

# A battery of image classification challenges reveals shared and distinct object categorization behavior across monkeys, humans, and deep networks

Han Zhang<sup>1, 2</sup>, Zhihao Zheng<sup>1</sup>, Jiaqi Hu<sup>1, 2</sup>, Qiao Wang<sup>1</sup>, Mengya Xu<sup>1</sup>, Zhaojiayi Zhou<sup>1</sup>, Zixuan Li<sup>1, 2</sup>, and Gouki Okazawa<sup>1, 2, \*</sup>

<sup>1</sup>Institute of Neuroscience, Key Laboratory of Brain Cognition and Brain-inspired Intelligence Technology, Center for Excellence in Brain Science and Intelligence Technology, International Center for Primate Brain Research, Chinese Academy of Sciences, Shanghai, China

<sup>2</sup>University of Chinese Academy of Sciences, China

\*Correspondence: okazawa@ion.ac.cn

## Abstract

Humans categorize objects at multiple levels of abstraction—animate versus inanimate, big versus small, and many other attributes. Despite its apparent challenge, the advent of deep neural networks (DNNs) has demonstrated that complex visual processing alone can support such classification without language or human-specific knowledge. This raises a natural question: to what extent can non-human primates, without language, perform such categorization? Although basic object-recognition behavior in monkeys such as similarity judgment has been extensively studied, their ability to classify objects across diverse rules remains poorly characterized. Here, we developed a task paradigm that enabled us to train monkeys on a large battery of binary classification tasks using natural object images, spanning more than 10 rules, such as animate versus inanimate, natural versus man-made objects, and mammalian versus non-mammalian animals. Monkeys acquired each rule in a few days, generalized the learned rules to new images, and exhibited error patterns consistent with human judgments. At the same time, their classification performance correlated more strongly with that of visual DNNs trained without language input, whereas human performance was better explained by language-informed DNNs. These results provide an important benchmark for the capacity of biological neural networks to perform image classification without language and human-specific knowledge.

## Introduction

Humans can readily judge whether a visually perceived object belongs to various categories across multiple levels of abstraction, such as animate, mammal, big, and so on. Earlier theories proposed that basic-level categories (e.g., “dog” and “cat”) are classifiable because exemplars share visual and other features (Rosch et al. 1976), whereas more abstract categories (i.e., superordinate categories; e.g., animate) tend to lack defining features (Murphy and Medin 1985; Rosch et al. 1976) and rely more on semantic knowledge and experience (Murphy 2004; Ralph et al. 2017). However, the rise of deep neural network models (DNNs) has transformed our view of how apparently complex object-classification tasks can be solved by sophisticated visual processing alone. Since then, many studies have compared internal object representations between DNNs and biological brains (Khaligh-Razavi and Kriegeskorte 2014). In particular, research in macaque monkeys has demonstrated strong similarities between ventral-pathway neural activity and DNNs (Yamins et al. 2014), reinforcing the notion that visual learning and processing underlying core object recognition are shared across humans and animals (Kar and DiCarlo 2024).

However, surprisingly little is known about whether, and to what extent, monkeys can classify objects into a variety of high-level categories as humans do, because most previous behavioral studies focused on basic object-recognition performance, such as similarity judgment (Koopman et al. 2017; Rajalingham et al. 2015). For example, Rajalingham et al. (2015) systematically compared human and monkey behavior using a match-to-sample task with 3D object stimuli. With this design, they were able to measure monkeys’ confusion rates across numerous object identities, revealing a striking similarity to human performance. But this approach does not readily scale to higher-level categorization, as the task design cannot easily prompt animals to attend to more abstract features; indeed, the observed confusion patterns did not appear to clearly reflect superordinate structures such as animacy (Rajalingham et al. 2015; but see Cherian et al. 2025). In contrast, earlier psychological studies in humans have demonstrated profound influences of language and conceptual knowledge on object categorization (Murphy 2004). Thus, critical gaps between humans and monkeys may emerge when they are required to report higher-level categories. This question is timely, as growing efforts are being made to develop DNNs whose representations of abstract object concepts align with those of humans (Du et al. 2025; Muttenthaler et al. 2025; Vong et al. 2024).

Prior to the proliferation of DNNs, multiple studies actually challenged monkeys with various image classification tasks, but they had their own limitations. For example, a series of studies by Fabre-Thorpe et al. (Delorme et al. 2000; Fabre-Thorpe 2003; Fabre-Thorpe et al. 1998; Fize et al. 2011) demonstrated that monkeys can report the presence or absence of an animal in rapidly presented natural photographs using a Go/No-Go paradigm. Other studies have shown that monkeys can classify trees vs. non-trees (Vogels 1999), humans vs. non-humans (Schrier and Brady 1987), monkeys vs. non-monkeys (Yoshikubo 1985), food vs. non-food (Fabre-Thorpe et al. 1998; Santos et al. 2001), cats vs. dogs (Eldridge et al. 2018; Freedman et al. 2001; Minamimoto et al. 2010), and cars vs. trucks (Minamimoto et al. 2010). But, each study typically tested only one or, at most, two classification rules, mostly using basic-level categories, and to our knowledge, there has been no systematic quantification of which classification rules monkeys can learn and which they

71 cannot. Moreover, it was obviously impossible for earlier studies to compare monkey behavior  
72 with DNNs. Therefore, a systematic characterization of monkey behavior is still much needed;  
73 yet classification tasks are challenging, as monkeys must be trained on each rule individually, and  
74 performing a test battery becomes impractical if training for each task requires extended durations.

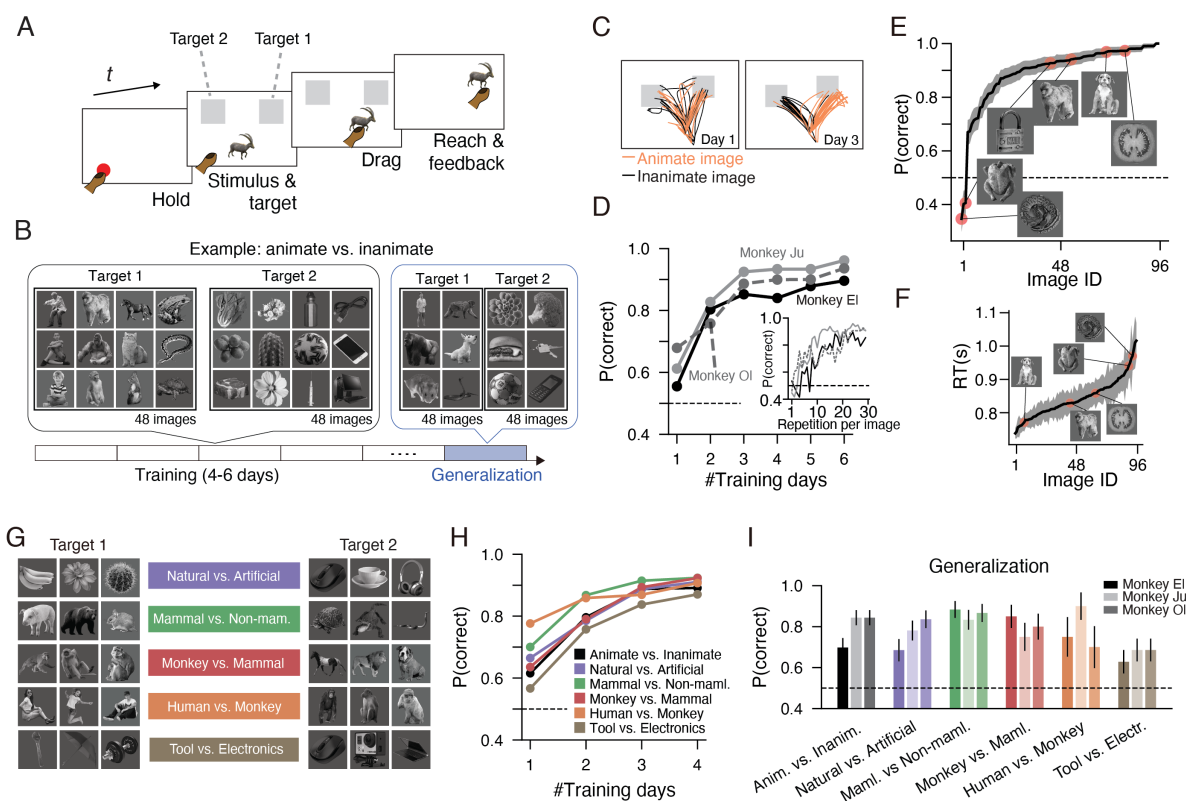
75 Here, we overcome these challenges by devising a novel binary decision-making paradigm  
76 that can train monkeys quickly, allowing us to test a series of image classification challenges with  
77 three monkeys (more than 10 tasks, 26 including subtasks, ~315K trials total; see Movies 1-3  
78 for examples). The monkeys learned many classifications, such as animate vs. inanimate, natural  
79 vs. artificial, mammal vs. non-mammal, big vs. small, and their performance was correlated  
80 with that of humans. At the same time, they failed to learn some more abstract rules that humans  
81 could readily acquire. Overall, their behavior was better explained by purely visual DNNs without  
82 language input, whereas human performance was better explained by language-guided DNNs.  
83 While our paradigm does not answer whether monkeys internally represent these concepts (e.g.,  
84 “mammal”), their performance provides an important benchmark for the extent to which biological  
85 neural networks can succeed in image-classification challenges without access to language and  
86 human-specific knowledge.

## 87 **Results**

### 88 **Monkeys quickly learn many visual classification tasks defined at various lev-** 89 **els of abstraction**

90 Our task is a version of binary decision-making paradigms, in which monkeys categorize an object  
91 image into one of two options according to a hidden rule (Fig. 1A). The experiments were con-  
92 ducted using a touchscreen system installed on the monkey’s home cage (Supplementary Fig. 1A).  
93 During each trial, after the monkey touched a red dot, an object image appeared together with two  
94 gray target boxes. The monkey had to touch the object image and drag it to one of the two target  
95 boxes. The image position moved along with the monkey’s finger until it reached a box. Upon  
96 reaching, the two boxes disappeared, and the same object image was revealed under the correct  
97 box, together with a juice reward for the correct choice. We reasoned that direct interaction with  
98 an object image encouraged monkeys to rely on the object’s properties to make a behavioral choice.  
99 Furthermore, this dragging motion required monkeys to slow down their choice reports ( $> 0.5$  sec  
100 to reach a target), potentially suppressing random impulsive choices and allowing them to make  
101 more deliberate decisions (Kaufman et al. 2015). The categorization rule remained the same within  
102 a session, and the same rule continued across sessions until the monkey’s performance plateaued.  
103 The monkeys had to infer this rule based on the feedback they received on each trial.

104 This task design enabled rapid, high-throughput training of monkeys. We first trained this task  
105 structure using amorphous shape stimuli for a month (Supplementary Fig. 1B, C) and then started  
106 the first object classification experiment: animate versus inanimate classification (Fig. 1B-F). The  
107 training images consisted of a variety of gray-scale images without background (Fig. 1B), but all



**Figure 1: Monkeys rapidly learn various abstract object concept classification tasks.** (A) Object drag task. In each trial, after the monkey held a fixation dot, an object image and two gray target boxes appeared. The monkey had to touch and drag the image to the correct box to receive a reward. Two boxes were associated with different object concepts. (B) We used grayscale photographs of objects with no background in this first set of experiments. We used ~100 training images to train monkeys and measured their generalization performance with test images. (C) Example reach trajectories showed that monkeys tended to choose targets randomly on Day 1, but by Day 3, they chose correct targets in many trials. (D) All monkeys' correct rates rapidly improved during the first three days of the animate vs. inanimate task, reaching 85-90% and then plateaued (chance level: 50%). Before the performance plateaued, monkeys saw each training image 10-20 times (inset). (E) After learning of the animate vs. inanimate task (post-Day 3), monkeys made more than 90% correct choices for many images, but performed poorly on certain images. The plot is the average of the three monkeys, with images ordered by correct rates. The shading indicates S.E.M. across trials. (F) Reaction times (RTs), measured as the time from stimulus onset to reaching a target, also varied across images. The plot is the average of the three monkeys, with images ordered by RTs. (G) We tested a variety of classification rules in succession. Each task used 60-96 grayscale training images. (H) Monkeys learned all tasks with similar learning rates. The plot showed the average of the three monkeys. (I) After training on each task, the monkeys could generalize the rule to new stimuli they saw for the first time. The error bars represent S.E.M. across images. All images were used either under the CC0 license or the Pixabay content license.

108 three tested monkeys achieved high classification performance within a few days of training (Fig.  
109 1C-D; averaging 760 trials per day). They started to show above-chance performance from Day  
110 1 ( $p < 0.011$  for all monkeys, binomial test), and their correct rate monotonically improved over  
111 the days, plateauing at approximately Days 4-6 with a performance of  $\sim 90\%$  (Fig. 1D). The same  
112 training images ( $n = 96$ ) were repeatedly presented during these sessions, but the monkeys started  
113 to show above-chance performance only after a few repetitions per image (Fig. 1D inset;  $p < 0.05$   
114 at 3, 4, and 9 repetitions for monkeys Ol, Ju, and El, respectively, binomial test).

