- suggesting that these representations are primarily bottom-up<sup>5,25</sup> and task-independent<sup>26</sup>.  
75 Furthermore, a significant effort to model the transformations that follow the retinal responses  
76 (driven by the image) and culminate into the pattern of activity in IT has recently come in the form  
77 of a set of artificial neural networks (ANNs) that can partly explain the neural responses along  
78 these pathways<sup>13,27,28</sup>. Therefore, a reasonable approach to probe the mechanisms underlying  
79 the visual processing of scene context is to ask if macaques also mimic human context-driven  
80 behavior. If so, one could empirically probe the underlying neural mechanisms and compare  
81 current ANNs' ability to explain those representations.

82

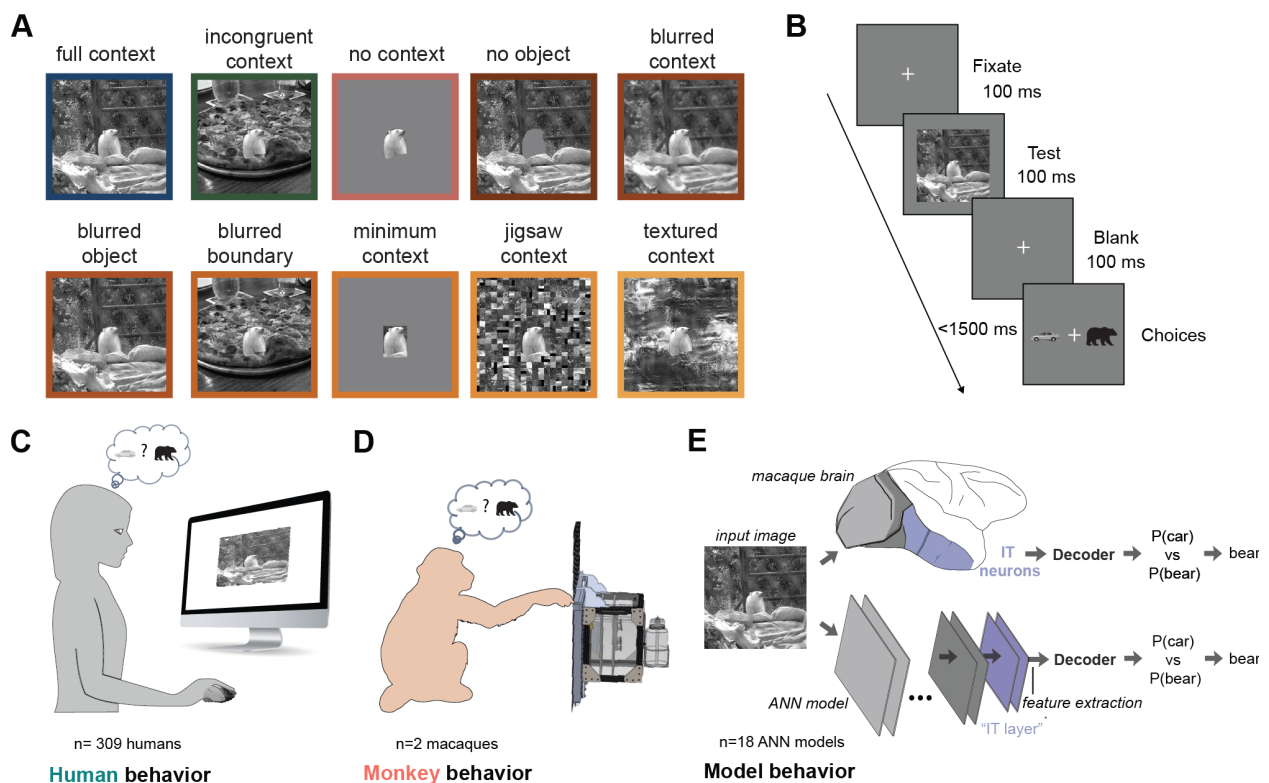
83 Interestingly, while current ANNs have been able to partially explain neural responses in V1<sup>29</sup>,  
84 V2<sup>30</sup>, V4<sup>27,31,32</sup>, and IT<sup>27,33,34</sup>, and many aspects of object recognition behavior<sup>22</sup>, recent studies  
85 have also shown that these models are heavily biased by the visual context during their training<sup>13</sup>  
86 which lead to their misalignment with human behavior. These models also develop specific biases  
87 (e.g., shape-texture bias) that do not align with human strategies<sup>35-37</sup>. With the increasing  
88 evidence of discrepancies between ANNs and human behavior, it is critical to figure out how these  
89 models can be improved. The ability to probe context-dependent behavioral biases in monkeys  
90 and their underlying neural mechanisms allows us to develop strong constraints that can guide  
91 future model development.

92

93 In this study, we first developed quantitative behavioral metrics (coarse to fine-grained) to  
94 evaluate the psychophysical effects of contextual changes during object discrimination. We then  
95 conducted a thorough comparative analysis of the behavior of humans and monkeys. We further  
96 performed large-scale neural recordings across the macaque IT cortex to probe the strength of  
97 the image-driven IT responses and explain the observed behavioral variances. We contrasted the  
98 IT representations with those retrieved from the current most human-aligned ANNs. Our results  
99 unveil a nuanced understanding of how context influences object recognition in biological and  
100 artificial systems, which highlights significant parallels but also divergences in how humans,  
101 monkeys, and ANNs process visual context information.

## 102 Results

103 We investigated the behavioral effects of scene context on humans and macaques during  
 104 recognition of real-world objects, such as cars, animals, and fruits. We introduced multiple  
 105 variations of the contextual information to further our understanding of what aspects of the object's  
 106 surrounding impact recognition. These variations include incongruent context, no context, and  
 107 blurred context, among many others (**Fig 1A**). We developed a binary delayed match to sample  
 108 object discrimination task (**Fig 1B**), where the participants, humans (**Fig 1C**) and monkeys (**Fig**  
 109 **1D**), identified the Target object shown in a sample test image (with varying contexts) when  
 110 probed with two object choices (a target and a distractor). We quantified context-driven  
 111 behavioral responses in both species with multiple quantitative metrics and assessed how well  
 112 these metrics matched each other. Next, to probe the nature of the neural representations that  
 113 could support these behavioral patterns, we examined how well the shared variance in their  
 114 behavior is explainable by neural data from the inferior temporal (IT) cortex and the IT-like sub-  
 115 units of current ANN models of primate vision (**Fig 1E**).  
 116  
 117

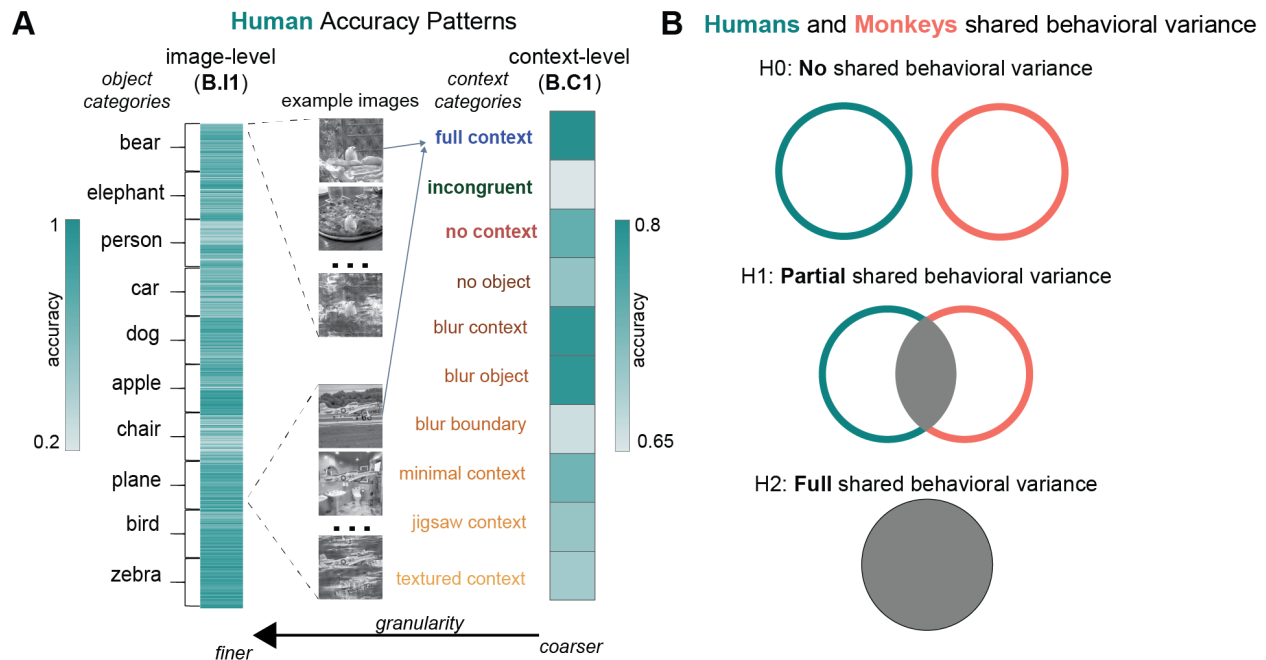


118 **Fig 1. Comparing the influence of context in object discrimination performance across humans,**  
 119 **monkeys and artificial neural networks (ANNs).** **A.** Example of the ten contextual manipulations for one  
 120 image of the set used for the experiments (details in Methods). The frames around each image indicate the  
 121 color associated with that context type (only used for reference in the article, not in the actual experiments).  
 122 **B.** Binary object discrimination task, showing the timeline of events for each trial. Subjects fixate on a  
 123 cross, then the test image containing one of ten possible objects and contextual manipulations is shown at the  
 124 center of the visual field (subtending 8 degrees of visual angle) for 100 ms. After a 100-ms delay, a  
 125

126 canonical view of the target object (the same category as, but not a template match to, the test image) and  
127 a distractor object (one of the other nine objects) appears. The human or monkey indicates which object  
128 was present in the test image by clicking on one of the two choices. **C.** Schematic of the human behavioral  
129 task for 309 participants recruited from Amazon MTurk. **D.** Schematic of the monkey behavioral task for  
130 two context-trained adult macaques. **E.** Schematic of the model behavioral task for eighteen pre-trained  
131 ANN models (bottom, details in Table 1) and the neural data (top). To make the artificial models compatible  
132 with the specific primate binary object discrimination task, their most IT-similar feature representations were  
133 extracted and used to train the decoder - a multiclass SVM classifier - calculating the cross-validated  
134 probabilities for each object class in a one-vs-all manner. The model output is then the object class with the  
135 highest one-vs-all probability. Similarly, the most reliable neural responses ( $n=122$  neural sites) from two  
136 context-naive monkeys were used to train the decoder and obtain the object class probabilities.

## 137 **Quantifying Context-Driven Changes in Object Recognition through Behavioral** 138 **Metrics**

139 To characterize how scene context influences the behavior of biological and artificial visual  
140 systems during object recognition, we developed quantitative metrics beyond the overall  
141 performance accuracy across all images. These metrics include the behavioral signature at the  
142 context level (B.C1, Behavioral, Context-Level 1-dimensional; see Methods) and a more fine-  
143 grained image level (B.I1, Behavioral, Image-Level 1-dimensional). The **context-level**  
144 **performance metric, B.C1** (human performances shown in **Fig 2A - right**), assesses the overall  
145 object discriminability within each context category (C). It does so by pooling accuracies across  
146 all images of a given context type (C) and all combinations of target and distractor pairs for those  
147 images (see Methods). This approach provides a broad understanding of how context influences  
148 recognition performance on a categorical level. In contrast, the image-level metric, B.I1 (detailed  
149 in Methods, human performance shown in **Fig 2A, left**), focuses on the discriminability of  
150 individual images, assessing how well the system distinguishes each object (O) from all others  
151 per image across varying contexts. This finer-grained metric allows for a more detailed analysis  
152 of performance variations at the image level. Expanding upon this foundation, we then seek to  
153 estimate the shared behavioral variance between humans and monkeys (behavioral signatures  
154 shown in **Fig S1A**), as depicted in **Fig 2B**. This comparative analysis could reveal one of the  
155 following scenarios. First, given species level differences<sup>38</sup>, we might observe that monkeys do  
156 not process visual context in the same way as humans and, therefore, exhibit no shared variance  
157 with humans (H0; **Fig 2B - top panel**). Second, it is possible that monkeys only share a fraction  
158 of variance with humans (H1; **Fig 2B - middle panel**). Lastly, it is also possible that within our set  
159 of tasks, images, and contextual variations – monkey and human behavior fully align with each  
160 other (H2; **Fig 2B - lower panel**). These conditions can be independently assessed for each of  
161 our behavioral metrics, and we expect that finer-grained metrics will enable us to more rigorously  
162 quantify the boundaries of the shared behavior between these two systems.

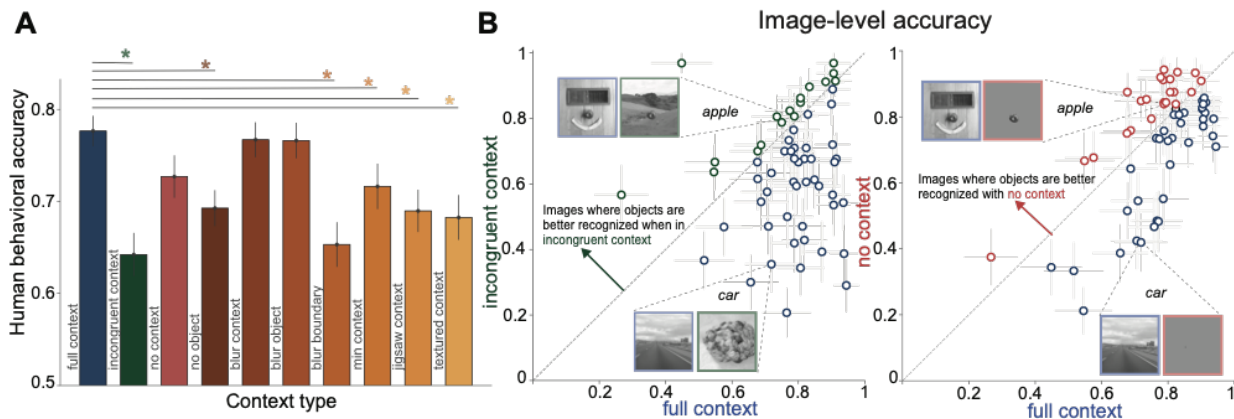


163  
 164 **Fig 2. Behavioral metrics to quantify context-driven variations in task performance.** **A.** Human  
 165 accuracy patterns at an image-level (fine granularity, B.I1, left) and context-level (coarse granularity, B.C1,  
 166 right). Each element of the B.I1 vector represents the overall accuracy (averaged across all tasks) for an  
 167 image. A few example images are shown in the middle panel grouped by object category. The context-level  
 168 signature (B.C1, right panel), is obtained by averaging the B.I1 values for all images of each specific context  
 169 type (see **Fig 1** for examples of all the context types). The light and dark teal colors indicate lower and  
 170 higher performances (see color scale next to each signature). **B.** The three hypotheses on the human-  
 171 monkey shared behavioral variance.  
 172