115 Although their overall performance did not reach a 100% correct rate in six days, we found that  
116 they achieved a nearly 100 % correct rate for some images and performed poorly for some other  
117 images (Fig. 1E). Images with a nearly 100% correct rate were often prototypical examples of  
118 animate and inanimate categories such as monkeys, humans, and key locks, while images with poor  
119 correct rates tended to be atypical examples, such as a snake in a coil or a roasted chicken (defined  
120 as inanimate food). The patterns of correct rates were also largely consistent across monkeys  
121 (average of three comparisons:  $R = 0.49$ ,  $p < 2.1 \times 10^{-4}$  for all comparisons). The reaction  
122 times (RTs) also tended to be longer for these stimuli (Fig. 1F), and the drag trajectories for  
123 these stimuli occasionally showed patterns that resembled those of the animals were changing  
124 their minds (Supplementary Fig. 2; Kaufman et al. 2015).

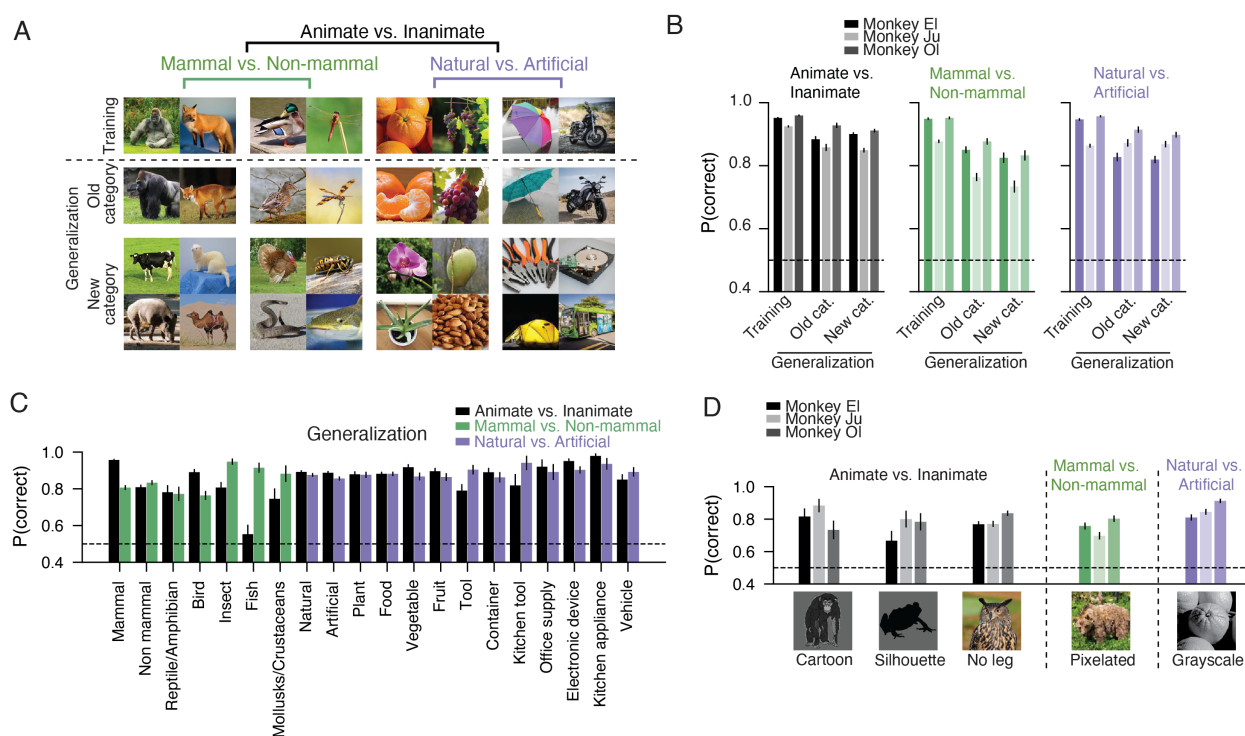
125 We then confirmed that the monkeys could generalize the rules that they learned to new images  
126 seen for the first time (Fig. 1I). After the learning sessions, we introduced a new set of images  
127 ( $n = 96$ ) that largely overlapped basic-level categories with the learned images but had different  
128 appearances (Fig. 1B). The monkeys' correct rates for these new images were high from the first  
129 exposure ( $p < 6.6 \times 10^{-5}$ , binomial test).

130 Encouraged by these results, we tested a variety of classification rules based on both superordi-  
131 nate and basic-level object concepts (Fig. 1G). These tasks were natural (plants, food) vs. artificial  
132 (man-made) objects, mammal vs. non-mammal animals (reptiles/amphibians), monkeys vs. mam-  
133 mals, humans vs. monkeys, and tools vs. electronics. The monkeys also had no difficulty learning  
134 these tasks in a few days (Fig. 1H;  $p < 10^{-10}$  for all monkeys in all tasks on Day 4, binomial test)  
135 and generalizing the rule to new images that they had never seen before (Fig. 1I;  $p < 6.1 \times 10^{-6}$   
136 in all tasks aggregated across monkeys, binomial test). When we switched to the next rule, we did  
137 not give the monkeys any explicit cues, but they learned to adopt the new rule quickly.

## 138 **Monkey behavior depends on combinations of visual/categorical features**

139 How did monkeys perform these categorizations? There are a couple of trivial hypotheses that have  
140 to be excluded. For example, monkeys might have simply formed associations between individual  
141 images and choices, and their choices for new stimuli in the generalization set (Fig. 1I) might be  
142 based on similar stimuli they had learned before. Alternatively, monkeys might have just learned  
143 a specific feature useful for categorization, such as the presence of faces or legs in the animate vs.  
144 inanimate task. Below, we show several control experiments in our test battery that excluded these  
145 possibilities.

146 One test involved a larger number of generalization images whose object categories were the  
 147 same or different from those shown during training. We used images from the THINGS database  
 148 (Hebart et al. 2019; Stoinski et al. 2024) and selected a subset of the concrete object concepts  
 149 in the database (e.g., dog, spider, lemon, and table) to create three tasks (Fig. 2A): animate vs.  
 150 inanimate (373 object categories, 3,683 images in total), natural vs. artificial objects (230 object  
 151 categories, 1,915 images in total), and mammal vs non-mammal animals (143 object categories,  
 152 1,668 images in total). For each task, we selected 100 images as a training set and then introduced  
 153 the remaining images to test generalization. Importantly, the generalization images, which were  
 154 shown only once, could be from object categories that were shown during training (“old”; Fig.  
 155 2A, second row) or were never shown (“new”; Fig. 2A, bottom two rows). This enabled us to  
 156 test whether the generalization performance depended on the similarity to the trained images (i.e.,  
 157 whether the “old” categories showed greater performance).



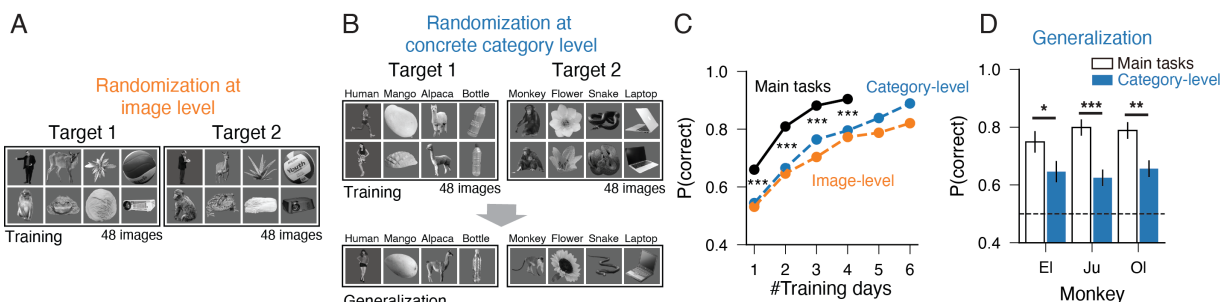
**Figure 2: Behavior cannot be explained by exemplar-based strategies or some accidental features.** (A) We assessed monkey performance in three tasks using large-scale image sets from the THINGS database (Hebart et al. 2019). We used 100 images to train monkeys (top row) and then presented 1,500-3,600 generalization images, containing “old” object categories present in training images (second row) and “new” categories that were never shown (bottom two rows). Each task used different training images. (B) All monkeys performed well on generalization images (75-93%) for both the “old” and “new” object categories. (C) Performance was generally high across a range of object categories. (D) Monkeys also performed well on control images that lacked certain visual features, such as cartoons without naturalistic textures, silhouettes without internal structures such as faces, and grayscale images without color. (B-D) Error bars indicate S.E.M. across object categories. All images were used either under CC0 license or Pixabay content license.

158 The monkeys performed these tasks without difficulty, despite experiencing much greater vari-  
159 ability in image appearance than that in the previous image sets (Fig. 2B; Movies 1-3;  $p < 10^{-10}$   
160 for all monkeys in the three tasks, binomial test), and importantly, their generalization perfor-  
161 mances to “new” object categories were largely comparable to that of the novel images of the “old”  
162 object categories they learned during training (Fig. 2B;  $U = 2,360 - 14,902$ ,  $p = 0.086 - 0.98$   
163 across monkeys and tasks, two-tailed Mann-Whitney  $U$ -test,  $BF_{01} > 3.7$  except for the mammal  
164 vs non-mammal in monkey OI whose  $BF_{01}$  was 1.9). Their performances were also largely similar  
165 across a range of object categories in the datasets (Fig. 2C), except for “fish” images during the  
166 animate vs. inanimate task. Thus, it is unlikely that the monkeys generalized the rules by judging  
167 the similarity to specific exemplars they learned during the training.

168 To test if the monkeys relied on a particular feature in the images (e.g., face, leg-like shape,  
169 certain colors, textures), we also presented control images that lacked these features (Fig. 2D). In  
170 the animate vs. inanimate task, we presented cartoon images that lacked naturalistic textures, and  
171 silhouette images that retained only the outlines of objects without internal structures such as faces.  
172 In the THINGS image set, there were close-up images of animals without legs. The monkeys  
173 performed well on these images (Fig. 2D; 73-88 % correct,  $p < 6.7 \times 10^{-3}$  for all monkeys,  
174 binomial test). In contrast, they could not generalize the rule to “textform” images (Kramer et al.  
175 2023; Long et al. 2018), which retained only mid-level image features (Supplementary Fig. 3).  
176 In the natural vs. artificial task, color could have been strong cues (natural objects tend to be  
177 greenish), but we found that the monkeys performed well on grayscale images (81-91 % correct,  
178  $p < 10^{-10}$  for all monkeys, binomial test). Finally, fur texture can be diagnostic for the mammal  
179 vs. non-mammal task, but the monkeys performed the task well even after distorting texture (see  
180 Methods; 70-80 % correct,  $p < 10^{-10}$  for all monkeys, binomial test). Thus, they likely did not  
181 rely on a single specific feature to solve the tasks.

182 As yet another test of whether monkeys learned specific stimuli instead of learning a common  
183 rule, we devised two versions of control tasks in which we associated random object images with  
184 targets (Fig. 3A-B). In one version, all new images ( $n = 96$ ) were fully randomly assigned to  
185 two categories, thus monkeys had to remember stimulus-target associations individually. To our  
186 surprise, the monkeys still gradually learned these many associations over a few days, but their  
187 learning speed was much slower than the main tasks (Fig. 3C). In the second version, we randomly  
188 assigned object categories to two targets, but different images within the same concrete category  
189 (e.g., two images of bottle) were always associated with the same target. Nonetheless, the learning  
190 speed was again slower than the main tasks (Fig. 3C), and importantly, the monkeys performed  
191 poorly in generalizing these category associations to new images (Fig. 3D). Thus, their behavior  
192 cannot be purely based on the similarity to individual images they have learned before to categorize  
193 new images.

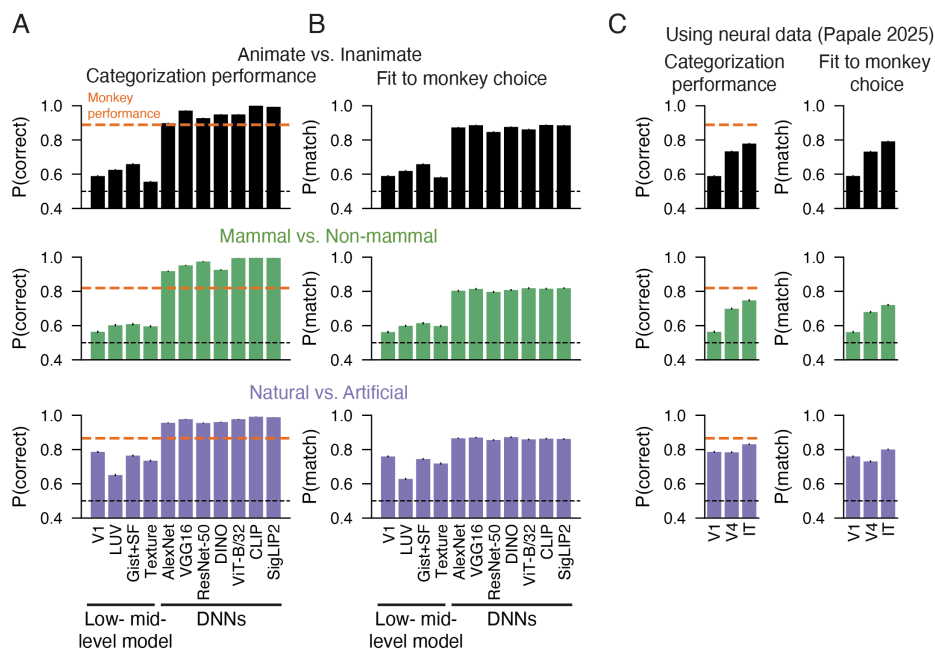
194 We then performed a comparison with monkey choices and outputs of a range of vision mod-  
195 els from low-level image features to state-of-the-art DNNs (Fig. 4). Low- and mid-level image  
196 models included V1-level filter outputs (Pinto et al. 2008), luminance- and color-based statistics,  
197 spatial-frequency summary statistics, and texture statistics (Portilla and Simoncelli 2000) derived  
198 from correlation structures of V1-level outputs. DNNs included AlexNet (Krizhevsky et al. 2012),  
199 VGG16 (Simonyan and Zisserman 2014), ResNet-50 (He et al. 2015), and Vision Transformer



**Figure 3: Absence of a common conceptual rule results in poorer performance.** (A-B) We assessed how well monkeys could learn stimulus-response associations in the absence of shared conceptual rules in two control tasks. In one task (A), two targets were associated with randomly selected images from our grayscale cropped image set (48 images per target). Thus, the monkeys had to remember the associations for individual images. In the other task (B), two targets were randomly associated with concrete object concepts (e.g., apple, horse, bin; referred to here as the “concrete category”; 16 concrete categories for each target). During the generalization test, different images of the same concrete categories were shown. (C) Although the monkeys could gradually learn both control tasks, performance was consistently lower than on the main concept classification tasks (the average across the six tasks shown in Fig. 1;  $p < 10^{-10}$  for all days, aggregated across monkeys, binomial test). (D) Their generalization performance for the concrete category randomized task (B) was also lower than that of the six main tasks ( $p < 0.016$  for all three monkeys, binomial test). Error bars are S.E.M. across images. Asterisks indicate statistical significance: \* $p < 0.05$ , \*\* $p < 0.01$ , \*\*\* $p < 0.001$ .

200 (ViT-B/32) (Dosovitskiy et al. 2020), which were pretrained on ImageNet. We also tested an un-  
 201 supervised vision model, DINO (Caron et al. 2021), and two recent large-scale models (CLIP  
 202 (Radford et al. 2021) and SigLIP2 (Tschannen et al. 2025)) that were trained to match visual and  
 203 text inputs, thus providing additional language supervision. For all the models, we first reduced  
 204 their output dimensions using principal component analysis (PCA), constructed a linear classifier  
 205 for either the true categories or the monkey’s choices for the training set, and then tested their  
 206 performance on the generalization set (see Methods for details). The results did not depend on  
 207 specific details of the dimension-reduction and regression methods (Supplementary Fig. 4).

208 We found that DNNs consistently outperformed low-level models in both task performance  
 209 and fitting monkey behavior. The classification performance of the low- and mid-level models was  
 210 typically  $\sim 60\%$ , and they were all significantly lower than the monkey performance ( $p < 10^{-10}$   
 211 for all comparisons between models and average monkey performance, binomial test, Bonferroni-  
 212 corrected for multiple comparisons). Correspondingly, even when these models were fit directly  
 213 to monkey behavior, the rates of choice agreement between the models and the monkeys were low  
 214 ( $\sim 60\%$  in many cases). By contrast, all DNNs performed our tasks with high accuracy ( $>90\%$ )  
 215 and fit the monkey behavior better than the low- and mid-level visual models ( $p < 10^{-10}$  for  
 216 all comparisons between models; binomial test, Bonferroni-corrected for multiple comparisons).  
 217 In addition, we leveraged previous large-scale neural recordings from visual cortices using the  
 218 same THINGS image set (Papale et al. 2025) and fitted neural responses to monkey behavior. As  
 219 expected, fitting performances were higher for later stages of the ventral visual pathway (Fig. 4C).



**Figure 4: Monkey performance is better explained by DNNs than by low- or mid-level vision models.** (A) Category classification accuracy of various vision models, including deep neural networks (DNNs). We constructed a linear classifier using the output of each model based on the training images and plotted their performance on the generalization images we used for the monkeys. The orange dashed lines represent the average classification performance of the three monkeys. (B) Accuracy of predicting monkey choices. A linear classifier was trained to classify monkey choices instead.  $P(\text{match})$  indicates the probability that the monkey and model choices matched. DNNs outperformed lower-level vision models. (C) Accuracy of category and monkey-choice classification, based on neural responses in ventral pathway areas, was higher at later stages along the pathway. We employed neural responses collected by Papale et al. (2025), who used the same THINGS database.

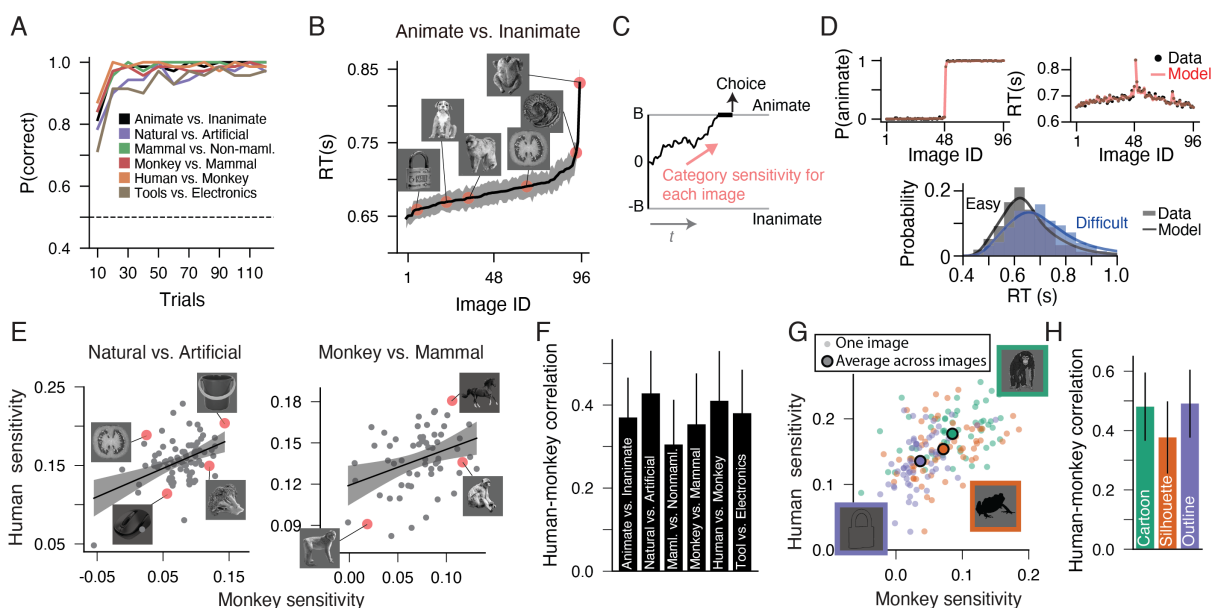
220 Altogether, these results support that the monkeys used high-level visual features, compara-  
 221 ble to the later stages of DNNs or the ventral visual pathway, rather than using trivial, low-level  
 222 features in the image sets.

## 223 Monkey categorization was also moderately correlated with human behavior

224 We next sought to compare monkey and human behavior in the same tasks and settings shown in  
 225 Fig. 1 to determine whether they share similar categorization patterns. When human participants  
 226 performed the same tasks, they understood the rules within tens of trials and thereafter showed  
 227 near-perfect performance across all tasks (Fig. 5A), making it impractical to compare patterns of  
 228 correct rates. However, they still showed different reaction times (RTs) across images (Fig. 5B,  
 229  $F(95, 570) = 2.88, p < 10^{-10}$  for animate vs inanimate,  $F > 2.14, p < 4 \times 10^{-6}$  for the other  
 230 tasks, repeated-measures one-way ANOVA). Thus, we expected that difficult stimuli would cause

231 longer RTs in humans and lower choice accuracy (and longer RTs) in monkeys (Fig. 1E-F).

232 We therefore constructed a metric of stimulus difficulty that embraces both choice accuracy and  
 233 RTs. We used the drift-diffusion model (Fig. 5C), which has been widely used to explain choices  
 234 and RTs in various experimental settings (Carlson et al. 2014; Heidari-Gorji et al. 2021; Luo et al.  
 235 2025; Mack and Palmeri 2011). The model assumes a one-dimensional stochastic internal state  
 236 that meanders between two decision bounds corresponding to two choice options until it reaches  
 237 one of them (Fig. 5C). An easy stimulus causes a large drift (higher sensitivity) toward the correct  
 238 bound, leading to higher accuracy and shorter RTs. The slope of the drift is termed “category  
 239 sensitivity” here, whose absolute value reflects how easily an image can be categorized correctly.  
 240 Note that we did not assume that an evidence accumulation process underlies our task (Stine et al.  
 241 2020). Rather, we used the model as a means of deriving a single metric of stimulus difficulty (the



**Figure 5: Human behavioral responses were correlated with monkey performance.** (A) Humans quickly learned the tasks in Fig. 1 within 20-30 trials. The plots are the average of seven participants per task. (B) Different average reaction times (RTs) across images, averaged across participants. (C) We fitted the drift-diffusion model (DDM) to the choices and RTs. The model makes a decision based on a stochastic state biased by the category sensitivity of each image. (D) DDM accurately fits choices, RTs, and RT distributions. The plots are from the animate vs. inanimate task. Image IDs were sorted by category sensitivity. RT distributions were generated using the trials aggregated across the top (easy) and bottom (difficult) 10 images. See Supplementary Fig. 6 for all tasks. (E) Example comparison of category sensitivity between humans and monkeys. Each dot represents an image, whose absolute category sensitivity was averaged across participants. The line is a linear regression, and shading indicates the standard error of the prediction line. (F) Positive correlations between humans and monkeys. Error bars indicate S.E.M. across images. (G) The category sensitivity for control images —cartoon, silhouette, and outline— was in the same order between monkeys and humans. (H) Positive correlations were also evident within each control image set.

242 absolute value of category sensitivity) for each image from the choices and RTs. Critically, the  
243 model fits both human and monkey data well (Fig. 5D, Supplementary Fig. 6;  $R = 0.99 - 1.0$   
244 for choices and  $R = 0.81 - 0.99$  for RTs across monkeys, humans, and tasks), demonstrating its  
245 utility.

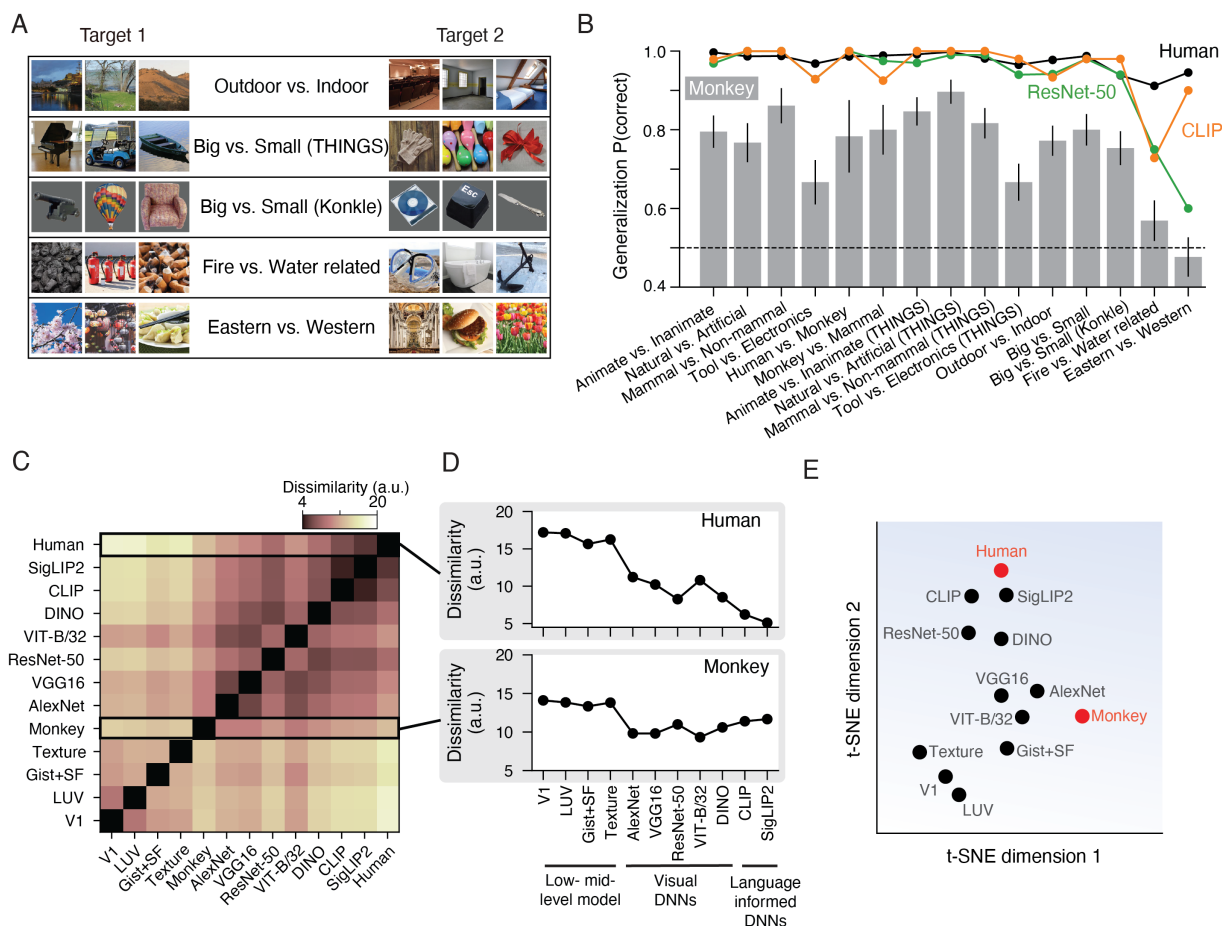
246 For all tasks, the absolute value of the category sensitivity of humans and monkeys across  
247 images had moderate positive correlations (Fig. 5E, F). Humans obviously had higher category  
248 sensitivity for all images, as expected from their near-perfect performance. Nonetheless, there  
249 was a tendency for images with higher category sensitivity for humans to also be higher for mon-  
250 keys (Fig. 5E). The correlations of sensitivities between humans and monkeys were significantly  
251 positive ( $R = 0.30 - 0.43$ ;  $p = 7.4 \times 10^{-5} - 6.0 \times 10^{-3}$ ) for all tasks (Fig. 5F).