### 173 Object context induces significant changes in human behavior

174  
 175 Humans (309 participants on Amazon Mechanical Turk) participated in a binary object  
 176 discrimination task (**Fig 1B**, for details, see Methods). Our results show that varying the context  
 177 of the image changes the performance of the human participants. For instance, consistent with  
 178 previous research<sup>10,13,14</sup> humans show a significant reduction in accuracy for incongruent  
 179 compared to congruent contexts ( $\Delta$ Accuracy =  $0.13 \pm 0.21$ ; Lilliefors test: full context  $p=0.004$ ,  
 180 incongruent context  $p=0.371$ , non-normal distribution,  $p>0.005$ ; Wilcoxon rank-sum test  
 181 statistic=4.4,  $p=0.0001$ ; **Fig 3A**: blue vs green bars). The effect of contextual manipulations  
 182 resulted in a consistent pattern of behavior (with a trial-split reliability of approximately 0.74, see  
 183 **Fig S2A**, reliability across context types in **Fig S2B**). This high self-reliability is critical to ensure  
 184 that contextual effects can be compared across animals, across species, and from biological  
 185 systems to ANN models. The decline in accuracy for incongruent (compared to congruent) context  
 186 was not solely due to the abrupt transition from the background to the object; even when the  
 187 context/object boundary was blurred (termed blurred boundary), we observed the same effect.  
 188 Predictably, removing the object, retaining only its silhouette, also led to reduced accuracy;  
 189 however, performance remained well above chance, indicating that the context alone (with the

190 object outline) provided enough information for accurate object discrimination. Moreover, when  
191 the context was removed or minimized, there was again a decrease in performance, confirming  
192 that humans also rely on the surroundings for object recognition. The blurring process itself  
193 seemed to have minimal influence on human responses, as the kernel size used was relatively  
194 small (see Methods). Using a synthesized texture (textured context), which retained the visual  
195 attributes of the original context, also adversely affected human behavior. Our results align with  
196 extensive research on context modulated human behavior<sup>7,13</sup> and notably extend beyond the  
197 scope of previous work. In particular, we provide quantitative results from a forced binary choice  
198 task for a wider range of context variations. We define two behavioral signatures, allowing a  
199 coarse and fine-grain comparison within the human population and, importantly, across species -  
200 similarity with rhesus macaques. The cross-species consistency, coupled with access to the  
201 macaques' neural circuits, provides a path for studying the neural mechanisms underlying  
202 contextual processing.  
203  
204



205  
206 **Fig 3. Context-driven changes in human behavioral task performance.** **A.** Contextual manipulations  
207 produce significant changes in human visual recognition. Accuracy (mean  $0.71 \pm 0.05$ ) for each contextual  
208 manipulation (B.C1, **Fig 2**), with standard error across images. Statistics are shown for full context  
209 compared to other context variations (\* denotes independent t-test,  $p < 0.05$ ). **B.** Left: Image-level accuracies  
210 (from the object discrimination task, with standard error across image trials) for full and incongruent context,  
211 each dot represents the human behavioral accuracy for the same object embedded in either full or  
212 incongruent context. Example images are shown where the object is better predicted in each context  
213 variation. Right: Similar as left, but comparing the accuracy for objects embedded in full context vs removing  
214 the context. Note that the car object is very small ( $< 1$  degrees of visual angle) and hard to see without a  
215 lot of zoom (images are presented at the center of the visual field subtending 8 DOV angle, **Fig 1**).  
216

## 217 Object context induces significant changes in monkey behavior

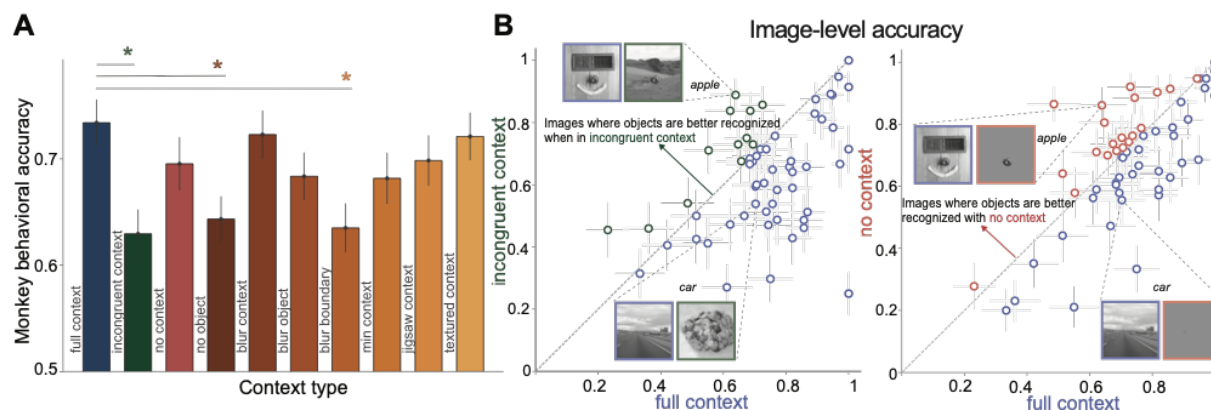
218  
219 To establish macaques as an appropriate animal model to probe the neural mechanisms of  
220 human context processing it is critical to first ask whether macaques behaviorally show similar  
221 contextual effects to humans. To ensure that macaques are familiar with the scene context per  
222 object category, we first explicitly trained them with images in full context (from the Microsoft



223 COCO dataset, 160 images per object for ten objects). Macaques showed robust cross-validated  
224 accuracy during such training (**Fig S1B**).

225 Once the monkeys (n=2) were fully trained (i.e., reached  $\geq 80\%$  performance) in their home cages  
226 (see learning curve **Fig S1B**), we measured their object discrimination performances with the  
227 same contextually manipulated images as humans (**Fig 1A**). Monkey behaviors were highly  
228 reliable (as measured by trial split-half reliability,  $r=0.76$ , see Methods, **Fig S3A**), and correlated  
229 with each other at both the context level (corrected Pearson  $R=0.98$ , corrected by both monkeys'  
230 self-consistency, see Methods, **Fig S4A**), and at the image-level (corrected Pearson  $R=0.83$ , **Fig**  
231 **S4B**). Similar to humans, monkeys also showed a significant reduction in accuracy for  
232 incongruent compared to congruent contexts ( $\Delta$ Accuracy =  $0.104 \pm 0.18$ ; Lilliefors test: full context  
233  $p=0.173$ , incongruent context  $p=0.58$ , normal distribution,  $p>0.005$ ; independent t-test,  $t(59) =$   
234  $3.305$ ,  $p=0.001$ , **Fig 4A**: blue vs green bars). **Fig 4B** compares the trial averaged image by image  
235 accuracy between full and incongruent context (left), as well as full and no context (right). At the  
236 individual image level, we observe some images for which the object placed in an incongruent  
237 context was better recognised than when the same object was embedded in a congruent context  
238 (see example of an apple in **Fig 4B**, left). Similarly, some objects were better recognized when  
239 fully removing the context compared to keeping the full congruent context (see apple example in  
240 **Fig 4B** right).

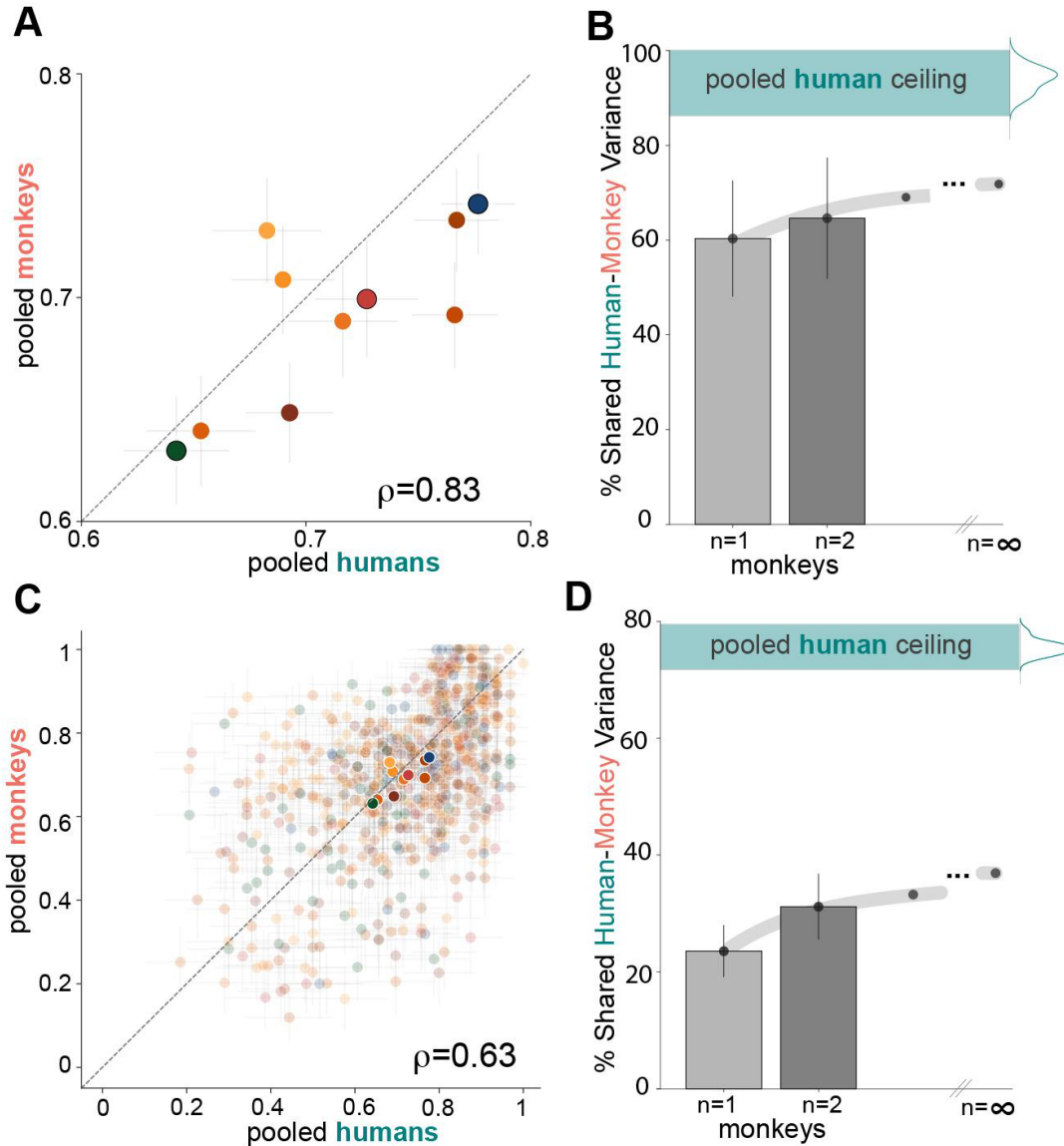
241  
242



243 **Fig 4. Context-driven changes in monkey behavioral task performance.** **A.** Contextual manipulations  
244 produce significant changes in monkey behavior. Accuracy (mean  $0.68 \pm 0.04$ ) for each contextual  
245 manipulation (B.C1), with standard error across images. As in **Fig 3A**, statistics are shown for full context  
246 compared to other context variations (\* denotes independent t-test,  $p<0.05$ ). Chance = 0.5. **B.** Left: Image-  
247 level accuracy for incongruent context versus full context (format as in **Fig 3B**), each dot represents the  
248 pooled monkeys' behavioral accuracy for the same object embedded in incongruent (y-axis) or full (x-axis)  
249 context with standard error across image trials. We show example images where the object is better  
250 recognized in each context variation. Right: Similar as left, but for no context (y-axis) versus full context (x-  
251 axis).  
252

253 **Humans and monkeys share significant variance in context driven changes in**  
254 **object recognition**

255 We directly compared monkey and human performance for the same images and task. Our results  
256 show a remarkable consistency between monkeys and humans at the context level (C1 corrected  
257 Pearson  $R=0.83$ , **Fig 5A**). For example, both monkeys and humans performed best in the full  
258 context condition (blue point) and worst in the incongruent context condition (green point).  
259 However, the majority of points in **Fig 5A** fall below the diagonal, indicating that humans  
260 outperformed monkeys in most context conditions ( $\Delta$  (human - monkey)  $=0.02\pm 0.03$ , Lilliefors  
261  $p=(0.517, 0.487)$  : normal distribution; paired t-test,  $t(9) = 1.86$ ,  $p=0.1$ ). The two exceptions were  
262 the jigsaw and textured context conditions, where monkeys slightly outperformed humans. To  
263 quantify the variability across humans, we calculated the human ceiling by comparing the shared  
264 variance between two separate pools of human subjects (teal band in **Fig 5B**). We then compared  
265 the shared human-monkey variance to this human ceiling. Since we are comparing a pooled  
266 population of 309 humans to the  $n=2$  monkey pool, we looked at the effects of monkey pool size  
267 on its consistency with human data. As the number of monkeys in the pool increased from one to  
268 two, the shared human-monkey variance increased by 4.3% (gray bars in **Fig 5B**). Extrapolating  
269 to an infinite pool of monkeys using a "pseudo" human consistency function (sigmoid) derived  
270 from subsampling the human pool, we estimate that the asymptotic shared variance between  
271 monkeys and humans would reach approximately 80% of the human ceiling. Next, we compared  
272 monkey and human performance at the individual image level (**Fig 5C**). Again, we found a  
273 significant correlation between monkeys and humans (I1 corrected Pearson  $R=0.63$ ), although  
274 the relationship was weaker than at the context level. The slope of the regression line in **Fig 5C**  
275 suggests that humans outperformed monkeys on average, but this difference was not as  
276 pronounced as at the context level ( $\Delta$  (human - monkey)  $=0.02\pm 0.18$ , Lilliefors  $p=(0.001, 0.001)$ ;  
277 non-normal distribution; Wilcoxon test: statistic = 79913.5,  $p=0.02$ ). The shared variance analysis  
278 at the image level (**Fig 5D**) revealed that humans were less consistent with each other compared  
279 to the context level (**Fig 5B**), as expected due to the increased granularity of individual images.  
280 This effect was even more pronounced for monkeys, with a larger drop in shared variance at the  
281 image level compared to the context level. Increasing the number of monkeys in the pool from  
282 one to two improved the shared human-monkey variance by 8.2% at the image level (**Fig 5D**).  
283 Extrapolating to an infinite pool of monkeys, we estimate that the asymptotic shared variance  
284 would reach approximately 70% of the human ceiling at the image level.



285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299

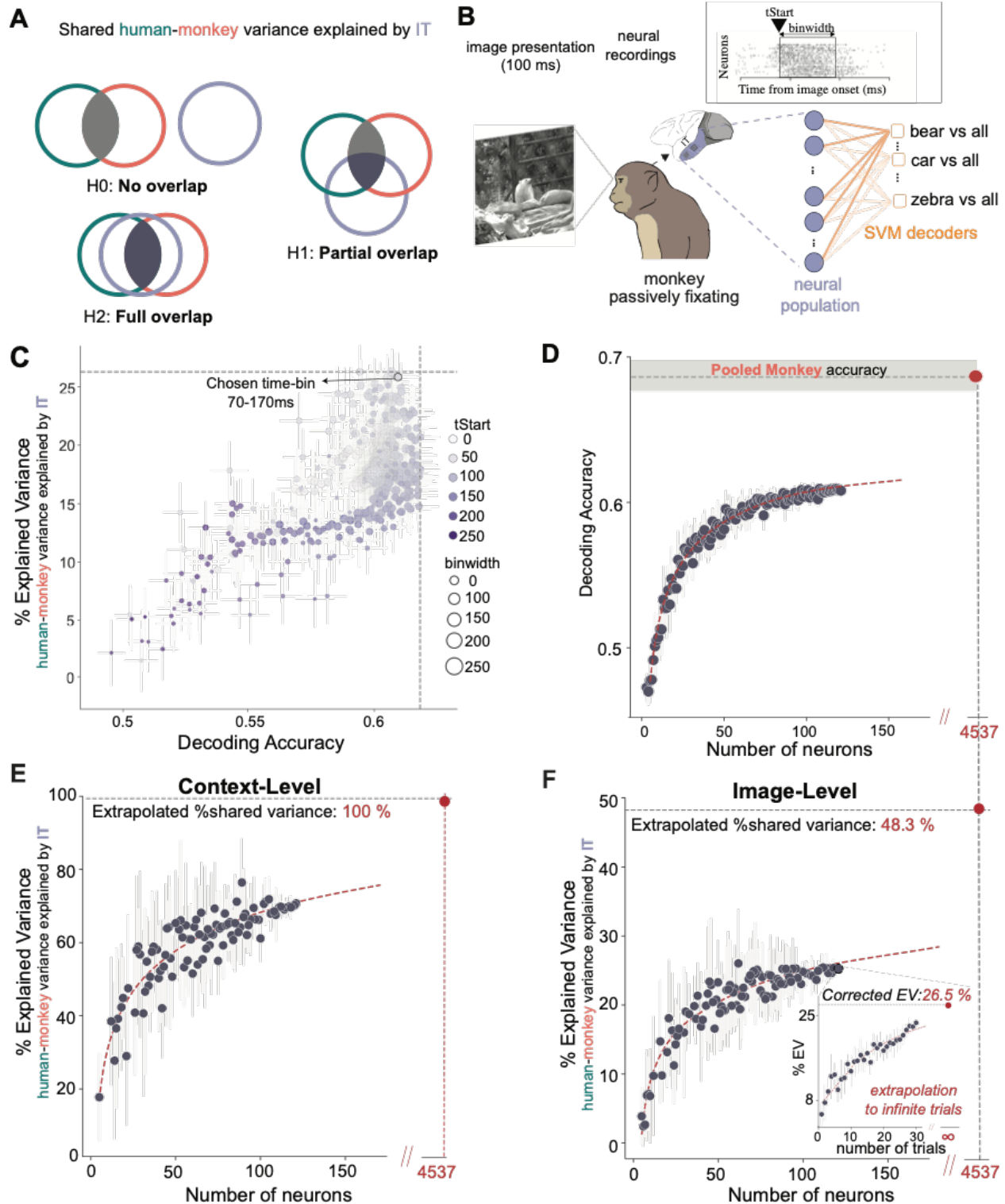
**Fig 5. Monkeys and humans show similar (but not identical) context-driven behavioral changes. A.** Context-level (B.C1) correlation between the pooled monkeys ( $n=2$ ) and pooled humans ( $n=309$ ). Each point represents the mean accuracy for a contextual variation with standard error across images of that context type (colors as in **Fig 1A**, pooled monkeys mean  $0.69 \pm 0.20$ , humans mean  $0.71 \pm 0.18$ ). The three main context types: full (blue), incongruent (green) and no context (red), are shown with a black stride around the filled point. The value  $\rho$  indicates the noise corrected correlation coefficient (Pearson R). **B.** Shared human-monkey explained variance at a context-level (mean with standard deviation across context types), as a function of the number of monkeys used for pooling. The asymptotic value for an infinite pool of monkeys is obtained by extrapolating the “pseudo” human consistency function (Methods). The human self-consistency ceiling is shown as a teal band. **E.** Image-level correlation (B.I1) for the pooled monkeys and humans, each low opacity point shows the performance (mean accuracy) for an image with standard error over image trials, the higher opacity points are the B.C1 mean (from A), colors map to context types as defined in **Fig 1A**. **F.** Similar to C but mean shared variance at an image-level with standard deviation across image subsamples.

300 **Population activity across the IT cortex in a context-naive monkey fully explains**  
301 **the shared behavioral variance between humans and monkeys at the overall**  
302 **context-level**

303  
304 Our behavioral results (**Fig 5**) demonstrate that humans and macaques share a significant  
305 proportion of variance (context level shared EV=62.43%) induced by context variations during  
306 object discrimination. To understand the neural mechanisms behind these contextual influences,  
307 we require a more detailed examination of the neural networks involved. Previous studies have  
308 shown that IT population responses in monkeys (passively viewing images, see Methods) can be  
309 linearly combined to sufficiently explain human object category (and category-orthogonal) based  
310 behavioral patterns<sup>1,5,25</sup>. Therefore, we aimed to assess the extent to which the image-driven  
311 responses in the IT cortex of context-naive macaques could account for the variance observed  
312 between humans and monkeys. Similar to the expected observations while comparing human  
313 and monkey behavior (**Fig 2B**), we hypothesized that there could be no overlap (*H0*; **Fig 6A**),  
314 partial overlap (*H1*; **Fig 6A**), or full overlap (*H2*; **Fig 6A**) between the neural predictions and  
315 primate behavior.

316  
317 We performed chronic neural recordings using Utah arrays across the IT cortex in two macaques  
318 that passively viewed the images (used in the behavioral tasks) presented for 100 ms each (**Fig**  
319 **6B**, see Methods). We combined the most reliable neural sites (n=122; see criteria in Methods,  
320 30 sites from monkey 1, 92 sites from monkey 2) across the two monkeys to generate a pooled  
321 neural population for further analysis. Similar to previous methods<sup>5,23,28</sup>, we used linear  
322 classification-based algorithms (**Fig 6B**) to decode the object category for each image from the  
323 pooled neural data and estimated the neural predictions for the behavioral metrics (explained  
324 above, e.g. C1, I1).

325  
326 We first asked how well the macaque neural responses can predict the shared variance between  
327 humans and macaques at the B.C1 level. Therefore, we performed a partial correlation analysis  
328 between human and macaque C1 behavioral patterns while controlling for the IT population  
329 activity-based predictions of B.C1. To account for the irreducible noise in the neural data, we  
330 corrected the partial correlation by extrapolating it to an infinite number of trials for the neural data  
331 (see inset **Fig 6F**). Interestingly, the neural data (122 sites) explained 75% of the context-level  
332 shared monkey-human C1 variance (**Fig 6E**). To further address the data limitations arising from  
333 the limited number of neural recordings, we extrapolated the neural decoding accuracy to match  
334 the monkey accuracy (logarithmic function, **Fig 6D**). This extrapolation led to an estimation of  
335 4357 neural sites needed to reach monkey accuracy. A logarithmic extrapolation of the explained  
336 variance (EV) to 4357 neural sites indicates that IT would fully explain the human-monkey B.C1  
337 variance if we had more neural recordings (**Fig 6C**).



338

339 **Fig 6. Context naive macaque IT fully explains the human-monkey shared behavioral variance at a**  
 340 **context-level but only partially at an image-level. A.** The hypotheses for how much of the human-  
 341 monkey shared explained variance (HM-EV) can be explained by IT. **B.** The neural data was recorded  
 342 while the monkey was passively fixating on the center of an image (8 degrees of visual angle) presented at  
 343 the animal's center of gaze for 100 ms. The object category decoding was done by training a multi-class

344 SVM classifier (one vs all for each object category) tested in a cross-validated way on the same images  
345 and tasks as those presented to humans and monkeys. For a given image, the decoding output is the object  
346 class with the highest one vs all probability. All behavioral predictions from the decoder were for images  
347 where the object was not seen in any phase of the model training, making sure we never show an image  
348 of the same object (regardless of the contextual manipulation) during the fitting and testing. We decoded  
349 the object category from each possible time-bin of the neural data by varying the tStart (start of the time-  
350 bin with respect to image onset, in ms) and binwidth (length of the time-bin, in ms) of the obtained neural  
351 population vector (0-300 ms per image presentation). **C.** Results from decoding all time-bins (filtered with  
352 self-consistency >0.1) from the neural data, color indicates the bin start, size indicates the bin length. The  
353 percent of image-level explained variance from the shared human-monkey variance is shown (y-axis, with  
354 standard error across image subsamples) as a function of the decoding accuracy for each bin (x-axis, with  
355 standard error across images). We used the 70-170 ms time-bin for all subsequent analyses. **D.** Decoding  
356 accuracy with standard deviation (one-vs-all accuracy, chance level = 0.5) across neuron subsamples for  
357 the 70-170 ms time-bin, as a function of the number of neurons. An extrapolation (dashed red curve)  
358 estimates the decoded accuracy from a neural population of 4537 recorded neural sites would reach the  
359 overall pooled monkey accuracy (0.69, gray band shows monkey accuracy mean with standard error across  
360 the 600 images). **E.** Context-level variance explained by the neural data, from the HM-EV. The EV is  
361 obtained by subtracting the HM-EV when controlling for the neural data (partial correlation) from the full  
362 HM-EV and normalizing by the full HM-EV (see [Methods](#)). We show the EV as a function of the number of  
363 neurons used for decoding, showing an extrapolation to 4537 neurons would fully explain the B.C1 HM-EV.  
364 Each point shows an average (with standard deviation error bar) across ten different subsamples of neurons  
365 used, corrected by extrapolating to an infinite number of trials for those specific neurons. **F.** Similar to E,  
366 but for image-level shared variance. The inset shows the correction for the EV for 122 neurons by  
367 extrapolating to an infinite number of trials as done for context-level, each point shows an average (with  
368 standard deviation error bar) across ten different subsamples of trials. The extrapolation of the EV to the  
369 number of neurons needed to reach monkey accuracy (see decoding accuracy extrapolation in D) gives a  
370 ceiling of 48.3% of image-level human-monkey behavioral variance that can be explained by the context  
371 naive monkey IT neural data.

## 372 **Population activity across the IT cortex in a context-naive monkey only partially** 373 **explains the shared context-driven behavioral variance between humans and** 374 **monkeys at the image-level**

375  
376 To further stress test whether IT responses from untrained (task-naive) monkeys can explain finer  
377 grained behavioral patterns, we next turned to predictions for the I1 level (image-level shared  
378 variance). As shown in **Fig 6F**, the recorded reliable neural population (122 neural sites) explains  
379 only a fraction of the image-by-image behavioral variance (up to 25%). This result suggests that  
380 the context-naive IT population may not capture all the necessary information to fully predict the  
381 shared human-monkey behavioral patterns at the image level. To address the possibility that the  
382 limited explained variance might be due to the restricted number of recorded neurons, we applied  
383 the same extrapolation method as used for the context-level EV (**Fig 6E**). We estimated that  
384 approximately 4357 neural sites would be needed to match the pooled monkey behavioral  
385 accuracy (**Fig 6D**). However, despite this extrapolation, the neural data from context-naive IT  
386 could not fully explain the image-by-image shared primate variance, reaching a ceiling of only  
387 48.3% (**Fig 6F**). The discrepancy between the context-level and image-level explained variance  
388 highlights the complexity of the neural mechanisms underlying context-dependent object

389 recognition and the limitations of using context-naive neural responses to predict fine-grained  
390 behavioral patterns. In summary, while the context-naive IT population activity can fully explain  
391 the shared human-monkey behavioral variance at the context level (**Fig 6E**), it only partially  
392 accounts for the variance at the image level (**Fig 6F**). This finding underscores the need for further  
393 investigation into the neural mechanisms that shape the shared behavioral patterns between  
394 humans and monkeys in the presence of contextual variations.

### 395 **Low-level image-based features do not explain the shared human-monkey** 396 **behavioral variance**

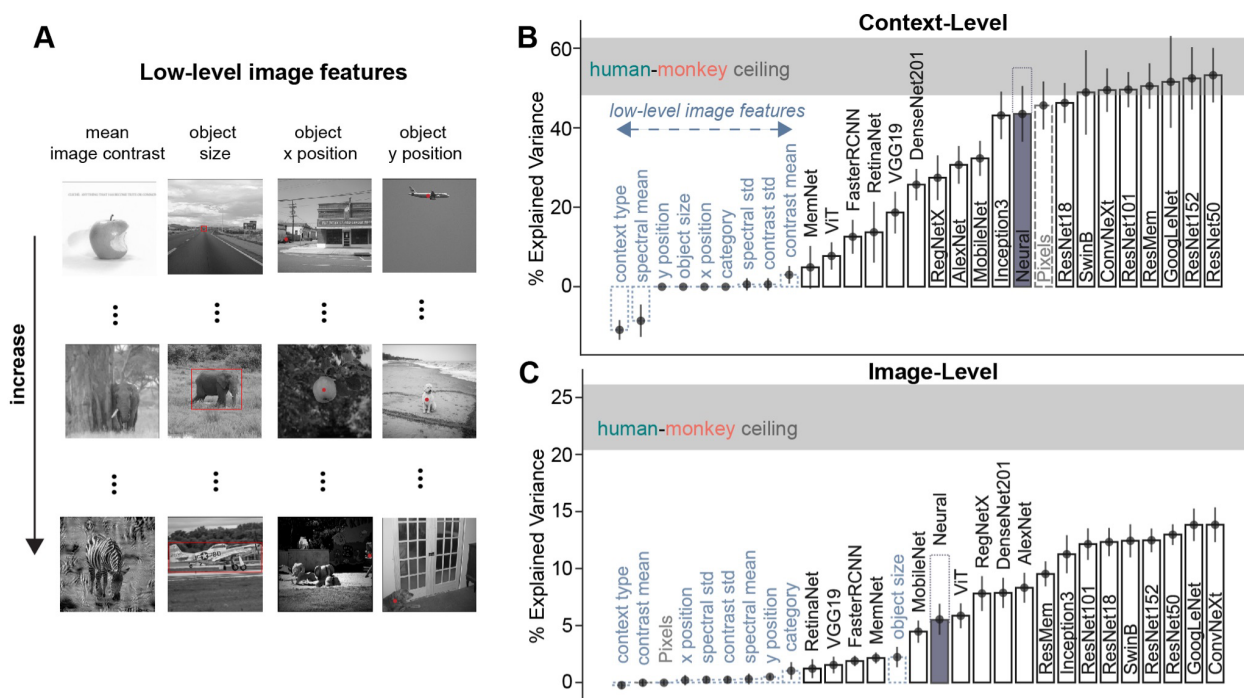
397 So far, we have observed that human and macaques share a significant amount of behavioral  
398 variance both at the coarse (B.C1) level and the finer-grained (B.I1) level. The image-driven task  
399 naive IT responses can fully explain the C1 variance but not the I1 level variance. We next asked  
400 how much of these results can be explained by low-level image features. For every image, we  
401 extracted a range of basic image features, such as object size, location, and category, spectral  
402 mean and standard deviation(std), and contrast mean and standard deviation (std) (**Fig 7A**). The  
403 low-level features were chosen to capture basic properties of the images that could potentially  
404 influence object recognition performance, such as the object's saliency (contrast) and its  
405 placement within the scene (location). We observed that these low-level features do not explain  
406 the context-level (**Fig 7B**,  $55.39 \pm 7.08\%$  mean  $\pm 95\%$  CI of noise ceiling, max low-level feature EV  
407 =  $3.03 \pm 2.2\%$ ) or image-level (**Fig 7C**,  $23.4 \pm 2.9\%$  mean  $\pm 95\%$  CI of noise ceiling, max low-level  
408 feature EV =  $2.47 \pm 1.47\%$ ) measured shared behavioral variance. Among all the low-level features  
409 tested, object size showed the most consistency with the shared human-monkey variance at the  
410 image level, aligning with prior studies highlighting its influence on human behavior (Zhang et al.,  
411 2020). In particular, the positive correlation with human and monkey performance was more  
412 significant for smaller object sizes and diminished for objects with size beyond 5 degrees of visual  
413 angle (**Fig S5**). Its effect, however, is marginal, accounting for only 10% of the shared human-  
414 monkey image-level variance. This suggests that while object size plays a role in shaping the  
415 shared behavioral patterns, it alone cannot fully explain the observed consistency between  
416 humans and monkeys. Taken together, we infer that low-level image features alone are  
417 insufficient to explain the shared behavioral patterns observed in humans and monkeys, indicating  
418 the need for more complex, higher-order processing to fully account for the context-dependent  
419 object recognition performance.