252 We also checked the two species' responses to pictures outside the range of natural images  
253 and again found the similarity between them (Fig. 5G, H). In the animate vs. inanimate task,  
254 we showed cartoon, silhouette, and outline images (Fig. 5G; part of results shown in Fig. 2D).  
255 The monkeys showed higher performance for cartoon and silhouette images than for outline im-  
256 ages, which was consistent with slower human RTs for outline images. We also found significant  
257 positive correlations of category sensitivities between humans and monkeys within each control  
258 set (Fig. 5H). Note that DNNs are also known to have poorer classification performance for out-  
259 line images than for silhouette images, even though they must retain the same shape information  
260 (Supplementary Fig. 5B; Baker et al. 2018).

## 261 **Triangular comparison across monkeys, humans, and artificial networks**

262 Thus far, monkeys have generalized all categorization rules we have tested, but are there any rules  
263 they cannot generalize? To explore this, we created five additional categorization tasks to challenge  
264 them further (Fig. 6A). As often investigated in the vision science literature, we first tested outdoor  
265 vs. indoor scenes (Zhou et al. 2017) and big vs. small object categorization (Konkle and Oliva  
266 2012; Long et al. 2016). Big vs. small objects were determined on the basis of physical object size  
267 (objects bigger than a chair were classified as “big”; Konkle and Oliva 2012) and had no correlation  
268 with image size. We had two versions of big vs. small, one using the set from Konkle et al. (2012)  
269 and one using images from the THINGS database. We then created two more tasks that would  
270 rely heavily on human knowledge: “fire-related” vs. “water-related” objects, such as lighter, fire  
271 extinguishers, faucets, and bathtubs, etc., and “Western culture” vs. “Eastern culture” objects,  
272 such as mooncakes, Chinese knots, crowns, and Easter eggs, etc. For each task, we first trained  
273 the monkeys on training images for 4-6 days and then assessed their generalization performance  
274 on test images as in Fig. 1B.

275 The monkeys successfully generalized rules in some tasks and failed in others, and visual  
276 DNNs exhibited similar patterns (Fig. 6B). The monkeys successfully generalized the outdoor vs.  
277 indoor rule to new stimuli ( $p < 4.3 \times 10^{-7}$  for all monkeys, binomial test) and also performed  
278 well on the two versions of the big vs. small task ( $p < 2.4 \times 10^{-6}$  for all monkeys in both tasks,  
279 binomial test), as did visual DNNs such as ResNet-50. In contrast, the monkey performance was



**Figure 6: Triangular comparison across monkeys, humans, and artificial networks.** (A) To determine the limit of monkeys' ability to classify concepts, we tested more concept tasks. As with the other tasks in Fig. 1, we trained the three monkeys with 90-120 images and assessed their generalization performance with 90-120 new images. (B) The monkeys ( $n = 3$ ) performed well on tasks such as big vs. small objects and indoor vs. outdoor scenes, but they failed on some abstract tasks, such as fire- vs. water-related objects and Western vs. Eastern objects. Visual deep neural networks (DNNs) such as ResNet-50 also showed poorer performance on these tasks, whereas humans ( $n = 9$ ) performed almost perfectly across all tasks. The monkey performance for the three THINGS tasks in Fig. 2 was replotted using only the first 100 test stimuli to facilitate comparison with the other tasks. See Supplementary Fig. 5A for the performance of other DNNs. Error bars indicate S.E.M. across images. (C) Finally, we concatenated the generalization performances of images across all tasks to generate the matrix of dissimilarity in behavioral performance across humans, monkeys, and models. (D) A closer look at the dissimilarity revealed that the monkey behavior was most similar to visual DNNs without language input, whereas the human behavior was most similar to language-informed DNNs. (E) t-SNE visualization of the dissimilarity matrix showed a spectrum from low-level visual models to language-informed DNNs, with the monkey behavior placed in the middle.

280 substantially lower for the fire- vs. water-related task (46-64% correct). And their performances  
 281 were at chance level for the Western vs. Eastern object task ( $p = 0.38 - 0.86$ ,  $BF_{01} = 4.9 - 7.4$ )

282 across monkeys, binomial test). Notably, visual DNNs such as ResNet-50 trained on the ImageNet  
283 also displayed poorer performance on the fire vs. water task (53-75 % correct across models) and  
284 on the Western vs. Eastern task (46-60 % correct across models; see Supplementary Fig. 5A). The  
285 performance of a language-informed model (e.g., CLIP) was better and above chance in both tasks  
286 ( $p < 6.9 \times 10^{-6}$ , binomial test). The human participants performed these tasks almost perfectly,  
287 suggesting their superiority over monkeys and DNNs.

288 Finally, we combined these generalization performance metrics across all tasks to visualize  
289 the overall similarity of behavioral responses among monkeys, humans, and vision models. For  
290 monkeys and humans, the average correct rate was calculated for each generalization image across  
291 subjects and then concatenated across tasks. For the models, we used the likelihood of correct  
292 classification for each image based on a logistic classifier trained on the training images. A matrix  
293 of dissimilarities among humans, monkeys, and models based on these values (Fig. 6C) revealed  
294 that monkey choices were closer to DNNs trained without language input (e.g., AlexNet; Fig. 6D  
295 bottom), whereas human choices were closer to the language-informed DNNs (CLIP and SigLIP2;  
296 Fig. 6D top). Consistent with this, t-SNE visualization revealed a spectrum from low-level visual  
297 models (e.g., the V1 model) to language-informed DNNs, with the monkey data located in the  
298 middle, proximal to visual DNNs (Fig. 6E).

## 299 Discussion

300 Macaque monkeys serve as a primary model for investigating the human visual system, allowing  
301 characterization of neural mechanisms at the cellular level and providing insight into the evolution  
302 of object vision (Orban et al. 2004; Yamins and DiCarlo 2016), but it is poorly investigated whether  
303 monkeys can classify objects according to various levels of categorization rules. We developed a  
304 task design that enabled us to test monkeys using a battery of image-classification challenges to  
305 quantify their capacities and limitations. Monkeys generalized a surprising variety of categoriza-  
306 tion rules to novel images, and their behavior could not be explained by learning individual objects  
307 (Figs. 2, 3) or by classification based on low-level image features (Figs. 2D, 4). Instead, mon-  
308 key behavior was best fit by visual DNN outputs (Fig. 4) and was also partially correlated with  
309 human behavioral responses (Fig. 5). At the same time, monkeys failed in some abstract rules  
310 (Fig. 6B). Overall, whereas human behavior was proximal to language-informed DNNs, monkey  
311 behavior was closer to DNNs without language input (Fig. 6C-E). The present study represents an  
312 unprecedentedly systematic characterization of animals' object-classification performance.

313 The widespread use of macaque monkeys in ventral pathway research was guided by the  
314 premise that they recognize objects just as humans do, but there is not even consensus on how  
315 to measure their object-recognition behavior. One approach is to use tasks that require monkeys  
316 to report the similarity of object images, such as match-to-sample (Rajalingham et al. 2015, 2018)  
317 and odd-one-out (Cherian et al. 2025; Hebart et al. 2020). Since these tasks do not require training  
318 of individual categories, the similarity can be measured across many combinations of object cat-  
319 egories to generate confusion matrices. The downside, however, is that the methods lack explicit  
320 control over which properties of object images the animals rely on to report similarity. Indeed,

321 confusion matrices measured in match-to-sample (Rajalingham et al. 2015) do not appear to be  
322 strongly influenced by higher-level categories such as animacy (but see Cherian et al. 2025 for the  
323 odd-one-out case). The second approach is to explicitly train animals to classify two categories  
324 (e.g., dog vs. cat) and measure their generalization performance, but in this case, each task must  
325 be trained individually. Consequently, previous studies performed only one or two tasks chosen  
326 according to their experimental purposes (Eldridge et al. 2018; Fabre-Thorpe et al. 1998; Minami-  
327 moto et al. 2010; Vogels 1999). Our innovation is that by leveraging a rapid training paradigm, we  
328 scaled the second approach into a large test battery. While we focused on various natural-object  
329 concepts used in vision research, the rules do not have to be confined to verbally describable cate-  
330 gories or natural-object images.

331 It was a surprise to us that the monkeys learned each task rule quickly, although the tasks  
332 required them to discover a common classification rule from many object images with different  
333 appearances, which contrasts with conventional visuomotor association tasks that use only a few  
334 images (Cromer et al. 2011). We also observed rather surprising consistency in the monkeys' learn-  
335 ing curves across a variety of classification tasks (Fig. 1H). Many studies posited that basic-level  
336 categories are most accessible due to the family resemblance of image features within a category  
337 (Grill-Spector and Kanwisher 2005; Rosch et al. 1976), whereas other studies have reported that  
338 superordinate categories are processed faster in the human brain (Clarke and Tyler 2015; Macé et  
339 al. 2009). But these differences in speed might simply depend on the amount of evidence needed  
340 to make a decision (Mack and Palmeri 2011), and the underlying decision-making mechanisms  
341 might not qualitatively differ across levels of object classification.

342 Our paradigms do not answer whether the monkeys were truly aware of object concepts such  
343 as “animacy” and “natural”, but the results still provide insights into the mechanisms that underlie  
344 the formation and representation of such concepts in the brain. Recent works have proposed the  
345 similarity of conceptual representations between humans and DNNs by showing alignment in their  
346 internal or latent representations (Du et al. 2025; Muttenthaler et al. 2025; Vong et al. 2024). How-  
347 ever, to our knowledge, there is no established definition of what “object concepts” entails. One  
348 may consider that multimodal associations (Quiroga et al. 2005) or arbitrary sensory associations  
349 (Loggia et al. 2025) that cannot be explained by mere perceptual similarity as a key trait of object  
350 concepts. Earlier psychological studies also emphasize the importance of knowledge structure,  
351 such as recognizing the functions of objects to form concepts (Murphy 2004). However, many of  
352 the object concepts we used here were decodable from natural images through visual DNNs and  
353 could be classified by monkeys. These results rather support the view that many object concepts  
354 are already structurally embedded in retinal images to some extent (Vong et al. 2024) and that the  
355 primate brain can readily learn to extract them even without non-visual knowledge. While there  
356 should be additional contributions from knowledge and multimodal associations, it is conceivable  
357 that such embedding of concepts in images can substantially facilitate the formation of object con-  
358 cepts in the human brain.

## 359 **Acknowledgment**

360 We thank Takafumi Minamimoto and Cheng Xue for the comments on earlier versions of the  
361 manuscript. We thank Roozbeh Kiani for providing natural object images and also thank the  
362 Japan Monkey Centre and the Center for the Evolutionary Origins of Human Behavior, Kyoto  
363 University, for providing monkey images used in the experiments. This work was supported by  
364 Strategic Priority Research Program of the Chinese Academy of Sciences (XDB1010202), the  
365 National Science and Technology Innovation 2030 Major Program (Grant No. 2021ZD0203703),  
366 National Natural Science Foundation of China (Grant No. 32371077 and No. W2432019), and  
367 Shanghai Municipal Science and Technology Major Project (Grant No. 2019SHZDZX02).

## 368 **Competing interests**

369 The authors declare no competing interests.

## 370 **Data and code availability**

371 Data and code used in this study will be made publicly available upon official publication.

## 372 **Methods**

### 373 **Subjects and apparatus**

#### 374 **Monkey experiments**

375 We collected behavioral data from three adult rhesus monkeys (*Macaca mulatta*; monkey Ju, fe-  
376 male, age 13; monkey El, female, age 12; monkey Ol, male, age 11). These monkeys were raised  
377 in a typical laboratory environment and had no prior experience performing relevant behavioral  
378 tasks. The dataset collected in this study consisted of 26 tasks including subtasks, amounting to  
379 ~315K trials from the three monkeys. Experimental procedures conformed to the National In-  
380 stitutes of Health *Guide for the Care and Use of Laboratory Animals* and were approved by the  
381 Animal Care Committee at Center for Excellence in Brain Science and Intelligence Technology  
382 (Institute of Neuroscience), Chinese Academy of Sciences.

383 Experiments were conducted using a custom-built touchscreen device attached to each mon-  
384 key's cage (Supplementary Fig. 1A), similar to those developed by other laboratories (Butler and  
385 Kennerley 2019; Ramezanzpour et al. 2024). The device box consisted of a commercial touchscreen  
386 (AOSIMAN, ASM-125UCT, 12.5 inches; 60 Hz refresh rate; 3840 × 2160 pixel screen resolution),  
387 a reward juice delivery system with a TTL-triggered water pump, a monitoring camera, and a con-  
388 trol computer. The monkey reached the touchscreen through a hole in the front panel attached to  
389 the cage and viewed the screen through another hole, which opened at approximately the level of  
390 the monkey's head. Below the hole was a reward juice outlet in a "cone"-shaped socket of a size  
391 suitable to cover the monkey's mouth (Kawaguchi et al. 2019). The monkeys performed the tasks  
392 while their mouth was in this socket and thus had a stable head position. Under these conditions,  
393 the viewing distance to the monitor was approximately 15-18 cm. The task was controlled by a  
394 custom-made MATLAB (MathWorks, MA, USA) program with Psychtoolbox.

#### 395 **Human experiments**

396 We recruited in total 33 human participants (20-40 years old; 10 males and 23 females, students or  
397 employees at the Chinese Academy of Sciences). Our participant sampling did not consider sex,  
398 as it was considered unlikely to influence the results of this study. They were assigned to perform  
399 different tasks, as explained below. All participants had normal or corrected-to-normal vision and  
400 were naïve to the purpose of the experiments. Written informed consent was obtained from the  
401 participants. All experimental procedures were approved by the Institutional Review Board of the  
402 Center for Excellence in Brain Science and Intelligence Technology (Institute of Neuroscience),  
403 Chinese Academy of Sciences.

404 During the experiments, the participants were seated in a height-adjustable chair, with their  
405 chin and forehead supported by a chinrest. The chinrest ensured a stable viewing distance (37 cm)  
406 to a touchscreen monitor (LENOVO, T24t-20, 23.8 inches; 60 Hz refresh rate; 1920 x 1080 pixel

407 screen resolution). The view distances were different from those used in monkey experiments, but  
408 we maintained similar visual angles of the object images (see below). The other task parameters  
409 were also kept identical to those of the monkey experiments.

## 410 **Object drag task**

411 We developed a binary object categorization task (“object drag task”) for the touchscreen device  
412 (Fig. 1A). During the task, the participants viewed an object stimulus and moved it to one of two  
413 target boxes by swiping it with their fingers. The two target boxes were associated with distinct  
414 object concepts (e.g., animate vs. inanimate objects).

415 Each trial started when the participants touched a red fixation dot. Following this, a stimulus  
416 (10 – 11 deg) appeared together with two gray square target boxes. The positions of the stimulus  
417 and target boxes were adjusted for each participant so that they could move a stimulus smoothly  
418 to both targets with similar amounts of time. In the monkey experiments, the distance between  
419 the stimulus and target boxes was 18 – 24 deg, and the distance between the two target boxes  
420 was 24 – 37 deg. Following the stimulus presentation, the participants were required to touch the  
421 stimulus within 1 s. After the touch, the participants were required to drag the stimulus to one  
422 of the two target boxes within 5 s. In the case of timeout, the trial was aborted (1.3 % of trials  
423 for monkeys). When the image reached one of the two boxes, the image and the target boxes  
424 disappeared, and instead, the same image appeared at the correct target location, providing visual  
425 feedback to the participants. At the same time, an auditory tone indicated correct and incorrect  
426 responses. The monkeys received a juice reward for a correct response and a short timeout (1-2 s)  
427 for an incorrect response.

428 A rule for a given task (e.g., animate vs. inanimate) did not change within a daily experimental  
429 session for the monkeys (see below for rules). The same rule usually lasted for one to two weeks  
430 until the participants learned the task and completed a generalization test (Fig. 1B). At the start of  
431 a new rule, they were not informed of the rule change, but the new task involved different stimuli,  
432 which may have reduced the interference from the previous task and facilitated the learning of the  
433 new rule.

## 434 **Image sets**

### 435 **Background-cropped image sets**

436 In the first series of our experiments (Fig. 1), we used grayscale object images with the background  
437 removed (examples shown in Fig. 1B, G). Some of the images (5%) were obtained from the  
438 existing databases (Kiani et al. 2007), whereas the remaining images were obtained by manually  
439 removing background from natural photographs of objects obtained through search engines such  
440 as Bing, Pixabay, and Google Image. We also obtained photographs of monkeys from the Japan

441 Monkey Centre and the Center for the Evolutionary Origins of Human Behavior, Kyoto University.  
442 Images were converted to grayscale because the distribution of colors was highly skewed in some  
443 objects (e.g., monkey images are always brownish). Because some original images do not allow  
444 redistribution, we replaced them with similar images with the Creative Commons Zero (CC0)  
445 license or other royalty-free licenses when shown in the figures. The sources of the human images  
446 shown in Figs. 1, 3, and Supplementary Fig. 3 are PurePNG (<https://purepng.com/>; authors: LPuo,  
447 hadiaaltaf, Momo, FMPure).

448 When collecting images for these sets, we focused on eight object groups: humans, monkeys,  
449 non-primate mammals, non-mammal (reptiles/amphibians), food, plants, tools/everyday objects,  
450 and electronics. Although these eight groups span across different levels of category hierarchy  
451 (e.g., “monkey” can be considered a basic category, whereas “food” is a superordinate category),  
452 we considered that these groups covered a range of object categories ideal for behavioral experi-  
453 ments with monkeys.

454 Using these eight object groups, we created six tasks (Fig. 1B, G): animate vs. inanimate, nat-  
455 ural vs. artificial, non-primate mammal vs. non-mammal (reptile/amphibian), human vs. monkey,  
456 monkey vs. non-primate mammal, tool/everyday object vs electronics. In each task, 60-96 images  
457 were used to train the monkeys, and 20-96 images were used to test their performance on novel  
458 stimuli after learning (generalization test) (Fig. 1B). Training images could have a small number  
459 of overlaps across tasks, whereas the generalization images were never shown to the monkeys in  
460 any previous tasks.

## 461 **Large-scale THINGS image sets**

462 To complement the experiments using the controlled but small stimulus sets above, we created  
463 large-scale image sets based on the THINGS database (Hebart et al. 2019, 2020; Stoinski et al.  
464 2024). The THINGS database consists of 1,814 concrete object concepts (such as dog, spider,  
465 lemon, and table), each of which contains at least 12 natural photographs.

466 For the purpose of our experiments, we selected a subset of object categories from the THINGS  
467 database and created three tasks (Fig. 2A): animate vs. inanimate (373 object categories, 3,683  
468 images in total), natural vs. artificial objects (230 object categories, 1,915 images in total), and  
469 mammal vs non-mammal animals (143 object categories, 1,668 images in total). Among the im-  
470 ages, we selected 100 images (50 images per target; no overlap in object categories) for training  
471 the monkeys on each task. During the generalization test, we showed the remaining images only  
472 once each. Due to copyright issues, the images shown in the figures are replaced with similar  
473 images with the CC0 license from the THINGSplus database (Stoinski et al. 2024) or images with  
474 the Pixabay Content license from Pixabay.

475 To inspect the patterns of monkey behavior, we further classified these objects into high-level  
476 object categories such as “bird”, “plant”, and “fruit” (Fig. 2C). These high-level categories were  
477 mostly taken from the THINGS database (Hebart et al. 2019), which uses human curation to de-  
478 termine them. In addition, we included several high-level categories that we defined ourselves

479 (mammals, non-mammals, natural, artificial, reptiles/amphibians, fish, and mollusks/crustaceans)  
480 as they could be informative for assessing monkey performance.

### 481 **Modified control image sets**

482 After performing the animate vs. inanimate task using the grayscale images (Fig. 1), we tested  
483 whether the monkeys could classify images outside the range of natural photographs (Figs. 2D, 5G,  
484 H). We created three control image sets. First, an “outline” image set was created by extracting the  
485 outlines of the object images using *bwboundaries* function in MATLAB. The contour was colored  
486 black, while its interior was filled with the background color of the task screen (gray). Second, we  
487 created a “silhouette” image set, where the inside of a contour was also filled with black. Third,  
488 we created a “cartoon” image set by selecting cartoon images of objects from a publicly available  
489 database (the “Multipic” database; Duñabeitia et al. 2018). In all the sets, we selected source  
490 images that were not shown during the training to ensure that the monkeys could not solve the task  
491 by remembering specific object images. Each set contained 60 images.

492 We also presented two control image sets after performing the categorization of large-scale  
493 THINGS images. First, during the natural vs. artificial task, we showed THINGS images con-  
494 verted into grayscale to confirm that the monkeys did not simply rely on color. Second, during  
495 the mammal vs. non-mammal task, we presented THINGS images whose texture was distorted by  
496 using *crystallization* function (cell size, 10) in Photoshop 2024 (Adobe Inc., San Jose, CA, USA).  
497 We had 400 images for each set, selected from the generalization image set. These images were  
498 shown once during the generalization test, but we consider it is unlikely that the monkeys remem-  
499 bered specific images to solve these control tasks because they were shown only once and mixed  
500 with many other images.

### 501 **Image sets for more conceptual tasks**

502 To further challenge the monkeys’ ability to learn abstract rules, we created additional conceptual  
503 tasks (Fig. 6). Each task used 90-120 training images and 90-120 generalization images.

504 **Outdoor vs. indoor scenes** We selected scene images from Zhou et al. (2017). The indoor  
505 scenes included pictures of a variety of buildings, including small scenes such as bedrooms and  
506 kitchens and also spacious scenes such as stadiums and concert halls. The outdoor scenes included  
507 both natural scenes, such as forests, and urban scenes, such as buildings. Because the sky and  
508 natural objects (e.g., a tree) could be strong cues indicating outdoors, we preferentially selected  
509 images that did not contain them (e.g., the front face of a building). Nevertheless, outdoor scenes  
510 tend to have certain colors (greenish). Therefore, we created a grayscale version of the generaliza-  
511 tion images as a control and confirmed that the monkeys also performed well without color (73 %  
512 correct).

513 **Big vs. small objects** In this task, the monkeys had to classify objects based on their physical  
514 size. Objects that are typically bigger or smaller than the size of a chair were classified as big or  
515 small, respectively (Konkle and Oliva 2012; Long et al. 2016). We created two sets of stimuli. The  
516 first set used images from the THINGS database, which had a natural background. Small objects  
517 included many household items such as spoons and nails, and big objects included furniture, ve-  
518 hicles, and others. We ensured that pictures of most big objects were taken indoors so that this  
519 categorization could not be confused with outdoor vs. indoor. The second set used images from  
520 the big vs. small images created by Konkle et al. (2012), which had no background. A comparison  
521 of these two sets allowed us to confirm that the monkey performance did not depend on a specific  
522 choice of images.

523 **Fire- vs. water-related objects** To create abstract rules that are unlikely to be correlated with  
524 any visual features, we selected object images that are semantically associated with fire and water  
525 from the THINGS database. The fire-related category included objects such as stoves, cigarettes,  
526 coal, fire trucks, and radiators. The water-related category included objects such as hoses, bathtubs,  
527 paddles, and squirt guns. The images of fire or water itself were avoided as they have distinctive  
528 features.

529 **Eastern- vs. Western-culture objects** As another abstract categorization task, we created an  
530 image set of objects associated with Eastern and Western cultures. We collected the images mainly  
531 by using Internet search engines, because a sufficient number of images could not be found in the  
532 THINGS database. The Eastern category included objects such as Chinese lanterns, dumplings,  
533 old pavilions, bamboo, and chopsticks. The Western category included objects such as crowns,  
534 donuts, the Eiffel Tower, tulips, and violins.

### 535 **Texform images**

536 Texforms are synthetic images generated by matching the mid-level texture features of natural  
537 object images (Long et al. 2016, 2017, 2018). These features consist of the magnitudes and cor-  
538 relation structures of V1-like oriented bandpass filter outputs and some other image statistics at  
539 each local image region (Freeman et al. 2013; Portilla and Simoncelli 2000). Humans can reliably  
540 classify the animacy of original objects from texform images (Long et al. 2017); thus, we tested  
541 whether monkeys could generalize the animate vs. inanimate task to texform images (Supplemen-  
542 tary fig. 3). We used the original object and texform images from Long et al. (2018).