420

### 421 **ANNs fully explain the overall context-level human-monkey shared behavioral** 422 **variance**

423 Next, we tested whether the current best models of primate vision, a family of deep convolutional  
424 neural networks (DCNNs), vision transformers, or recurrent convolutional neural networks, can  
425 predict the behavioral variance observed on a context and image-by-image level. Using a  
426 multiclass SVM decoder, we mimicked the same object discrimination task presented to humans  
427 and monkeys. We used the most 'IT-like' layer features from each model (for details, see Table  
428 1), projected to a 3k lower dimensional space (via a Random Gaussian Projection). These model  
429 accuracy decodes showed sensitivity to contextual changes, with their accuracy varying across

430 different context types and images, as shown in their behavioral signatures (**Fig S9**). The ANN  
 431 models' ability to capture context-dependent performance variations suggests that they have  
 432 learned to extract and process contextual information in a manner that is relevant for primate  
 433 object recognition. Similarly, for the neural data, we did a partial correlation analysis for the pooled  
 434 monkey and human population behavioral patterns while controlling for each of the artificial  
 435 models' variance. Our results show that most models, including the Pixels control model, can  
 436 explain the context-level shared behavioral accuracy patterns of humans and monkeys.  
 437



438 **Fig 7. ANNs fully explain the human-monkey shared behavioral variance at a context-level but only**  
 439 **partially at an image-level.** **A.** Example of images with increasing “intensity” (top to bottom) of 4 example  
 440 low-level image features: mean image contrast, object size, object x and y position. The objects are noted  
 441 with a red dot or bounding box. **B.** The human-monkey shared variance explained by the low-level features  
 442 and ANN models at a context level, showing the mean fraction of explained shared variance with standard  
 443 deviation across different image subsamples from 20 bootstraps (choosing 600 images with repetition). The  
 444 low-level features are shown with dotted light blue bars. The ANN decoding was done in the same way as  
 445 for the neural population (multiclass SVM), only using the extracted model features from the ‘IT’ layer for  
 446 each model. Pixels’ (control model - flattened image pixels) performance is shown in a dashed gray bar.  
 447 The mean human-monkey shared variance ceiling is shown in gray, with standard deviation across different  
 448 image subsamples from 20 bootstraps (same as for the bars). We are noting the neural corrected EV  
 449 (purple) when using all the recorded reliable neural responses (122), and the extrapolation (from **Fig 6**)  
 450 with the dotted bar. **C.** Similar to B but showing the explained variance at an image-level.

452 **ANNs only partially explain the image-level human-monkey shared behavioral**  
 453 **variance**

454 While ANNs fully capture the shared human-monkey behavioral variance at the context-level,  
 455 their performance at the finer-grained, image-level is less comprehensive. Despite the models'  
 456 ability to explain the overall context-dependent behavioral patterns, they struggle to account for



457 the more intricate, image-specific variations in primate behavior. At B.I1 level, the models explain  
458 at most 70% of the image-level shared primate behavioral variance. This discrepancy is due to  
459 the finer grain accuracy variations within both context and object category types that are not  
460 consistently aligned between the primates and the artificial models. We found a strong correlation  
461 between the fraction of explained shared human-monkey variance and the decoding accuracy  
462 from the model features (Pearson  $R=0.78$  for B.C1 and  $0.95$  for B.I1, Fig S7), indicating that  
463 improving the model accuracy could allow them to fully explain the shared human-monkey B.I1  
464 behavioral variance. The control Pixels model - using the raw image pixel values, capturing the  
465 context-level shared behavioral patterns, was falsified at an image level. This reveals more  
466 complex image-level shared behavioral patterns that are not due to the raw image features. The  
467 image-level gap was consistent when comparing at an individual level - these models could not  
468 fully explain the (full) human, monkey or neural image-level behavioral patterns (**Fig S2, S3** and  
469 **S8**). This indicates that such models do not currently possess the mechanisms required to  
470 process scene context in a primate-like fashion.

471  
472 In summary, while low-level image features and current artificial neural networks can account for  
473 the overall context-level shared behavioral variance between humans and monkeys, they fall  
474 short in fully explaining the more intricate, image-level behavioral patterns. These findings  
475 highlight the need for further advancements in artificial neural network architectures and training  
476 paradigms to better capture the nuanced, context-dependent object recognition processes  
477 observed in primates.

## 478 Discussion

479 In this study, we highlighted the critical role of context in primate object recognition. The visual  
480 object recognition abilities of monkeys that were initially trained to categorize objects in their  
481 natural context were strongly modulated when we deliberately varied the contextual cues. Our  
482 findings reveal that both humans and monkeys exhibit a significant sensitivity to contextual cues,  
483 which goes beyond low-level image attributes. Indeed, macaques shared a significant variance in  
484 their context-driven behavioral error patterns with humans. Thus, we established rhesus  
485 macaques as a viable animal model for investigating scene context in human visual recognition,  
486 paving the way for further studies into the neural underpinnings of contextual modulation.  
487 However, our analysis also revealed that at the image-level, monkeys do not entirely mimic  
488 human behavioral patterns, suggesting potential limitations in the depth and duration of their  
489 training or inherent species-level differences in sensory processing and cognition. In addition, we  
490 observed that the population activity distributed across the IT cortex of naive monkeys that were  
491 not explicitly trained with objects in context do not fully explain the context-driven behavioral  
492 patterns of the context-trained monkeys. Furthermore, our ANN-based simulations further reveal  
493 the substantial impact of context on the predictive behavior of current ANN models. Notably,  
494 ANNs exhibit limitations in their explanatory power for image-level comparisons with primates  
495 under varying contexts, indicating a clear need for model enhancements to accurately mimic the  
496 complex influence of context in primate visual recognition.  
497

### 498 **Context modulates visual object recognition in humans and monkeys**

499 Our results underscore the critical role of context in primate object recognition, aligning with an  
500 extensive corpus of literature on visual cognition<sup>11,13,39</sup>. This research has established that human  
501 visual object recognition capabilities are modulated by contextual cues. Such cues are informed  
502 by our understanding of object occurrence statistics, which dictate notions of congruency or  
503 incongruency within a given scene. Interestingly, our findings reveal that monkeys, much like  
504 humans, exhibit sensitivity to these statistical cues. Across diverse context manipulations, we  
505 observed substantial decrements in object discrimination accuracy compared to fully congruent  
506 scenes - up to 13% in humans and 10.4% in monkeys for incongruent contexts. These striking  
507 parallels between the two primate species underscore the viability of macaques as a model  
508 system for probing the neural computations underlying context processing. Critically, our results  
509 extend beyond prior work by demonstrating context sensitivity across a broad range of  
510 manipulations and employing rigorous, multi-faceted behavioral metrics designed to quantify  
511 performance changes induced by contextual cues. The tight concordance points to potential  
512 shared cognitive mechanisms, such as knowledge of object co-occurrence statistics, relative  
513 sizes, and positional regularities, which could account for the context facilitation effects observed  
514 in both species.

### 515 **Goodness of monkeys as a model of human contextual processing**

516

517 While the largely consistent effects at the coarser, context-level, validate macaques as a model,  
518 some discrepancies remain. Our analysis revealed that at the image-level, monkeys do not  
519 entirely mimic human behavioral patterns. To familiarize them with various contexts—such as  
520 cars on roads, bears in the jungle, and chairs in rooms—we trained these monkeys extensively  
521 with natural photographs (from the MS COCO image dataset) until their performance plateaued,  
522 as shown in **Fig S1B**. Despite reaching high performance, the mismatch between human and  
523 monkey responses suggests that the depth and duration of training might not have been sufficient.  
524 Enhancing the training regimen could potentially lead to a better alignment with human context-  
525 level behavior, reducing the disparity observed in image-level variance. However, potential  
526 confounds like the limited stimulus set size and specific task demands cannot be ruled out either.  
527 Importantly, instances where context manipulations like blurring had relatively small impact on  
528 performance in both species provide insights into boundary conditions that inform and constrain  
529 models of contextual reasoning. Another critical consideration is the inherent species-level  
530 idiosyncrasies and differences in brain structures between humans and monkeys<sup>38</sup>. These  
531 biological distinctions might inherently limit the degree to which monkeys can model human  
532 contextual processing. While further training might narrow the behavioral gap, some level of  
533 divergence might always persist due to fundamental differences in sensory processing, visual  
534 experience, and cognition between the two species. Understanding and acknowledging these  
535 limitations is vital as we continue to refine monkeys as models for human visual processing.  
536 Future research should explore both the potential and the boundaries of this animal model, aiming  
537 to optimize training strategies and deepen our understanding of the species-specific factors that  
538 influence contextual processing. Through this nuanced approach, we can better leverage the  
539 strengths of monkeys as models while being mindful of their inherent limitations.

#### 540 **Insufficiency of ANN models to explain primate context-driven behavior**

541 Deep ANNs are currently the best models of human vision and also show remarkable  
542 performance in computer vision tasks<sup>3,28</sup>. These models have been trained extensively on images  
543 of objects in context from large datasets (typically ImageNet). Our findings show that while these  
544 ANNs were able to fully explain the context level (B.C1; the coarser metric) shared primate  
545 variance, they failed to completely capture the finer grain image level accuracy patterns (B.I1).  
546 Even simple pixel-based models could predict the broad variations in B.C1 (**Fig 7**), underscoring  
547 the limitation of such coarse metrics in capturing the nuanced differences in visual context  
548 processing. However, a shift in focus to finer, image-by-image level variations revealed a more  
549 intricate picture. At this granular level, we discerned the primary distinctions between humans,  
550 monkeys, and ANNs. While monkeys show partial overlap with human behavior, a significant  
551 portion of this image-level variance remains unexplained by current ANN models. This gap  
552 highlights a critical area where artificial systems diverge from natural primate visual processing,  
553 suggesting that while ANNs can mimic some aspects of primate vision, they still lack certain  
554 mechanisms that drive the nuanced, context-driven behaviors observed in humans and monkeys.  
555 These observations not only challenge the sufficiency of broad behavioral metrics in capturing  
556 the essence of visual context processing but also point to image-level analyses as a more  
557 sensitive and discriminating tool for understanding the subtleties of primate vision. The partial  
558 alignment yet notable divergence of current ANN models from primates points towards key

559 computational mechanisms underlying context integration during object recognition that may still  
560 be lacking in artificial systems. Aspects like rapid integration of segmented objects with contextual  
561 associations and scene statistics, combination of high-resolution foveal and low-resolution  
562 peripheral representations, oculomotor sampling routines tuned for context (however, see<sup>13, 14</sup>),  
563 or other dynamic processes could be critical for human-level contextual reasoning. Pinpointing  
564 and distilling such mechanisms from the primate brain represent exciting future directions. We  
565 tested a range of models (Table 1), to gain further insight into the model architectures that could  
566 explain the B.I1 primate shared behavioral patterns better. We observed that deep ANNs with  
567 residual connections, as well as inception modules, are most aligned with human-monkey  
568 behavior (**Fig S6**). This indicates that allowing for feed-forward long-range dependencies between  
569 features (e.g., low-level features like edges with higher-level features) and preserving the finer-  
570 grained information from earlier layers (which can be lost due to the depth of models) by using  
571 bypass connections could benefit the alignment of these ANNs with primate behavior.  
572 Furthermore, ANN decoding accuracy (signatures in **Fig S9**) predicts the fraction of explained  
573 monkey-human shared variance (**Fig S7**), indicating that by improving the model's decoding  
574 accuracy, we could come closer to bridging the I1 explainability gap.

### 575 **Role of IT cortex in processing scene context**

576 The inferior temporal (IT) cortex is integral to visual object processing<sup>4,5,23,40</sup>, yet our findings  
577 indicate that responses from context-naive monkeys may not fully encapsulate the representation  
578 of scene context akin to that in humans or context-trained monkeys. This shortfall calls for a  
579 nuanced approach in future investigations into the IT cortex's role in context processing. One  
580 explanation for this is that our data might be sample-limited, affecting the breadth and depth of  
581 our inferences. Constraints such as the extent of IT neural data sampling, the diversity of images,  
582 trials, objects, and context variations might have curtailed our ability to fully capture IT's  
583 capabilities in context processing. To address these limitations, we conducted extrapolation  
584 analyses (**Fig 6**) to estimate the scaling laws governing our data, aiming to predict how increasing  
585 our sample might influence our findings – further corroborating the insufficiency of naive IT-based  
586 decodes to explain human behavior. Secondly, the lack of refined representational capacity in the  
587 IT cortex of naive monkeys might be due to insufficient exposure to varied contextual cues –  
588 improving which might amplify the IT cortex's ability to represent scene context. Additionally,  
589 investigating the interaction of the IT cortex with other brain regions, both within the ventral stream  
590 like areas V4 and outside the ventral pathway such as the ventrolateral prefrontal cortex (vlPFC),  
591 and their correlation with behavior in trained and untrained monkeys could illuminate new aspects  
592 of neural processing. This exploration is crucial to discern whether other areas might compensate  
593 for or augment the IT cortex's function in context processing, thus providing a more holistic view  
594 of the neural networks at play in this intricate task. Together, these strategies will deepen our  
595 understanding of the IT cortex's role and pave the way for a more comprehensive grasp of the  
596 neural underpinnings of context processing in vision.

597  
598 By bridging behavioral, computational and neural levels of analyses<sup>41</sup>, we can develop integrated  
599 accounts reconciling the cognitive influences of context with their neural underpinnings and use  
600 them to inspire more neurally-grounded computational models. Overall, this multi-pronged

601 approach paves the way for a deeper understanding of how context facilitates robust object  
602 perception across primates.

## 603 **Methods**

### 604 **Visual Stimuli**

605 We generated an imageset comprising 600 grayscale images from 10 object categories (bear,  
606 elephant, person, car, dog, apple, chair, plane, bird, zebra). For each object category, we selected  
607 six natural images from the Microsoft Common Objects in Context (COCO) dataset, varying in  
608 object size and location, which were center cropped, converted to grayscale, and recalled to  
609 512x512 pixels. We then generated 10 different contextual variations for each image. The  
610 changes were made using the object segmentation for each image obtained from the COCO  
611 object annotation masks and (for some conditions) replacing the background with different  
612 backgrounds based on the contextual manipulation conditions. The main manipulations per  
613 context type are as follows: (1) Full context: No manipulation, serving as the reference image with  
614 the object in a congruent context; (2) Incongruent context: Context swapped with a different  
615 (wrong) context; (3) No context: Context removed by swapping with gray pixels; (4) No object:  
616 Object removed by swapping with gray pixels; (5) Blurred context: Gaussian blur with kernel size  
617 2 applied on the context; (6) Blurred object: Gaussian blur with kernel size 2 applied on the object;  
618 (7) Blurred incongruent boundary: Gaussian blur with kernel size 2 applied on the object-  
619 incongruent context boundary; (8) Minimal context: All context apart from the smallest bounding  
620 box around the object is removed; (9) Jigsaw context: 25x25 pixel context patches randomly  
621 shuffled around the object; and (10) Textured context: Context swapped with texture generated  
622 with Portilla & Simoncelli method<sup>42</sup> (5 iterations) on the baseline image. Each of these context  
623 conditions was applied to 60 images, with 6 images per object category, resulting in a total of 600  
624 images in the imageset.  
625

### 626 **Low-Level Image Features**

627 For every image, we extracted a range of basic image features, such as object size, location and  
628 category, spectral mean and standard deviation(std), and contrast mean and standard deviation  
629 (std). The standard contrast metric for gray-scale images was used, calculated by the highest and  
630 lowest pixel values. The contrast standard deviation was derived from the pixel-wise standard  
631 deviation of the grayscale image. From the COCO object annotations, we determined the object  
632 size, represented in degrees of visual angle, as the fraction of the full image size (considering the  
633 full image was presented at 8 degrees) covered by the smallest bounding square around the  
634 object. The x and y coordinates, relative to the image, captured the object's central position. Using  
635 the Fast Fourier Transform (FFT), we transformed the image in the spectral domain, and noted  
636 its spectral mean and standard deviation.

## 637 **Subjects**

### 638 Human Participants

639 A total of 309 human subjects participated in the binary object discrimination tasks. Observers  
640 completed 5–10-min tasks through Amazon Mechanical Turk (MTurk), an online platform in which  
641 subjects could complete experiments for a payment of \$15 CAD/hour. We confirm that this  
642 experimental protocol involving human participants was approved by and in concordance with the  
643 guidelines of the York University Human Participants Review Subcommittee.

### 644 Non human primates

645  
646 The nonhuman subjects in our experiments were four adult male rhesus monkeys (*Macaca*  
647 *mulatta*). 2 of these monkeys (monkey M and monkey B), were trained with objects in congruent  
648 context and could perform the object discrimination tasks. The other 2 (monkey P, and monkey  
649 K) were naive to the discrimination task, and were only trained to passively fixate on the screen.  
650 All data were collected, and animal procedures were performed, in accordance with the NIH  
651 guidelines, the Massachusetts Institute of Technology Committee on Animal Care, and the  
652 guidelines of the Canadian Council on Animal Care on the use of laboratory animals and were  
653 also approved by the York University Animal Care Committee.  
654

## 655 **Behavioral testing**

### 656 **Primate behavioral testing**

#### 657 Humans active binary object discrimination task

658  
659 We collected large-scale psychophysical data from 309 subjects using Amazon Mechanical Turk  
660 (MTurk), an online crowdsourcing platform. The reliability of MTurk for psychophysical  
661 experiments has been previously validated by comparing online and in-lab results. Each trial  
662 began with a brief presentation (100 ms) of a sample image, selected from a set of 600 images.  
663 After a 100 ms blank gray screen, subjects were shown a choice screen displaying the target and  
664 distractor objects, similar to the procedure used in<sup>22,23</sup>. Subjects indicated their choice by touching  
665 the screen or clicking the mouse on the target object. No information regarding the sex of the  
666 participants was collected.

#### 667 Macaques active binary object discrimination task

668  
669 We measured monkey behavior from 2 male rhesus macaques. Images were presented on a 24-  
670 inch LCD monitor (1920 × 1080 at 60 Hz) positioned 42.5 cm in front of the animal. Monkeys were  
671 head fixed. Monkeys fixated a white cross (0.2°) for 300 ms to initiate a trial. The trial started with

672 the presentation of a sample image (from a set of 640 images) for 100 ms. This was followed by  
673 a blank gray screen for 100 ms, after which the choice screen was shown containing a standard  
674 image of the target object (the correct choice) and a standard image of the distractor object. The  
675 monkey was allowed to view freely the choice objects for up to 1500 ms and indicated its final  
676 choice by holding fixation over the selected object for 400 ms. Trials were aborted if gaze was not  
677 held within  $\pm 2^\circ$  of the central fixation dot during any point until the choice screen was shown. Prior  
678 to testing in the laboratory, monkeys were trained in their home-cages to perform the delayed  
679 match to sample tasks on the same object categories (but with a different set of images).

680  
681

## 682 ANN behavioral testing

683 We evaluated eighteen ANN models, on the exact images shown to the macaques and humans.  
684 We focused on publicly available pre-trained PyTorch model architectures that have  
685 demonstrated significant success in computer vision benchmarks. Table 1 lists the models used  
686 and their characteristics.

687  
688

Model	Architecture	Layer used
<b>Image classification models trained on ImageNet</b>		
AlexNet <sup>43</sup>	Generic CNN	features.12
VGG-19 <sup>44</sup>	Generic CNN	features.27
MobileNet-v2 <sup>45</sup>	Generic CNN	features.15
ResNet-18 <sup>46</sup>	Skip connections CNN	layer4.1
ResNet-50 <sup>46</sup>	Skip connections CNN	layer4.2
ResNet-101 <sup>46</sup>	Skip connections CNN	layer4.2
ResNet-152 <sup>46</sup>	Skip connections CNN	layer4.2
DenseNet-201 <sup>47</sup>	Skip connections CNN	features.transition3.pool
ConvNetXt Large <sup>48</sup>	Skip connections CNN	avgpool
GoogLeNet <sup>49</sup>	Inception block CNN	inception5b
Inception-v3 <sup>50</sup>	Inception block CNN	Mixed_7c
RegNetX 32GF <sup>51</sup>	Generic CNN	trunk_output.block3.block



		3-12.activation
ViT-b32 <sup>52</sup>	Transformer	encoder.layers.encoder_layer_11.ln_2
Swin-b <sup>53</sup>	Transformer	features.7.1.norm2
<b>Image memorability models trained on LaMem</b>		
MemNet <sup>54</sup>	Generic CNN	pool5
ResMem <sup>55</sup>	Skip connections CNN	features.layer4.2
<b>Object detection models trained on Microsoft COCO</b>		
FasterRCNN (ResNet50 backbone) <sup>56</sup>	Skip connections RCNN	backbone.body.layer4.2
RetinaNet (ResNet50 backbone) <sup>57</sup>	Skip connections RCNN	backbone.body.layer4.2

689

690 **Table 1. Summary of the ANN models used grouped by training objective.**

691

692 To make these pre-trained models compatible with our specific 10-way object recognition task,  
 693 we used the extracted features from each model for every stimulus, from the most IT-like layers  
 694 (chosen based on BrainScore if that data was available, otherwise the most reasonable  
 695 penultimate layer) shown in Table 1. To ensure consistency in results across the models, given  
 696 the varying layer sizes for each, we standardized the dimension for every model down to 3,000  
 697 features. This was done by using Gaussian random projection with 3,000 components to project  
 698 the full extracted features space on a randomly generated linear subspace in such a way that  
 699 distances between the points are nearly preserved. We trained a multiclass SVM classifier using  
 700 these scaled features (standard scaling) to calculate the cross validated probabilities for each  
 701 object class (using 10 one-vs-all classifiers, 5 folds, 10 repetitions), mimicking the subjects' active  
 702 binary object discrimination task. All behavioral predictions from the decoder were for images  
 703 where the object was not seen in any phase of the model training regardless of the surrounding  
 704 context.

## 705 **Electrophysiological recording and data preprocessing**

### 706 **Passive Fixation Task**

707 During the passive viewing task, monkeys fixated a white cross (0.2°) for 300 ms to initiate a trial.  
 708 We then presented a sequence of 5 to 10 images, each ON for 100 ms followed by a 100 ms gray  
 709 (background, 'OFF') blank screen. This was followed by fluid (water) reward and an inter-trial

710 interval of 500 ms, followed by the next sequence. The animals (n = 2, male rhesus macaques)  
711 used in the passive fixation experiments study can be classified as “categorization task naive”,  
712 since they have not been explicitly trained to perform any object categorization tasks.  
713

## 714 Eye Tracking

715 We monitored eye movements using video eye tracking (SR Research EyeLink 1000). Using  
716 operant conditioning and water reward, our 2 subjects were trained to fixate a central white square  
717 (0.2°) within a square fixation window that ranged from ±2°. At the start of each behavioral  
718 session, monkeys performed an eye-tracking calibration task by making a saccade to a range of  
719 spatial targets and maintaining fixation for 500 ms. Calibration was repeated if drift was noticed  
720 over the course of the session.

721 Real-time eye-tracking was employed to ensure that eye jitter did not exceed ±2°, otherwise the  
722 trial was aborted, and data discarded. Stimulus display and reward control were managed using  
723 the MWorks Software (<https://mworks.github.io>).

## 724 Data Analyses

### 725 Behavioral Metrics

726 We developed two behavioral metrics, the hit rate at context level - B.C1 and more fine grained  
727 image level - B.I1 (as introduced in<sup>22</sup>). We obtained a biological or artificial signature for each  
728 system by applying each metric to its behavioral accuracies per image averaged across all trials.  
729 The one-versus-all context-level performance metric (B.C1) estimates the discriminability of all  
730 images of context category  $c$ , essentially pooling the accuracies across all images of context type  
731  $c$  and all object/distractor pairs within. Because we tested 10 context categories, the resulting  
732 B.C1 signature has 10 independent values.  
733

734 The one-versus-all image-level performance metric (B.I1) estimates the discriminability of each  
735 image, pooling across all distractors. Because we have an image test set of 600 images (60 per  
736 object, see above), the resulting B.I1 signature has 600 independent values. Given an image  $i$  of  
737 object category  $o$ , and all nine distractor objects ( $d \neq o$ ), we computed the average performance  
738 per image as:  $I1_i^o = (\sum_{d=1}^{10} PC_i^{o,d \neq o}) \div 9$ , where  $PC$  (percent correct) is the fraction of correct  
739 responses for the binary task between object categories  $o$  and  $d$ . Considering every image  $i_c$  of  
740 context type  $c$ , the B.C1 performance for each context type is the mean across the performance  
741 of all images (60 per context type):  $C1_c = (\sum_{i_c=1}^{60} I1_{i_c}^c) \div 60$

### 742 Human-monkey shared behavioral variance

743 To quantify the behavioral pattern similarity at a context and image level across humans and  
744 monkeys, we calculated the percent of shared behavioral variance (SV) for both signatures. The  
745 SV is obtained as the square of the correlation (Pearson’s R) of the pooled humans and pooled  
746 monkeys behavioral signature, corrected by the human and monkey signature internal

747 consistency. This was repeated 20 times choosing 600 images with repetition (bootstrap). The  
748 ceiling estimates in Fig 7B and 7C show the full range for the 20 bootstrap values for C1 and I1  
749 respectively.

## 750 Partial correlation analysis

751 To estimate the fraction of shared human-monkey variance that is explained by the models  
752 (including the Neural model), we calculated the partial correlation for the pooled humans and  
753 monkeys population - while controlling individually for each model. The partial correlation gives  
754 us the fraction of the primate shared variance that is independent of the model variance. The  
755 percentage of shared human-monkey variance explained by the model is then given by the  
756 formula:

757  $(R^2 - R_p^2) / R^2$ , where  $R^2$  is the human-monkey Pearson correlation,  $R_p$  is the human-monkey  
758 partial correlation, while controlling for the model (calculated as the product of the residuals of the  
759 model predictions).

760  
761 The neural correction of the partial correlation is done by fitting a sigmoid extrapolation to an  
762 infinite number of neural trials (see **Fig 6D** inset for 122 neurons). Both  $R^2$  and  $R_p$  are corrected  
763 by the (Spearman-Brown corrected) human and monkey split-half reliabilities  $\sqrt{\rho_h * \rho_m}$ ,  
764 however, due to the normalization by  $R^2$ , we did not need to account for the human or monkey  
765 noise.

## 766 Internal consistency

767 The reliability of each system (pooled human, monkey, and IT population) was assessed by  
768 calculating the trial split-half Spearman-Brown corrected correlation. For the pooled humans or  
769 monkeys, this was done by splitting all the accuracy trials per image in two halves, taking the  
770 mean for both halves (for each image), and computing the corrected Spearman correlation across  
771 all images for the two halves, repeated 100 times with different trial splits. The internal consistency  
772 for the decoding accuracy of the neural data was computed by calculating the decoding accuracy  
773 for each mean half of the neural trials and correlating the two obtained accuracies (across all  
774 images). The ceiling estimates shown in Fig 5B and 5D are the pooled human internal  
775 consistency, showing the full range of values (min-max).

776

## 777 Statistical Analyses

778 For each statistical analysis, we first tested the normality of the data. We used the Lilliefors test  
779 assuming normal distribution, with a threshold 5% (normal distribution:  $p > 0.05$ ).

780

781 To test for statistical significance with a normal distribution of the data, a paired (monkey-human  
782 comparison) independent (comparing contextual variations) T-test was performed. This is a test  
783 for the null hypothesis that two samples have identical average (expected) values. The t(DOF)-  
784 statistic value quantifies the difference between the arithmetic means of the two samples. It is

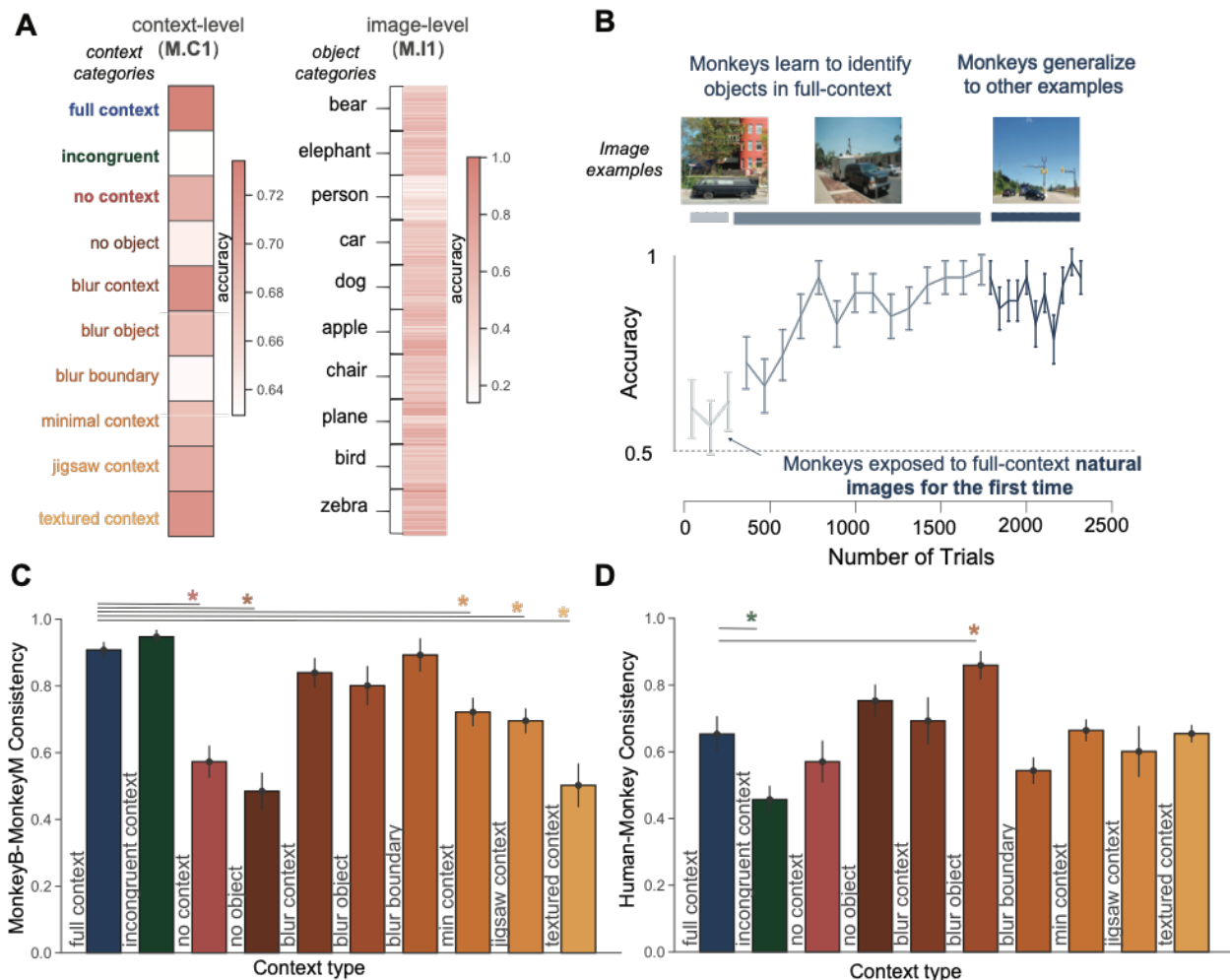
785 calculated as the mean of the difference of the two variables, divided by the standard error. The  
 786 p-value quantifies the probability of observing as or more extreme values assuming the null  
 787 hypothesis, that the samples are drawn from populations with the same population means, is true.  
 788

789 A Wilcoxon signed rank (paired variables) or ranksum (independent variables) was performed in  
 790 case of a non-normal data distribution. The null hypothesis is that two (paired or independent  
 791 respectively) samples come from the same distribution. In particular, it tests whether the  
 792 distribution of the differences  $x - y$  is symmetric about zero. It is a non-parametric version of the T-  
 793 test.