543 We first created a set of original object images (30 each for animate and inanimate images) to  
544 confirm that the monkeys could successfully classify these original images (Supplementary Fig.  
545 3A). We then created a texform set (30 each for animate and inanimate; Supplementary Fig. 3B).  
546 The texform images were chosen such that they were not made from the original images we have  
547 shown because the monkeys might have remembered the original images. After testing the gen-  
548 eralization performance of this texform set, we trained the monkeys on this texform classification

549 for 7-8 sessions (Supplementary Fig. 3C). Finally, we created another texform set (30 each for  
550 animate and inanimate; Supplementary Fig. 3D) and tested whether the monkeys could generalize  
551 the learned texform classification rule to new texform images. The set contained both big and small  
552 objects for animate and inanimate images, and we evenly distributed them for each set.

## 553 **Monkey experimental procedures**

### 554 **Object drag tasks using grayscale image sets**

555 We performed six tasks using grayscale image sets (Fig. 1B, G). The procedure for each task  
556 consisted of training and generalization test sessions. During the training sessions, an image from  
557 the training set was pseudo-randomly presented to the monkey in each trial. Training lasted for  
558 4-6 sessions (days) depending on the monkey's learning speed. In the early training phase of  
559 each task, we had a portion of trials in which we showed only the correct target box until the  
560 monkey performance reached above  $\sim 70\%$ . These trials were used to maintain the monkey's  
561 engagement and were excluded from the data analysis. Immediately after the training sessions,  
562 we had a generalization session, in which we started to show generalization images in 50% of the  
563 trials. We analyzed only the trials in which these generalization images were presented for the first  
564 time.

### 565 **Object drag tasks using large-scale THINGS image sets**

566 We performed three tasks using large-scale THINGS image sets. For each task, we first confirmed  
567 that the monkeys perform well on the same task with the grayscale image set. We then started to  
568 introduce 100 training images from the THINGS set (Fig. 2A). After we confirmed that the mon-  
569 key performance was stable for the 100 training images, we gradually introduced generalization  
570 images. Each generalization image was shown only once. We controlled the proportion of trials in  
571 which generalization images were presented, gradually increasing it up to 50%. In the other 50%  
572 of the trials, the training images were repeatedly shown to ensure that the monkeys continued to  
573 follow the same rule.

### 574 **More abstract classification tasks**

575 Experiments on more classification tasks (Fig. 6) also largely followed the procedure for the other  
576 tasks described above (Fig. 1B). Training for each task continued until the monkey performance  
577 plateaued but was terminated after six days regardless of their final performance to fairly compare  
578 their performance across tasks with similar training durations. The average performances for the  
579 outdoor vs. indoor, fire- vs. water-related, and Western vs. Eastern-culture tasks at the last day of  
580 training was 91%, 77%, and 88% correct, respectively. The big vs. small task had two versions:  
581 one using the THINGS database and one using the images from Konkle et al. (2012). We started

582 with the THINGS set (final performance during training: 94%) and, after its generalization test,  
583 switched to training on the Konkle set (final performance during training: 93%).

## 584 **Random stimulus-response association tasks**

585 We conducted two versions of control experiments in which monkeys were trained to associate  
586 random object images with targets to check how the absence of a common conceptual rule im-  
587 paired their performance. In one version (Fig. 3A), the associations of images and targets were  
588 randomized such that even images of the same object category could be associated with differ-  
589 ent targets. In the second version (Fig. 3B), two targets were associated with random concrete  
590 object concepts (e.g., “deer”, “banana”; we refer to them as “concrete categories” here to avoid  
591 confusion), but images in each concrete category were associated with the same target. Thus, if  
592 monkeys remembered the association of concrete categories and targets, they could solve the task  
593 and generalize the rule to novel stimuli selected from the same concrete categories.

594 Both versions of the task used the same number of grayscale object images as the animate  
595 vs. inanimate task (96 images in the training set) and followed the same training procedure. The  
596 training image set for the second version (Fig. 3B) contained 16 concrete categories for each target  
597 (2-6 images per concrete category, 48 images per target). As in the other grayscale object tasks  
598 (Fig. 1), these concrete categories were sampled from humans, monkeys, non-primate mammals,  
599 reptiles/amphibians, food, plants, tools/everyday objects, and electronics. We also constructed a  
600 generalization image set with the same number of images and object categories. All the images  
601 used in these control experiments differed from those used in the main tasks.

## 602 **Human experimental procedures**

### 603 **Object drag tasks using grayscale image sets**

604 We performed the six grayscale image tasks (Fig. 1B, G) with human participants (Fig. 5). Seven  
605 participants were assigned to each task. In each trial, an image from the training image set was  
606 randomly presented, and the participants had to choose a target as in the monkey experiments. The  
607 participants were not informed of the categorization rule and had to infer it from the feedback,  
608 but they quickly understood the rule in all tasks (Fig. 5A). Then, each participant continued the  
609 task for ~1,500 trials, resulting in 15-25 repetitions per image. We also conducted the animate vs.  
610 inanimate task using three control image sets (cartoons, silhouettes, outlines; Fig. 5G, H). Four  
611 participants were recruited for each set, and each participant performed ~1,500 trials. In total,  
612 83,416 trials were collected across all tasks and participants.

## 613 **Object drag tasks with generalization images**

614 We also measured the generalization performance of the human participants to novel stimuli in  
615 the various tasks we tested with the monkeys (15 tasks in total; Fig. 6). Nine participants were  
616 assigned to each task. For each task, we first asked the participants to categorize the training image  
617 set. We continued training until the participants reached 90% accuracy (20-trial sliding window)  
618 and viewed all the training images (65-183 trials). We then showed the generalization images only  
619 once. The same training and generalization images were used as those in the monkey experiments.  
620 In total, 13,754 trials for the generalization images were collected across all tasks and participants.

## 621 **Data analysis**

### 622 **Choice accuracy**

623 To plot accuracy as a function of the number of repeated stimulus presentations in Figure 1D  
624 inset, we computed the correct rate using all stimuli in the two-target trials at a given number of  
625 repetitions across training sessions. The correct rate for each image (Fig. 1E) was calculated using  
626 the sessions after the monkey performance reached around 75% correct (the last 3-5 sessions). The  
627 generalization performances for novel images and manipulated images (Figs. 1I, 2B-D, 3D, 6B)  
628 were all calculated using only the first trial of each stimulus.

## 629 **Model comparison**

### 630 **Category classification and behavioral fitting procedures**

631 We evaluated the performance of various vision models in classifying object concepts and fitting  
632 monkey behavior by constructing linear classifiers based on these model outputs (Fig. 4, Supple-  
633 mentary Figs. 4, 5). The outputs of many models are highly overparameterized ( $10^2 - 10^5$  param-  
634 eters), whereas the number of images we used in each task was limited ( $10^2 - 10^4$ ). Therefore, we  
635 first performed a dimension reduction on the model outputs and then employed a cross-validation  
636 approach to evaluate their performance. The cross validation was performed by constructing a  
637 classifier using the training set of each task and calculating the model's correct rate using the  
638 generalization set, mirroring the monkey experimental procedure. For dimension reduction, we  
639 performed principal component analysis (PCA) on the model outputs of the training set and se-  
640 lected the first 20 PCs. The number of dimensions did not qualitatively affect our conclusions. We  
641 then normalized each dimension and constructed a linear classifier of object concepts or monkey  
642 choices using a logistic regression. Finally, we projected the model outputs for the generalization  
643 set on the same PC space and calculated the proportion of stimuli correctly classified by the trained  
644 linear classifier.

645 A linear classifier was constructed by fitting the training data to the following linear logistic  
646 function:

$$\text{logit}[P_i(\text{category}1)] = M_i \cdot w + b \quad (1)$$

647 where  $P_i(\text{category}1)$  is the probability of classifying image  $i$  into category 1,  $M_i$  is the dimension-  
648 reduced model output for image  $i$ ,  $w$  is the linear weight, and  $b$  is a bias parameter. We then  
649 computed  $P_i(\text{category}1)$  for each of the generalization images and calculated the correct rate  
650 assuming that the model chooses category 1 if  $P_i(\text{category}1)$  is greater than 0.5. For the fitting  
651 to monkey behavior, we again performed the linear logistic regression but estimated  $w$  and  $b$  to  
652 maximize the fit performance to monkey choices in the training trials of the sessions after reaching  
653 a plateau performance. The fit performance was quantified as the proportion of generalization  
654 images for which the monkey and model choices matched. In the figures, we showed the fit results  
655 aggregated across the three monkeys.

656 We confirmed that changing the details of this model evaluation procedure did not qualita-  
657 tively affect our results. First, we tested different numbers of dimensions being reduced by PCA  
658 (Supplementary Fig. 4A). All the models quickly improved their cross-validated performance as  
659 the number of dimensions increased, and their performance reached a plateau before or near 20  
660 dimensions. Second, we tested another dimension reduction technique. We constructed a repre-  
661 sentational dissimilarity matrix of the images by computing Pearson's correlation of the model  
662 outputs for each image pair. We then performed a classical multi-dimensional scaling (MDS) to  
663 obtain a 20-dimensional space capturing the similarity pattern. The classification based on these 20  
664 dimensions was similar to the main results (Supplementary Fig. 4B). Third, we replaced logistic  
665 regression with a linear support vector machine (SVM) and found that the classification accuracy  
666 did not substantially change (Supplementary Fig. 4C).

## 667 **Model details**

668 **V1-like model** We followed Pinto et al. (2008) to create a population of V1-like responses,  
669 which is generated from a set of Gabor filters of different orientations and spatial frequencies  
670 with normalization and non-linear thresholding. In brief, we used Gabor filters ( $43 \times 43$  pixels)  
671 spanning 16 orientations and six spatial frequencies (1/2, 1/3, 1/4, 1/6, 1/11, 1/18 cycles/pixel)  
672 with a fixed Gaussian envelope (standard deviation of 9 cycles/pixel in both directions) and fixed  
673 phase (0) (96 filters in total).

674 **Luminance+color model** The model computes marginal statistics of color and luminance infor-  
675 mation at both the global and local levels. Each pixel intensity of an input image was converted  
676 into luminance ( $cd/m^2$ ) and CIE- $u'v'$  coordinate via the calibration data of the touchscreen moni-  
677 tors obtained using a spectrometer (X-Rite i1 Pro, X-Rite, Grand Rapids, MI, USA). At the global  
678 level, the mean, standard deviation, skewness, and kurtosis of the distributions of luminance,  $u'$ ,  
679 and  $v'$  were calculated for the entire image. At the local level, we segmented the image into  $4 \times$   
680 4 blocks and calculated these statistics for each segment. All of these values were normalized and  
681 then concatenated as the output of the model.

682 **Gist+SF model** The Gist model, adopted from Oliva et al. (2006), captures the spatial struc-  
683 ture of an image by computing the orientation and spatial-frequency powers of each local image  
684 region. Although the extracted features are similar to those of the V1-like model, this model cap-  
685 tures larger-scale image structures with more compact representations. In brief, each image was  
686 divided into  $4 \times 4$  segments and oriented Gabor filters (8 orientations) were applied over different  
687 scales (4 scales) in each segment. The average filter energy was subsequently calculated for each  
688 filter and segment, and all outputs were concatenated. To capture global image patterns, we also  
689 calculated the spatial frequency and orientation power (4 frequencies and 8 orientations) of the  
690 whole image. The images were converted to grayscale before the calculation and all the statistics  
691 were normalized.

692 **Texture model** The texture model consists of statistics that capture the correlation structures of  
693 V1-level filter outputs across different scales, orientations, and positions (Portilla and Simoncelli  
694 2000). These statistics explain part of the neural responses in macaque V2 and V4 (Freeman et  
695 al. 2013; Kim et al. 2019; Okazawa et al. 2015). We adopted the code of Portilla et al. (2000),  
696 which used steerable pyramids to compute the outputs of oriented bandpass filters. We chose 4  
697 orientations, 4 scales, and 7 spatial neighborhoods to compute spatial correlations. The texture  
698 model also included marginal statistics of the luminance distribution and the powers of the V1-  
699 level filter outputs.

700 **Contour model** For the background-cropped image sets, we computed statistics derived from the  
701 outline contour of an object image (Supplementary Fig. 5A). The statistics included the aspect ratio  
702 (horizontal vs. vertical extent of the contour), the angle of the longitudinal axis (the most elongated  
703 axis), the extent of elongation (the ratio between the most elongated axis and its orthogonal axis),  
704 the distance of the contour from the image origin at each 5-degree angle, the standard deviation  
705 of these distances, and the spatial frequency powers of the contour shape at 4 scales. Thus, these  
706 parameters quantified information such as to how much and in which direction the object contour  
707 is elongated and how smooth the contour is, as well as the coarsely sampled location of the contour  
708 itself.