794

795 We chose the threshold 5% ( $p < 0.05$ ) to reject the null hypothesis for all tests.

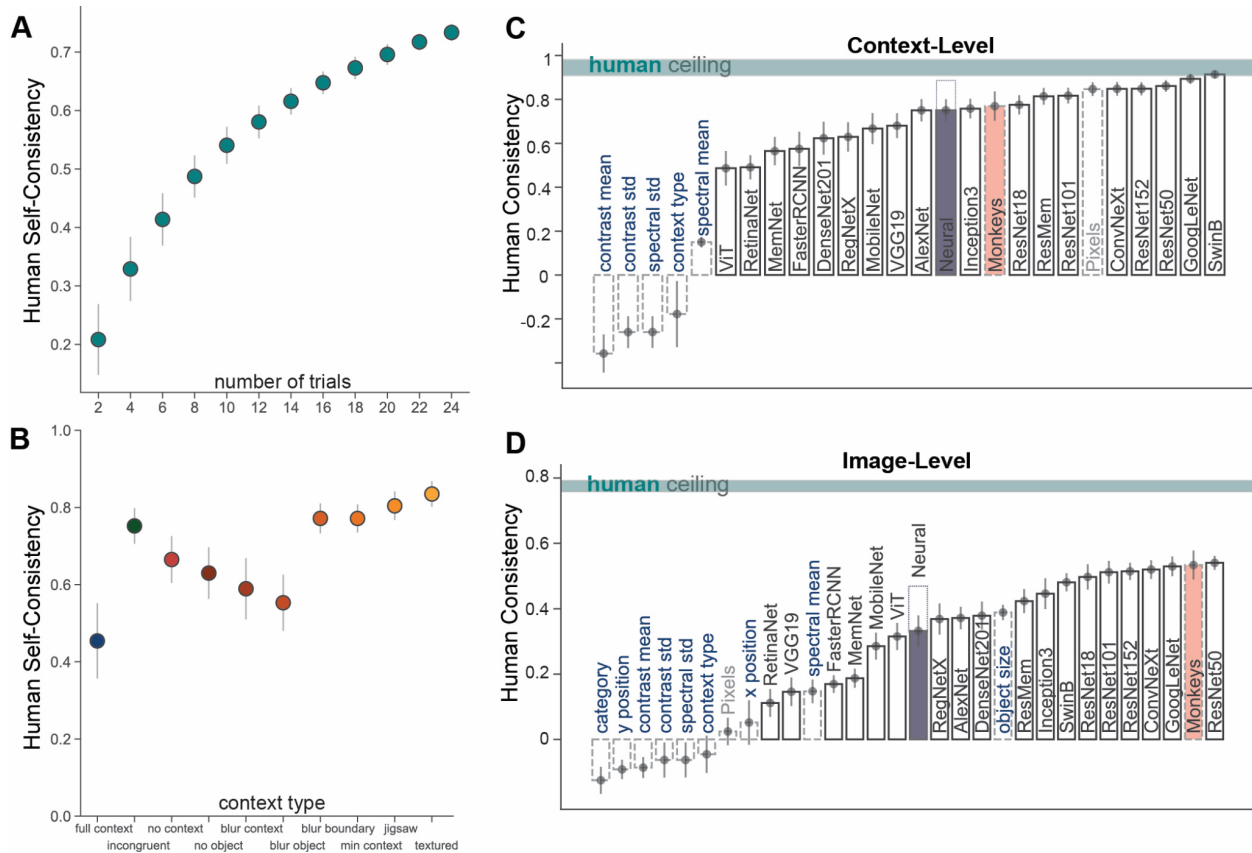
## 796 Supplementary Figures



797

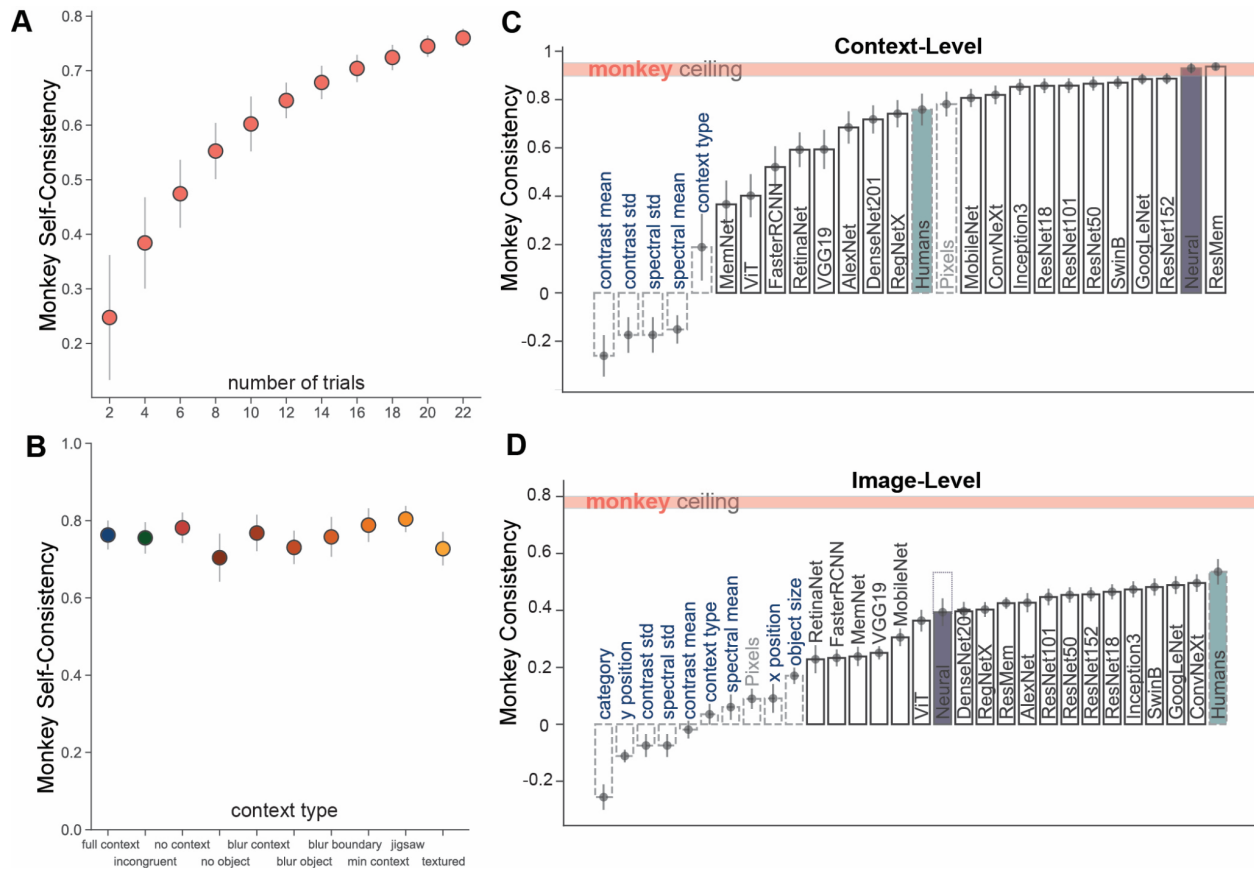
798 **Fig S1. Monkeys' behavioral signatures, context learning curve and image-level consistency per**  
 799 **context type.** **A.** Context-level (B.C1) and image-level (B.I1) behavioral signatures defined from the  
 800 monkeys' accuracy. Same as **Fig 2A**, but for the pooled monkeys. **B.** Training process for one macaque  
 801 in the task shown in **Fig 1B** (chance = 0.5). The monkey, who was previously not exposed to images with

802 objects in context, has low starting performance for full-context images (light blue curve). However, the  
 803 monkey quickly learns to recognize images in full context (blue curve). Furthermore, this ability  
 804 generalizes to new images (dark blue curve). Error-bars show the standard deviation across object  
 805 categories. **C.** Each bar shows the corrected Pearson correlation between the two monkeys (monkey M  
 806 and monkey B) for the images of a context category. Error-bars are standard errors across ten  
 807 subsamples of images within a context category. Statistics are shown for full context compared to each  
 808 other context variation (\* denotes t-test,  $p < 0.05$ ). **D.** Similar to C but shows the corrected Pearson  
 809 correlation between the pooled two monkeys and the pooled human population.  
 810  
 811



812 **Fig S2. ANNs and IT population are consistent with human behavior at a context-level and only**  
 813 **partially at an image-level.** **A.** Human population self-consistency: Spearman-Brown corrected split-half  
 814 correlation for increasing number of trials. The mean correlation for 100 different splits for each subset of  
 815 trials, with standard deviation across the splits. **B.** Human population self-consistency (Spearman-Brown  
 816 corrected split-half correlation) using all 24 trials per image, for images grouped by context category. The  
 817 mean human internal reliability across 100 splits with standard deviation for each context category (color  
 818 coding same as **Fig 1A**, and labeled on the x axis). **C.** The human consistency - Pearson R with the low-  
 819 level features, ANN models and the Neural model at a context level. The low-level features are shown with  
 820 light blue text (dashed bars). The mean human internal behavioral consistency ceiling is shown in gray,  
 821 with standard deviation across different image subsamples. We show the noise corrected (by the split-half  
 822 decoding consistency) neural consistency (purple bar), when using all the recorded reliable neural  
 823 responses (122), the extrapolated consistency (to 4537 neurons, as in **Fig 6E**) is shown with a dashed bar  
 824 on top. The noise corrected (by the monkey internal reliability) consistency with monkeys is shown in coral.  
 825 **D.** Similar to C but showing the consistency at an image-level.  
 826  
 827

828



829

830

831

832

833

834

835

836

837

838

839

840

841

842

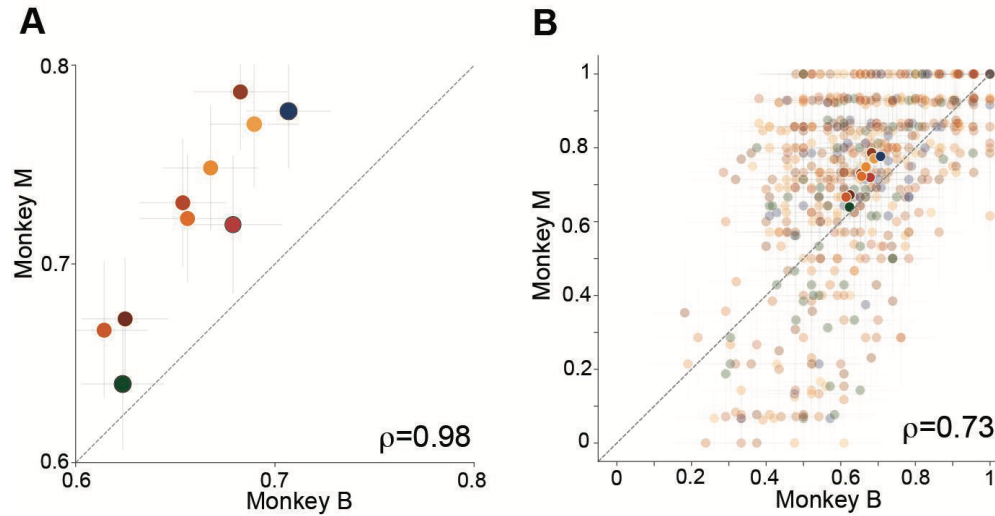
843

844

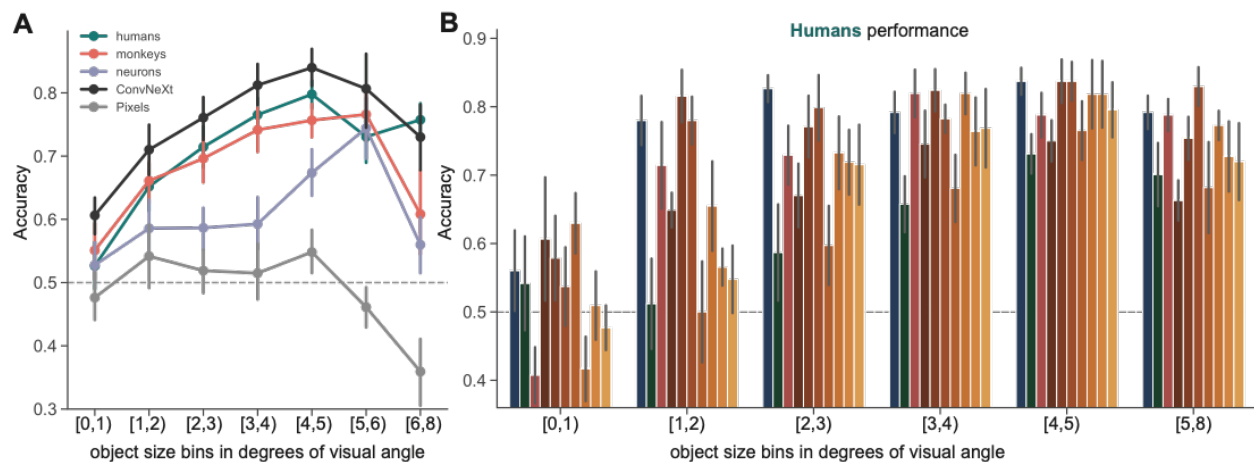
845

846

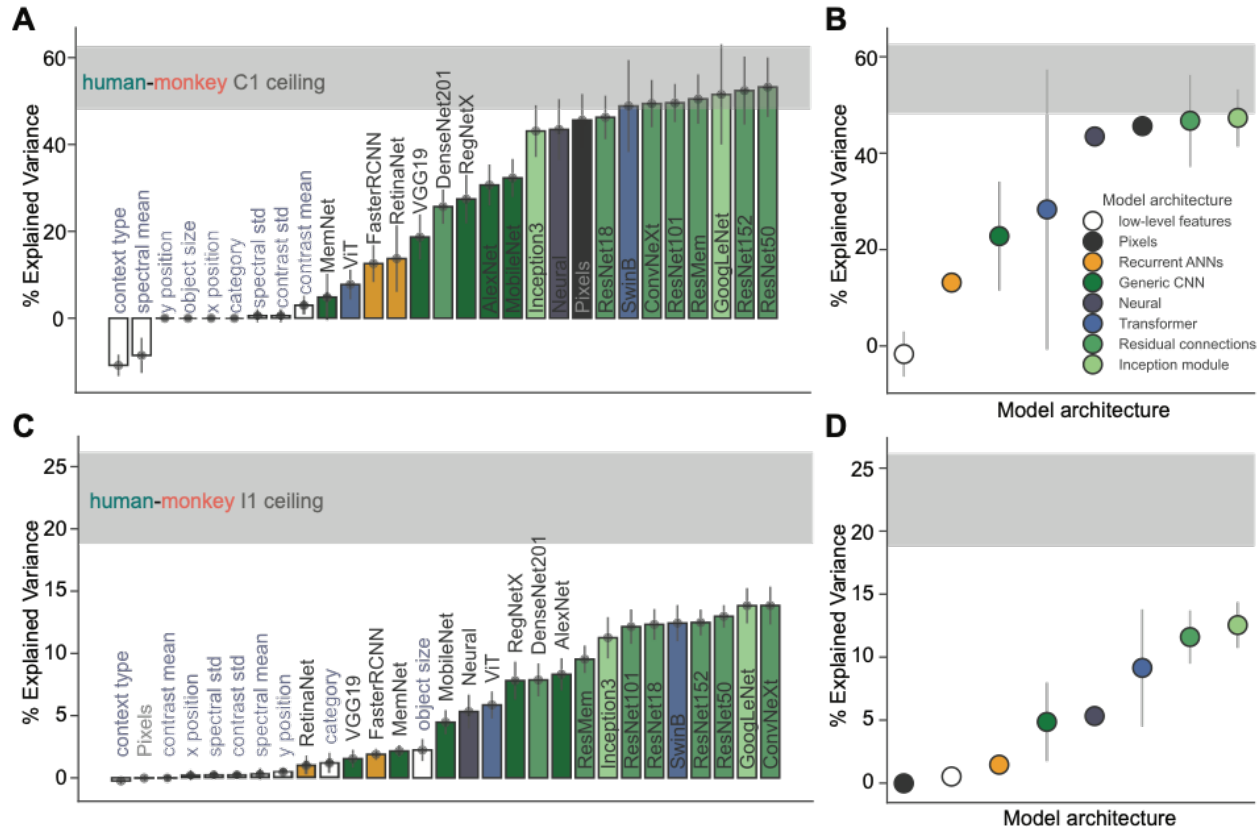
**Fig S3. ANNs and IT population are consistent with monkey behavior at a context-level and only partially at an image-level.** **A.** Monkey self-consistency: Spearman-Brown corrected split-half correlation for increasing number of trials. The mean correlation for 100 different splits for each subset of trials, with standard deviation across the splits. **B.** Monkey self-consistency (Spearman-Brown corrected split-half correlation) using all 22 trials per image, for images grouped by context category. The mean monkey internal reliability across 100 splits with standard deviation for each context category (color coding same as **Fig 1A**, and labeled on the x axis). **C.** The monkey consistency - Pearson R with the low-level features, ANN models and the Neural model at a context level. The low-level features are shown with light blue text (dashed bars). The mean monkey internal behavioral consistency ceiling is shown in gray, with standard deviation across different image subsamples. We are noting the noise corrected (by the split-half decoding consistency) neural consistency (purple) when using all the recorded reliable neural responses (122), the extrapolated consistency (to 4537 neurons, as in **Fig 6E**) is shown with a dashed bar on top (no extrapolation needed for context-level as the consistency is already within the monkey ceiling). The noise corrected (by the human internal reliability) consistency with humans is shown in teal. **D.** Similar to C but showing the consistency at an image-level.



847  
848 **Fig S4. Two monkeys show similar (but not identical) context-driven behavioral changes.** **A.** Context-  
849 level (B.C1) correlation between the two monkeys. Each point represents the mean accuracy for a  
850 contextual variation with standard error across images of that context type (colors as in **Fig 1A**, Monkey B  
851 mean  $0.66 \pm 0.03$ , Monkey M mean  $0.72 \pm 0.05$ ). The three main context types: full (blue), incongruent (green)  
852 and no context (red), are shown with a black stride. The value  $\rho$  indicates the noise corrected correlation  
853 coefficient (Pearson R). **B.** Image-level correlation (B.I1) for the two macaques, each low opacity point  
854 shows the performance(accuracy) for an image with standard error across trials, the higher opacity points  
855 are the B.C1 mean (from A), colors map to context types as defined in **Fig 1A**.  
856



857  
858 **Fig S5. Object size predicts primate and ANN average accuracy.** **A.** Average image-level accuracy (for  
859 all images) grouped in bins based on the object size (in degrees of visual angle, the full image is 8 degrees),  
860 with standard error across images in each bin. The performance is shown for humans (teal), pooled  
861 monkeys (coral), IT population (purple), Pixels (gray) and ConvNeXt (black) - the best model explaining  
862 the highest fraction of the human-monkey shared image-level behavioral variance (see Fig 7C). Chance  
863 level accuracy (0.5) is noted with the dashed gray line. The size bins are labeled with the minimum size  
864 and max size (not included) of the bin (eg. the size bin [1,2) contains all images where the object size is  
865 greater or equal to 1 degree and smaller than 2 degrees). **B.** Human accuracy for each object size grouped  
866 by context category (color coding for context type from Fig 1), with standard error across images in each  
867 bin. The human data in part A (teal) is the average over all context conditions shown in part B.

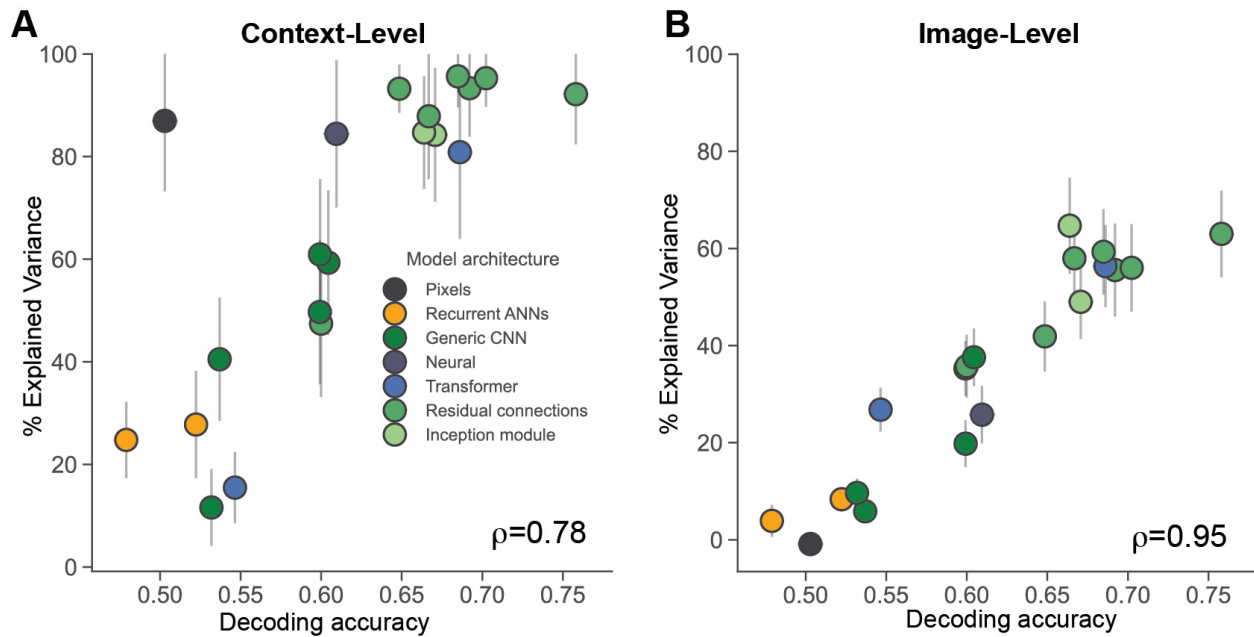


868  
869  
870  
871  
872  
873  
874  
875

**Fig S6. ANN model architecture effects on explained human-monkey shared variance.** **A.** Shows the same models as **Fig 7**, but the bars are color coded for the model architecture (see legend in **B**). Green is used for CNNs (with subgroups: models with Inception modules and residual connections), blue for visual transformers and yellow for recurrent neural networks. **B.** The model performance, grouped by model architecture, with the standard deviation across models. CNN models with residual or inception blocks share the most of the shared human-monkey variance. **C.** Same as **A** but for image-level EV. **D.** Same as **B** but for image-level EV.



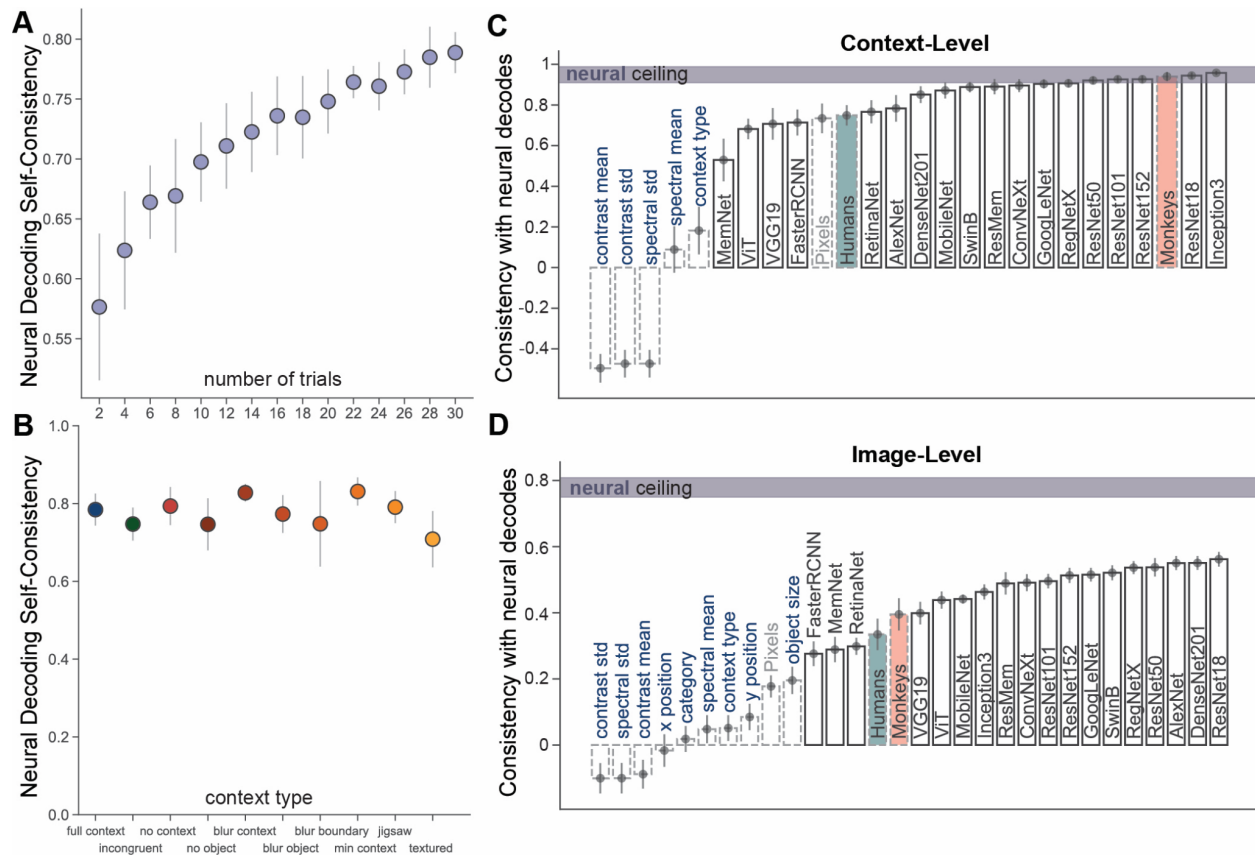
876  
877



878  
879

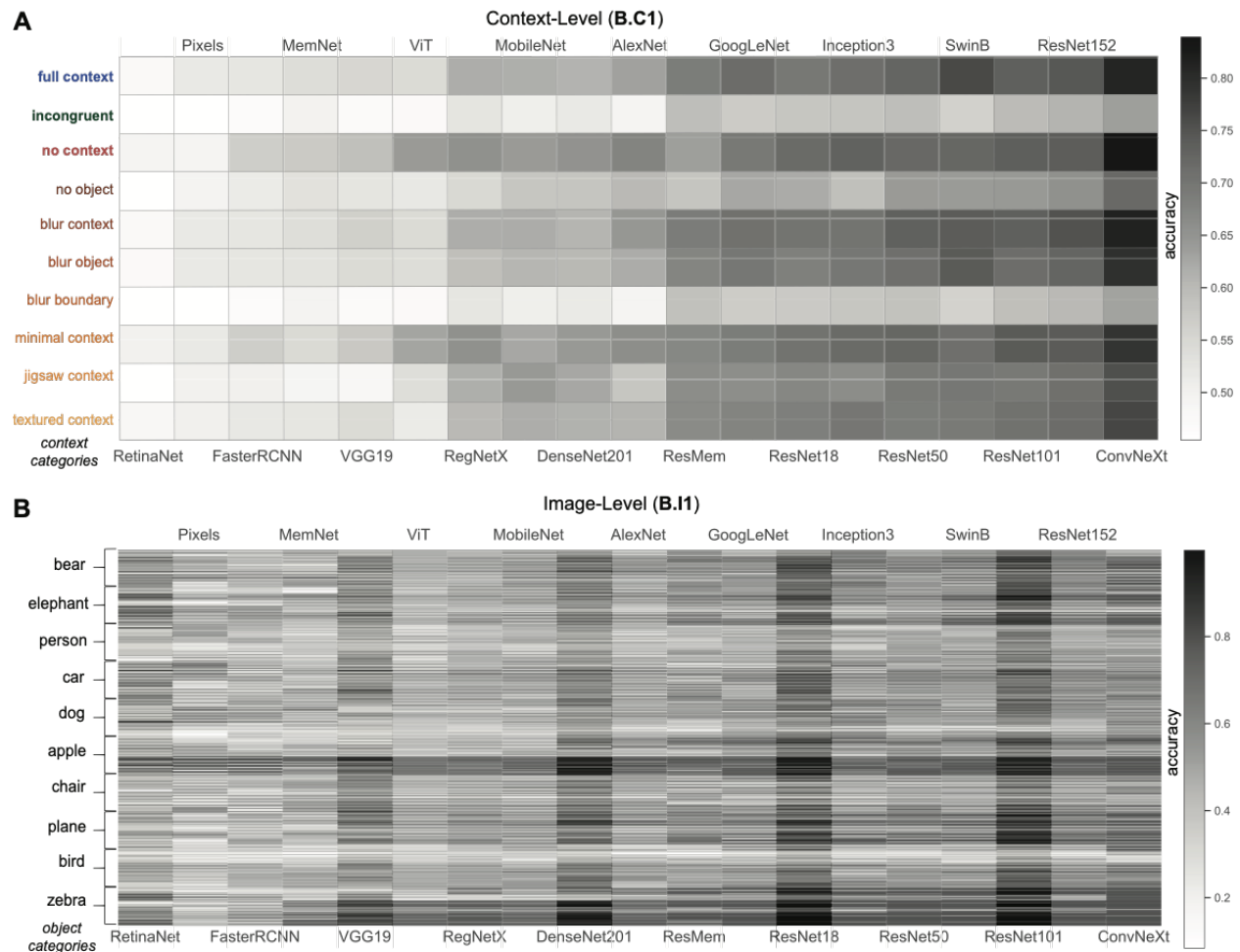
880 **Fig S7. The models' decoding accuracy predicts the fraction of explained shared human-monkey**  
881 **variance. A.** The percent of variance explained by each model from the shared human-monkey variance  
882 at a context-level as a function of the mean image decoding accuracy for each model. Each point is a  
883 different ANN model (Neural model in purple, Pixels in black), color coded by model architecture (as in **Fig**  
884 **S6**, see legend and Table 1). The y axis shows the normalized % EV (by the human-monkey shared  
885 variance ceiling) with standard deviation across image subsamples for each model. The x axis shows the  
886 mean decoding accuracy across all images (with standard error, x error bars are smaller than the points).  
887  $\rho$  notes the Pearson correlation between the accuracy and EV across models. **B.** Same as A but showing  
888 the %EV as a function of the decoding accuracy at an image-level.