709 **Deep neural network (DNN) models** We used seven DNN models commonly used in vision  
710 science research —AlexNet (Krizhevsky et al. 2012), VGG-16 (Simonyan and Zisserman 2014),  
711 ResNet-50 (He et al. 2015), Vision Transformer (ViT-B/32; Dosovitskiy et al. 2020), DINO(Caron  
712 et al. 2021), CLIP (Radford et al. 2021) and SigLIP2 (Tschannen et al. 2025). AlexNet, VGG16,  
713 and ResNet-50 are convolutional neural networks (CNNs) with different numbers of convolutional  
714 layers and network architectures. We used these networks pre-trained on 1,000 object categories in  
715 the ImageNet1k dataset. Using the hook mechanism in PyTorch (version 2.0.1), we extracted the  
716 activity of the penultimate layer for each image as model outputs. Vision Transformer was used to  
717 demonstrate how a DNN with a very different architecture from CNNs performs in our tasks. We  
718 extracted the activity of the “encoder” block of ViT-B/32 for each image. DINO is a self-supervised  
719 vision transformer that learns visual features from unlabeled images. We chose DINO with ViT-  
720 B/16 as the architecture. CLIP and SigLIP2 are recent DNNs that learn visual representations by

721 aligning feature spaces between image and text data (Radford et al. 2021; Tschannen et al. 2025).  
722 We chose CLIP and SigLIP2 with ViT-B/32 as the architecture of visual processing and extracted  
723 the activity of its visual encoder for each image.

724 **Ventral-pathway neural responses** We also compared our behavioral data with neural data from  
725 the THINGS ventral stream spiking dataset (TVSD; Papale et al. 2025). These neural data were  
726 obtained from two macaque monkeys viewing images from the same THINGS database we used,  
727 with a large number of microelectrodes implanted in V1, V4, and inferotemporal (IT) cortex. The  
728 dataset includes normalized multi-unit activity (MUA) averaged in a time window centered on the  
729 peak response of each region (25-125 ms for V1, 50-150 ms for V4, and 75-175 ms for IT). We  
730 used this normalized MUA data to assess how well neural population (V1: 1024 units, V4: 448  
731 units, IT: 576 units) could classify object categories and fit our monkey behavioral data, using the  
732 same procedure we adopted for the model outputs.

## 733 **The drift-diffusion model (DDM) fit**

734 To quantify the difficulty of categorizing each stimulus, we used DDM, which can jointly fit  
735 choices and RTs. The model assumes that the drift rate (category sensitivity) varies across im-  
736 ages and reflects stimulus difficulty, whereas the decision bound and non-decision time are stable  
737 across images for each participant and task. For each image, the drift-diffusion process forms the  
738 decision variable at time  $t$  ( $v(t)$ ) as

$$v(t) = \int_0^t dr_i + \eta(\tau) d\tau \quad (2)$$

739 where  $dr_i$  is the drift rate (category sensitivity) for image  $i$ .  $\eta(\tau)$  represents a noise term following  
740 a Gaussian distribution with a mean of 0 and an SD of 1. When  $v(t)$  reaches a positive ( $+B$ ) or  
741 negative ( $-B$ ) bound, the model commits to the decision associated with the bound (category 1  
742 or 2). The reaction times are modeled as the time required to reach a bound plus an additional  
743 delay (non-decision time) common to all images. We modeled the non-decision time as following  
744 a Gaussian distribution with a mean of  $T_0$  and an SD of  $T_0/3$ .

745 The model had two free parameters ( $B, T_0$ ) common to all images and the drift rate (category  
746 sensitivity;  $dr_i$  in Eq. 2) as a free parameter for each image (total  $n + 2$  free parameters;  $n$  is the  
747 number of images). Higher drift rates indicate that images were easier to classify. We fitted the  
748 model to each participant's choices and RT distributions in each task using maximum likelihood  
749 estimation (Luo et al. 2025; Zheng et al. 2025). Given a set of parameters, the RT distributions  
750 for the two choices were calculated using the aforementioned model formulation by numerically  
751 solving the Fokker-Planck equation. From these distributions, we derived the likelihood of ob-  
752 serving the participants' choices and RTs. Summing the likelihoods over all trials yielded the total  
753 likelihood of the parameter set. We used a simplex search method (*fminsearch* in MATLAB) to  
754 determine the parameter set that maximized the summed likelihood. As shown in Supplementary

755 Fig. 6, our model accurately fit both monkey and human RTs and correct rates, thereby justify-  
756 ing the use of the drift rate as a metric of stimulus difficulty. Figure 5E-H showed the drift rate  
757 averaged across participants.

## 758 Dissimilarity matrix

759 We use the generalization performance of images across all tasks to construct the dissimilarity ma-  
760 trix among humans, monkeys, and models (Fig. 6C). As in Fig. 6B, we focused on the first trial  
761 of each generalization image and calculated the average correct rate across subjects for humans  
762 and monkeys. For the models, we calculated the probability of correct classification for the same  
763 image based on the likelihood derived from a logistic classifier trained on the training images. We  
764 then concatenated all image performances across tasks. The dissimilarity was defined as Euclidean  
765 distance; thus, both the overall correct rates and the patterns of performance across images influ-  
766 enced the values. Using 1– Pearson’s correlation yielded similar results, but low- and mid-level  
767 models became more distant from each other because they exhibited distinct error patterns. t-SNE  
768 (Fig. 6E) also used Euclidean distance, and its perplexity parameter was fixed to 4.

## 769 References

- 770 Baker, N., Lu, H., Erlikhman, G., & Kellman, P. J. (2018). Deep convolutional networks do not  
771 classify based on global object shape. *PLoS Computational Biology*, *14*(12), e1006613.
- 772 Butler, J. L., & Kennerley, S. W. (2019). Mymou: A low-cost, wireless touchscreen system for  
773 automated training of nonhuman primates. *Behavior Research Methods*, *51*(6), 2559–2572.
- 774 Carlson, T. A., Ritchie, J. B., Kriegeskorte, N., Durvasula, S., & Ma, J. (2014). Reaction time  
775 for object categorization is predicted by representational distance. *Journal of Cognitive  
776 Neuroscience*, *26*(1), 132–142.
- 777 Caron, M., Touvron, H., Misra, I., Jégou, H., Mairal, J., Bojanowski, P., & Joulin, A. (2021).  
778 Emerging properties in self-supervised vision transformers. *arXiv*. [https://arxiv.org/abs/  
779 2104.14294](https://arxiv.org/abs/2104.14294)
- 780 Cherian, T., Jacob, G., & Arun, S. (2025). Do monkeys see the way we do? qualitative similarities  
781 and differences between monkey and human perception. *bioRxiv*, 2025–01. [https://doi.org/  
782 10.1101/2025.01.30.635614](https://doi.org/10.1101/2025.01.30.635614)
- 783 Clarke, A., & Tyler, L. K. (2015). Understanding what we see: How we derive meaning from  
784 vision. *Trends in Cognitive Sciences*, *19*(11), 677–687.
- 785 Cromer, J. A., Machon, M., & Miller, E. K. (2011). Rapid association learning in the primate  
786 prefrontal cortex in the absence of behavioral reversals. *Journal of Cognitive Neuroscience*,  
787 *23*(7), 1823–1828.
- 788 Delorme, A., Richard, G., & Fabre-Thorpe, M. (2000). Ultra-rapid categorisation of natural scenes  
789 does not rely on colour cues: A study in monkeys and humans. *Vision Research*, *40*(16),  
790 2187–2200.

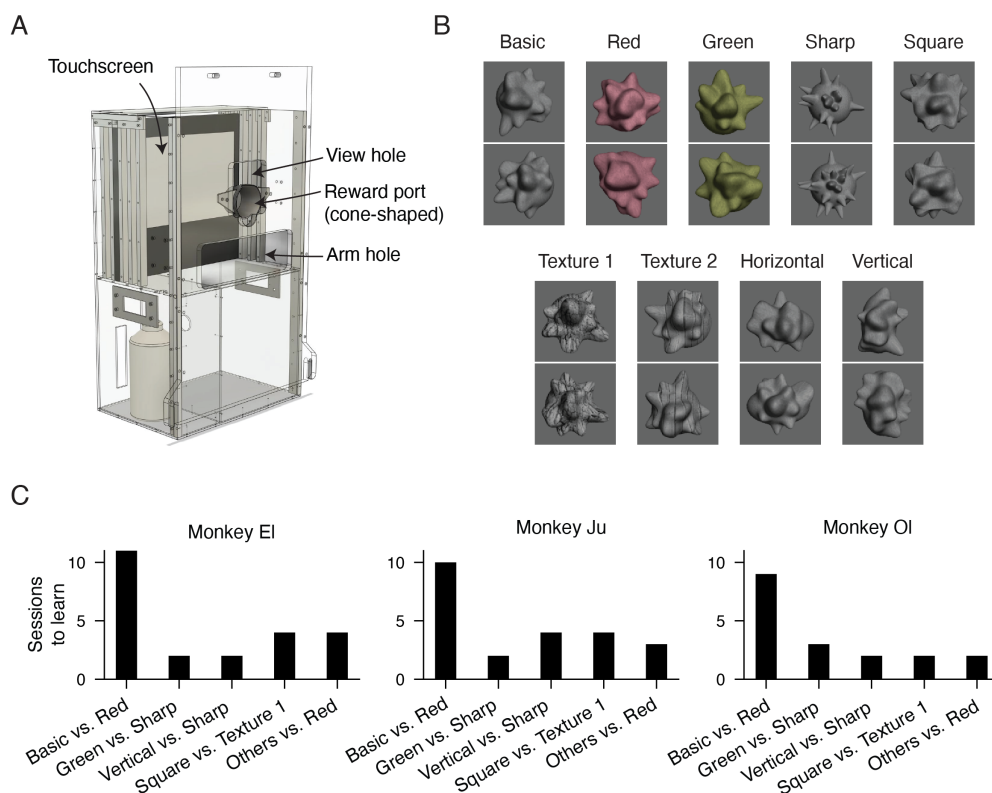
- 791 Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani,  
792 M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., & Houlsby, N. (2020). An image is  
793 worth 16x16 words: Transformers for image recognition at scale. *arXiv*. [https://doi.org/10.](https://doi.org/10.48550/arxiv.2010.11929)  
794 [48550/arxiv.2010.11929](https://doi.org/10.48550/arxiv.2010.11929)
- 795 Du, C., Fu, K., Wen, B., Sun, Y., Peng, J., Wei, W., Gao, Y., Wang, S., Zhang, C., Li, J., et al.  
796 (2025). Human-like object concept representations emerge naturally in multimodal large  
797 language models. *Nature Machine Intelligence*, 7, 860–875.
- 798 Duñabeitia, J. A., Crepaldi, D., Meyer, A. S., New, B., Pliatsikas, C., Smolka, E., & Brysbaert,  
799 M. (2018). Multipic: A standardized set of 750 drawings with norms for six european lan-  
800 guages. *Quarterly Journal of Experimental Psychology*, 71(4), 808–816.
- 801 Eldridge, M. A., Matsumoto, N., Wittig, J. H., Masseau, E. C., Saunders, R. C., & Richmond, B. J.  
802 (2018). Perceptual processing in the ventral visual stream requires area TE but not rhinal  
803 cortex. *eLife*, 7, e36310.
- 804 Fabre-Thorpe, M. (2003). Visual categorization: Accessing abstraction in non-human primates.  
805 *Philosophical Transactions of the Royal Society of London. Series B: Biological Sciences*,  
806 358(1435), 1215–1223.
- 807 Fabre-Thorpe, M., Richard, G., & Thorpe, S. J. (1998). Rapid categorization of natural images by  
808 rhesus monkeys. *Neuroreport*, 9(2), 303–308.
- 809 Fize, D., Cauchoix, M., & Fabre-Thorpe, M. (2011). Humans and monkeys share visual represen-  
810 tations. *Proceedings of the National Academy of Sciences, USA*, 108(18), 7635–7640.
- 811 Freedman, D. J., Riesenhuber, M., Poggio, T., & Miller, E. K. (2001). Categorical representation  
812 of visual stimuli in the primate prefrontal cortex. *Science*, 291(5502), 312–316.
- 813 Freeman, J., Ziemba, C. M., Heeger, D. J., Simoncelli, E. P., & Movshon, J. A. (2013). A functional  
814 and perceptual signature of the second visual area in primates. *Nature Neuroscience*, 16(7),  
815 974–981.
- 816 Grill-Spector, K., & Kanwisher, N. (2005). Visual recognition: As soon as you know it is there,  
817 you know what it is. *Psychological Science*, 16(2), 152–160.
- 818 He, K., Zhang, X., Ren, S., & Sun, J. (2015). Deep residual learning for image recognition. *arXiv*.  
819 <https://doi.org/10.48550/arxiv.1512.03385>
- 820 Hebart, M. N., Dickter, A. H., Kidder, A., Kwok, W. Y., Corriveau, A., Van Wicklin, C., & Baker,  
821 C. I. (2019). THINGS: A database of 1,854 object concepts and more than 26,000 natural-  
822 istic object images. *PLoS One*, 14(10), e0223792.
- 823 Hebart, M. N., Zheng, C. Y., Pereira, F., & Baker, C. I. (2020). Revealing the multidimensional  
824 mental representations of natural objects underlying human similarity judgements. *Nature*  
825 *human behaviour*, 4(11), 1173–1185.
- 826 Heidari-Gorji, H., Ebrahimpour, R., & Zabbah, S. (2021). A temporal hierarchical feedforward  
827 model explains both the time and the accuracy of object recognition. *Scientific Reports*,  
828 11(1), 5640.
- 829 Kar, K., & DiCarlo, J. J. (2024). The quest for an integrated set of neural mechanisms underlying  
830 object recognition in primates. *Annual Review of Vision Science*, 10, 91–121.
- 831 Kaufman, M. T., Churchland, M. M., Ryu, S. I., & Shenoy, K. V. (2015). Vacillation, indecision  
832 and hesitation in moment-by-moment decoding of monkey motor cortex. *eLife*, 4, e04677.
- 833 Kawaguchi, K., Pourriahi, P., Seillier, L., Clery, S., & Nienborg, H. (2019). Easily adaptable head-  
834 free training system of macaques for tasks requiring precise measurements of eye position.  
835 *bioRxiv*. <https://doi.org/10.1101/588566>