889  
890  
891  
892



893  
 894 **Fig S8. ANNs and primates are consistent with context-naive IT decoded behavior at a context-level**  
 895 **and only partially at an image-level.** **A.** Neural decoding accuracy self-consistency: Spearman-Brown  
 896 corrected split-half correlation of the decoding accuracy across all images for increasing number of neural  
 897 trials used for decoding ( $n=122$  neural sites used). The mean correlation for 20 different splits for each  
 898 subset of trials, with standard deviation across the splits. **B.** Neural decoding accuracy self-consistency  
 899 (Spearman-Brown corrected split-half correlation) using all 30 neural trials per image, per neuron (for the  
 900 122 neurons), for images grouped by context category. The mean neural decoding internal reliability across  
 901 20 splits with standard deviation for each context category (color coding same as **Fig 1A**, and labeled on  
 902 the x axis). **C.** The consistency with neural decode based predictions (Pearson R) with the low-level  
 903 features, ANN models, pooled humans and monkeys behavioral accuracy at a context-level. The low-level  
 904 features are shown with light blue text (dashed bars). The mean neural internal behavioral consistency  
 905 ceiling is shown in gray (split-half decoding reliability), with standard deviation across different image  
 906 subsamples. We are noting the (internal reliability corrected) human (teal) and monkey (coral) consistency.  
 907 **D.** Similar to C but showing the consistency at an image-level.

908  
 909



910  
 911 **Fig S9. Pixels and ANN models behavioral signatures.** **A.** Context-level behavioral signature (B.C1) for  
 912 each model (column), sorted by their overall average decoding accuracy - least accurate (left) to most  
 913 accurate (right). The average accuracy is shown for each context type (row) with the color indicating the  
 914 accuracy (increasing from white to black, see colorbar on the right). **B.** Similar to A but showing the image-  
 915 level signature for each model (B.I1), each line represents the image accuracy averaged across trials and  
 916 distractors, with the rows sorted and grouped by object category (see on the left).  
 917  
 918

## 919 Reference

- 920  
 921 1. Riesenhuber, M. & Poggio, T. Hierarchical models of object recognition in cortex. *Nature*  
 922 *Neuroscience* **2**, 1019–1025 (1999).  
 923 2. Yamins, D. L. & DiCarlo, J. J. Eight open questions in the computational modeling of higher  
 924 sensory cortex. *Current Opinion in Neurobiology* **37**, 114–120 (2016).  
 925 3. Kar, K. & DiCarlo, J. J. The Quest for an Integrated Set of Neural Mechanisms Underlying

- 926 Object Recognition in Primates. Preprint at <http://arxiv.org/abs/2312.05956> (2023).
- 927 4. Hung, C. P., Kreiman, G., Poggio, T. & DiCarlo, J. J. Fast Readout of Object Identity from  
928 Macaque Inferior Temporal Cortex. *Science* **310**, 863–866 (2005).
- 929 5. Majaj, N. J., Hong, H., Solomon, E. A. & DiCarlo, J. J. Simple Learned Weighted Sums of  
930 Inferior Temporal Neuronal Firing Rates Accurately Predict Human Core Object Recognition  
931 Performance. *Journal of Neuroscience* **35**, 13402–13418 (2015).
- 932 6. Liu, H., Agam, Y., Madsen, J. R. & Kreiman, G. Timing, Timing, Timing: Fast Decoding of  
933 Object Information from Intracranial Field Potentials in Human Visual Cortex. *Neuron* **62**,  
934 281–290 (2009).
- 935 7. Oliva, A. & Torralba, A. The role of context in object recognition. *Trends in Cognitive*  
936 *Sciences* **11**, 520–527 (2007).
- 937 8. Lauer, T., Cornelissen, T. H. W., Draschkow, D., Willenbockel, V. & Vö, M. L.-H. The role of  
938 scene summary statistics in object recognition. *Sci Rep* **8**, 14666 (2018).
- 939 9. Vö, M. L.-H. The meaning and structure of scenes. *Vision Research* **181**, 10–20 (2021).
- 940 10. Bar, M. & Aminoff, E. Cortical Analysis of Visual Context. *Neuron* **38**, 347–358 (2003).
- 941 11. Bar, M. Visual objects in context. *Nature Reviews. Neuroscience* **5**, 617–629 (2004).
- 942 12. Bar, M. *et al.* Top-down facilitation of visual recognition. *Proceedings of the National*  
943 *Academy of Sciences of the United States of America* **103**, 449–454 (2006).
- 944 13. Zhang, M., Tseng, C. & Kreiman, G. Putting visual object recognition in context. (2019)  
945 doi:10.48550/ARXIV.1911.07349.
- 946 14. Bomatter, P. *et al.* When Pigs Fly: Contextual Reasoning in Synthetic and Natural  
947 Scenes. *arXiv:2104.02215 [cs]* (2021).
- 948 15. Adesnik, H., Bruns, W., Taniguchi, H., Huang, Z. J. & Scanziani, M. A neural circuit for  
949 spatial summation in visual cortex. *Nature* **490**, 226–231 (2012).
- 950 16. Keller, A. J. *et al.* A Disinhibitory Circuit for Contextual Modulation in Primary Visual  
951 Cortex. *Neuron* **108**, 1181–1193.e8 (2020).

- 952 17. Mély, D. A., Linsley, D. & Serre, T. Complementary surrounds explain diverse contextual  
953 phenomena across visual modalities. *Psychological Review* **125**, 769–784 (2018).
- 954 18. Henry, C. A. & Kohn, A. Spatial contextual effects in primary visual cortex limit feature  
955 representation under crowding. *Nature Communications* **11**, 1687 (2020).
- 956 19. Fisher, T. G., Alitto, H. J. & Usrey, W. M. Retinal and Nonretinal Contributions to  
957 Extraclassical Surround Suppression in the Lateral Geniculate Nucleus. *The Journal of*  
958 *Neuroscience: The Official Journal of the Society for Neuroscience* **37**, 226–235 (2017).
- 959 20. Connor, C. E., Brincat, S. L. & Pasupathy, A. Transformation of shape information in the  
960 ventral pathway. *Current Opinion in Neurobiology* **17**, 140–147 (2007).
- 961 21. Rajalingham, R., Schmidt, K. & DiCarlo, J. J. Comparison of Object Recognition  
962 Behavior in Human and Monkey. *J. Neurosci.* **35**, 12127–12136 (2015).
- 963 22. Rajalingham, R. *et al.* Large-Scale, High-Resolution Comparison of the Core Visual  
964 Object Recognition Behavior of Humans, Monkeys, and State-of-the-Art Deep Artificial  
965 Neural Networks. *J. Neurosci.* **38**, 7255–7269 (2018).
- 966 23. Kar, K., Kubilius, J., Schmidt, K., Issa, E. B. & DiCarlo, J. J. Evidence that recurrent  
967 circuits are critical to the ventral stream’s execution of core object recognition behavior. *Nat*  
968 *Neurosci* **22**, 974–983 (2019).
- 969 24. Kar, K. & DiCarlo, J. J. Fast Recurrent Processing via Ventrolateral Prefrontal Cortex Is  
970 Needed by the Primate Ventral Stream for Robust Core Visual Object Recognition. *Neuron*  
971 **109**, 164-176.e5 (2021).
- 972 25. Hong, H., Yamins, D. L. K., Majaj, N. J. & DiCarlo, J. J. Explicit information for category-  
973 orthogonal object properties increases along the ventral stream. *Nature Neuroscience* **19**,  
974 613–622 (2016).
- 975 26. McKee, J. L., Riesenhuber, M., Miller, E. K. & Freedman, D. J. Task Dependence of  
976 Visual and Category Representations in Prefrontal and Inferior Temporal Cortices. *J.*  
977 *Neurosci.* **34**, 16065–16075 (2014).

- 978 27. Yamins, D. L. K. *et al.* Performance-optimized hierarchical models predict neural  
979 responses in higher visual cortex. *Proceedings of the National Academy of Sciences of the*  
980 *United States of America* **111**, 8619–8624 (2014).
- 981 28. Schrimpf, M. *et al.* *Brain-Score: Which Artificial Neural Network for Object Recognition Is*  
982 *Most Brain-Like?* <http://biorxiv.org/lookup/doi/10.1101/407007> (2018) doi:10.1101/407007.
- 983 29. Cadena, S. A. *et al.* Deep convolutional models improve predictions of macaque V1  
984 responses to natural images. *PLoS Comput Biol* **15**, e1006897 (2019).
- 985 30. Laskar, M. N. U., Sanchez Giraldo, L. G. & Schwartz, O. Deep neural networks capture  
986 texture sensitivity in V2. *Journal of Vision* **20**, 21 (2020).
- 987 31. Pospisil, D. A., Pasupathy, A. & Bair, W. ‘Artiphysiology’ reveals V4-like shape tuning in  
988 a deep network trained for image classification. *eLife* **7**, e38242 (2018).
- 989 32. Bashivan, P., Kar, K. & DiCarlo, J. J. Neural population control via deep image  
990 synthesis. *Science* **364**, eaav9436 (2019).
- 991 33. Kubilius, J. *et al.* Brain-Like Object Recognition with High-Performing Shallow Recurrent  
992 ANNs. Preprint at <http://arxiv.org/abs/1909.06161> (2019).
- 993 34. Nayebi, A. *et al.* Recurrent Connections in the Primate Ventral Visual Stream Mediate a  
994 Trade-Off Between Task Performance and Network Size During Core Object Recognition.  
995 *Neural Computation* **34**, 1652–1675 (2022).
- 996 35. Geirhos, R. *et al.* ImageNet-trained CNNs are biased towards texture; increasing shape  
997 bias improves accuracy and robustness. *arXiv:1811.12231 [cs, q-bio, stat]* (2019).
- 998 36. Geirhos, R. *et al.* Generalisation in humans and deep neural networks.  
999 *arXiv:1808.08750 [cs, q-bio, stat]* (2020).
- 1000 37. Paulun, V. C., Zheng, K. & Kar, K. Distributed population activity in the macaque inferior  
1001 temporal cortex but not current deep neural networks predict the ponzo illusion. *Journal of*  
1002 *Vision* **22**, 3354 (2022).
- 1003 38. Rossion, B. & Taubert, J. What can we learn about human individual face recognition

- 1004 from experimental studies in monkeys? *Vision Research* **157**, 142–158 (2019).
- 1005 39. Fize, D., Cauchoix, M. & Fabre-Thorpe, M. Humans and monkeys share visual  
1006 representations. *Proc. Natl. Acad. Sci. U.S.A.* **108**, 7635–7640 (2011).
- 1007 40. Tsao, D. Y., Schweers, N., Moeller, S. & Freiwald, W. A. Patches of face-selective cortex  
1008 in the macaque frontal lobe. *Nat Neurosci* **11**, 877–879 (2008).
- 1009 41. Marr, D. & Poggio, T. From understanding computation to understanding neural circuitry.  
1010 (1976).
- 1011 42. Portilla, J. A Parametric Texture Model Based on Joint Statistics of Complex Wavelet  
1012 Coefficients. *International Journal of Computer Vision* **40**, 49–70 (2000).
- 1013 43. Krizhevsky, A., Sutskever, I. & Hinton, G. E. ImageNet Classification with Deep  
1014 Convolutional Neural Networks. in *Advances in Neural Information Processing Systems* (eds.  
1015 Pereira, F., Burges, C. J. C., Bottou, L. & Weinberger, K. Q.) vol. 25 (Curran Associates, Inc.,  
1016 2012).
- 1017 44. Simonyan, K. & Zisserman, A. Very Deep Convolutional Networks for Large-Scale  
1018 Image Recognition. *arXiv:1409.1556 [cs]* (2015).
- 1019 45. Howard, A. G. *et al.* MobileNets: Efficient Convolutional Neural Networks for Mobile  
1020 Vision Applications. Preprint at <https://doi.org/10.48550/arXiv.1704.04861> (2017).
- 1021 46. He, K., Zhang, X., Ren, S. & Sun, J. Deep Residual Learning for Image Recognition.  
1022 Preprint at <https://doi.org/10.48550/arXiv.1512.03385> (2015).
- 1023 47. Huang, G., Liu, Z., van der Maaten, L. & Weinberger, K. Q. Densely Connected  
1024 Convolutional Networks. Preprint at <https://doi.org/10.48550/arXiv.1608.06993> (2018).
- 1025 48. Liu, Z. *et al.* A ConvNet for the 2020s. Preprint at  
1026 <https://doi.org/10.48550/arXiv.2201.03545> (2022).
- 1027 49. Szegedy, C. *et al.* Going Deeper with Convolutions. Preprint at  
1028 <https://doi.org/10.48550/arXiv.1409.4842> (2014).
- 1029 50. Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J. & Wojna, Z. Rethinking the Inception

- 1030 Architecture for Computer Vision. Preprint at <https://doi.org/10.48550/arXiv.1512.00567>  
1031 (2015).
- 1032 51. Radosavovic, I., Kosaraju, R. P., Girshick, R., He, K. & Dollár, P. Designing Network  
1033 Design Spaces. Preprint at <https://doi.org/10.48550/arXiv.2003.13678> (2020).
- 1034 52. Dosovitskiy, A. *et al.* An Image is Worth 16x16 Words: Transformers for Image  
1035 Recognition at Scale. Preprint at <https://doi.org/10.48550/arXiv.2010.11929> (2021).
- 1036 53. Liu, Z. *et al.* Swin Transformer: Hierarchical Vision Transformer using Shifted Windows.  
1037 Preprint at <https://doi.org/10.48550/arXiv.2103.14030> (2021).
- 1038 54. Tai, Y., Yang, J., Liu, X. & Xu, C. MemNet: A Persistent Memory Network for Image  
1039 Restoration. Preprint at <https://doi.org/10.48550/arXiv.1708.02209> (2017).
- 1040 55. Yang, Z. *et al.* ResMem: Learn what you can and memorize the rest. Preprint at  
1041 <https://doi.org/10.48550/arXiv.2302.01576> (2023).
- 1042 56. Ren, S., He, K., Girshick, R. & Sun, J. Faster R-CNN: Towards Real-Time Object  
1043 Detection with Region Proposal Networks. Preprint at  
1044 <https://doi.org/10.48550/arXiv.1506.01497> (2016).
- 1045 57. Lin, T.-Y., Goyal, P., Girshick, R., He, K. & Dollár, P. Focal Loss for Dense Object  
1046 Detection. Preprint at <https://doi.org/10.48550/arXiv.1708.02002> (2018).