- 836 Khaligh-Razavi, S.-M., & Kriegeskorte, N. (2014). Deep supervised, but not unsupervised, models  
837 may explain it cortical representation. *PLoS Computational Biology*, *10*(11), e1003915.
- 838 Kiani, R., Esteky, H., Mirpour, K., & Tanaka, K. (2007). Object category structure in response  
839 patterns of neuronal population in monkey inferior temporal cortex. *Journal of Neurophys-*  
840 *iology*, *97*(6), 4296–4309.
- 841 Kim, T., Bair, W., & Pasupathy, A. (2019). Neural coding for shape and texture in macaque area  
842 V4. *Journal of Neuroscience*, *39*(24), 4760–4774.
- 843 Konkle, T., & Oliva, A. (2012). A real-world size organization of object responses in occipitotem-  
844 poral cortex. *Neuron*, *74*(6), 1114–1124.
- 845 Koopman, S. E., Mahon, B. Z., & Cantlon, J. F. (2017). Evolutionary constraints on human object  
846 perception. *Cognitive Science*, *41*(8), 2126–2148.
- 847 Kramer, L. E., Chen, Y.-C., Long, B., Konkle, T., & Cohen, M. R. (2023). Contributions of early  
848 and mid-level visual cortex to high-level object categorization. *bioRxiv*. [https://doi.org/10.](https://doi.org/10.1101/2023.05.31.541514)  
849 [1101/2023.05.31.541514](https://doi.org/10.1101/2023.05.31.541514)
- 850 Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). Imagenet classification with deep convolu-  
851 tional neural networks. *Advances in Neural Information Processing Systems*, *25*.
- 852 Loggia, S. R., Feibes, H. E., Duffield, S., Kasi, K., & Conway, B. R. (2025). The representation of  
853 object concepts across the brain. *bioRxiv*. <https://doi.org/10.64898/2025.12.07.692815>
- 854 Long, B., Konkle, T., Cohen, M. A., & Alvarez, G. A. (2016). Mid-level perceptual features distin-  
855 guish objects of different real-world sizes. *Journal of Experimental Psychology: General*,  
856 *145*(1), 95.
- 857 Long, B., Störmer, V. S., & Alvarez, G. A. (2017). Mid-level perceptual features contain early cues  
858 to animacy. *Journal of Vision*, *17*(6), 20–20.
- 859 Long, B., Yu, C.-P., & Konkle, T. (2018). Mid-level visual features underlie the high-level categor-  
860 ical organization of the ventral stream. *Proceedings of the National Academy of Sciences,*  
861 *USA*, *115*(38), E9015–E9024.
- 862 Luo, T., Xu, M., Zheng, Z., & Okazawa, G. (2025). Limitation of switching sensory information  
863 flow in flexible perceptual decision making. *Nature Communications*, *16*(1), 172.
- 864 Macé, M. J.-M., Joubert, O. R., Nespoulous, J.-L., & Fabre-Thorpe, M. (2009). The time-course  
865 of visual categorizations: You spot the animal faster than the bird. *PloS one*, *4*(6), e5927.
- 866 Mack, M. L., & Palmeri, T. J. (2011). The timing of visual object categorization. *Frontiers in*  
867 *Psychology*, *2*, 165.
- 868 Minamimoto, T., Saunders, R. C., & Richmond, B. J. (2010). Monkeys quickly learn and generalize  
869 visual categories without lateral prefrontal cortex. *Neuron*, *66*(4), 501–507.
- 870 Murphy, G. (2004). *The big book of concepts*. MIT press.
- 871 Murphy, G. L., & Medin, D. L. (1985). The role of theories in conceptual coherence. *Psychological*  
872 *Review*, *92*(3), 289.
- 873 Muttenthaler, L., Greff, K., Born, F., Spitzer, B., Kornblith, S., Mozer, M. C., Müller, K.-R., Un-  
874 terthiner, T., & Lampinen, A. K. (2025). Aligning machine and human visual representa-  
875 tions across abstraction levels. *Nature*, *647*(8089), 349–355.
- 876 Okazawa, G., Tajima, S., & Komatsu, H. (2015). Image statistics underlying natural texture selec-  
877 tivity of neurons in macaque v4. *Proceedings of the National Academy of Sciences, USA*,  
878 *112*(4), E351–E360.
- 879 Oliva, A., & Torralba, A. (2006). Building the gist of a scene: The role of global image features in  
880 recognition. *Progress in Brain Research*, *155*, 23–36.

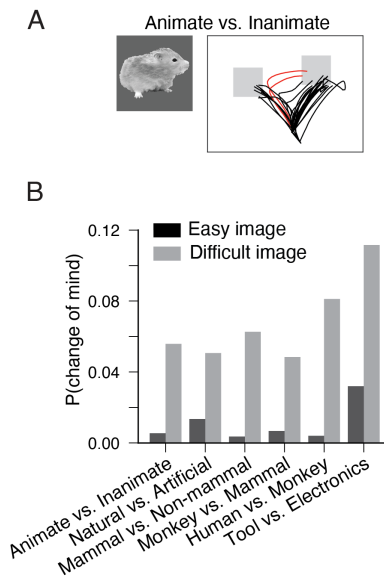
- 881 Orban, G. A., Van Essen, D., & Vanduffel, W. (2004). Comparative mapping of higher visual areas  
882 in monkeys and humans. *Trends in Cognitive Sciences*, 8(7), 315–324.
- 883 Papale, P., Wang, F., Self, M. W., & Roelfsema, P. R. (2025). An extensive dataset of spiking  
884 activity to reveal the syntax of the ventral stream. *Neuron*, 113(4), 539–553.
- 885 Pinto, N., Cox, D. D., & DiCarlo, J. J. (2008). Why is real-world visual object recognition hard?  
886 *PLoS Computational Biology*, 4(1), e27.
- 887 Portilla, J., & Simoncelli, E. P. (2000). A parametric texture model based on joint statistics of  
888 complex wavelet coefficients. *International Journal of Computer Vision*, 40, 49–70.
- 889 Quiroga, R. Q., Reddy, L., Kreiman, G., Koch, C., & Fried, I. (2005). Invariant visual representa-  
890 tion by single neurons in the human brain. *Nature*, 435(7045), 1102–1107.
- 891 Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A.,  
892 Mishkin, P., Clark, J., Krueger, G., & Sutskever, I. (2021). Learning transferable visual  
893 models from natural language supervision. *arXiv*. [https://doi.org/10.48550/arxiv.2103.](https://doi.org/10.48550/arxiv.2103.00020)  
894 00020
- 895 Rajalingham, R., Issa, E. B., Bashivan, P., Kar, K., Schmidt, K., & DiCarlo, J. J. (2018). Large-  
896 scale, high-resolution comparison of the core visual object recognition behavior of humans,  
897 monkeys, and state-of-the-art deep artificial neural networks. *Journal of Neuroscience*,  
898 38(33), 7255–7269.
- 899 Rajalingham, R., Schmidt, K., & DiCarlo, J. J. (2015). Comparison of object recognition behavior  
900 in human and monkey. *Journal of Neuroscience*, 35(35), 12127–12136.
- 901 Ralph, M. A. L., Jefferies, E., Patterson, K., & Rogers, T. T. (2017). The neural and computational  
902 bases of semantic cognition. *Nature Reviews Neuroscience*, 18(1), 42–55.
- 903 Ramezani, H., Giverin, C., & Kar, K. (2024). Low-cost, portable, easy-to-use kiosks to facili-  
904 tate home-cage testing of nonhuman primates during vision-based behavioral tasks. *Journal*  
905 *of Neurophysiology*, 132(3), 666–677.
- 906 Rosch, E., Mervis, C. B., Gray, W. D., Johnson, D. M., & Boyes-Braem, P. (1976). Basic objects  
907 in natural categories. *Cognitive Psychology*, 8(3), 382–439.
- 908 Santos, L. R., Hauser, M. D., & Spelke, E. S. (2001). Recognition and categorization of biologi-  
909 cally significant objects by rhesus monkeys (*macaca mulatta*): The domain of food. *Cogni-*  
910 *tion*, 82(2), 127–155.
- 911 Schrier, A. M., & Brady, P. M. (1987). Categorization of natural stimuli by monkeys (*macaca mu-*  
912 *latta*): Effects of stimulus set size and modification of exemplars. *Journal of Experimental*  
913 *Psychology: Animal Behavior Processes*, 13(2), 136.
- 914 Simonyan, K., & Zisserman, A. (2014). Very Deep Convolutional Networks for Large-Scale Image  
915 Recognition. *arXiv*. <https://doi.org/10.48550/arxiv.1409.1556>
- 916 Stine, G. M., Zylberberg, A., Ditterich, J., & Shadlen, M. N. (2020). Differentiating between inte-  
917 gration and non-integration strategies in perceptual decision making. *eLife*, 9, e55365.
- 918 Stoinski, L. M., Perkuhn, J., & Hebart, M. N. (2024). THINGSplus: New norms and metadata for  
919 the things database of 1854 object concepts and 26,107 natural object images. *Behavior*  
920 *Research Methods*, 56(3), 1583–1603.
- 921 Tschannen, M., Gritsenko, A., Wang, X., Naeem, M. F., Alabdulmohsin, I., Parthasarathy, N.,  
922 Evans, T., Beyer, L., Xia, Y., Mustafa, B., Hénaff, O., Harmsen, J., Steiner, A., & Zhai,  
923 X. (2025). Siglip 2: Multilingual vision-language encoders with improved semantic under-  
924 standing, localization, and dense features. *arXiv*. <https://arxiv.org/abs/2502.14786>

- 925 Vogels, R. (1999). Categorization of complex visual images by rhesus monkeys. part 1:  
926 Behavioural study. *European Journal of Neuroscience*, *11*(4), 1223–1238.
- 927 Vong, W. K., Wang, W., Orhan, A. E., & Lake, B. M. (2024). Grounded language acquisition  
928 through the eyes and ears of a single child. *Science*, *383*(6682), 504–511.
- 929 Yamins, D. L., & DiCarlo, J. J. (2016). Using goal-driven deep learning models to understand  
930 sensory cortex. *Nature Neuroscience*, *19*(3), 356–365.
- 931 Yamins, D. L., Hong, H., Cadieu, C. F., Solomon, E. A., Seibert, D., & DiCarlo, J. J. (2014).  
932 Performance-optimized hierarchical models predict neural responses in higher visual cor-  
933 tex. *Proceedings of the National Academy of Sciences, USA*, *111*(23), 8619–8624.
- 934 Yoshikubo, S. (1985). Species discrimination and concept formation by rhesus monkeys (macaca  
935 mulatta). *Primates*, *26*, 285–299.
- 936 Zheng, Z., Hu, J., & Okazawa, G. (2025). Spatiotemporal evidence accumulation through saccadic  
937 sampling for object recognition. *Journal of Neuroscience*, *45*(44).
- 938 Zhou, B., Lapedriza, A., Khosla, A., Oliva, A., & Torralba, A. (2017). Places: A 10 million im-  
939 age database for scene recognition. *IEEE Transactions on Pattern Analysis and Machine*  
940 *Intelligence*, *40*(6), 1452–1464.

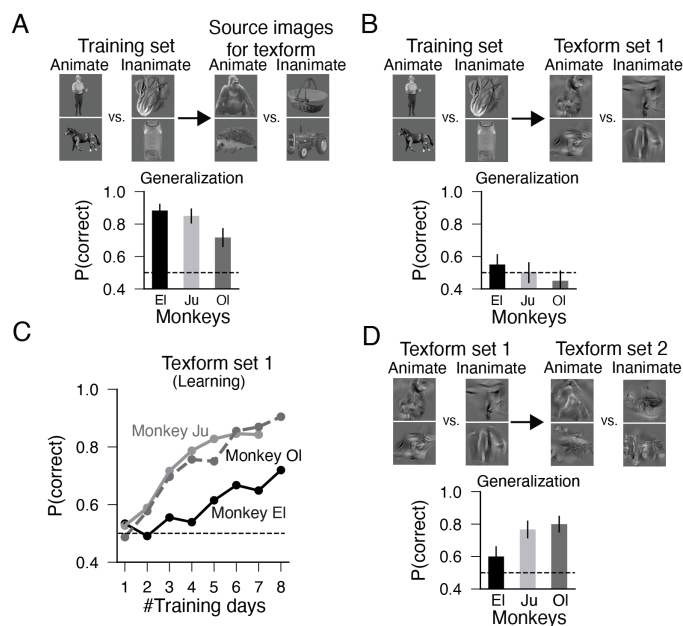
## Supplementary figures



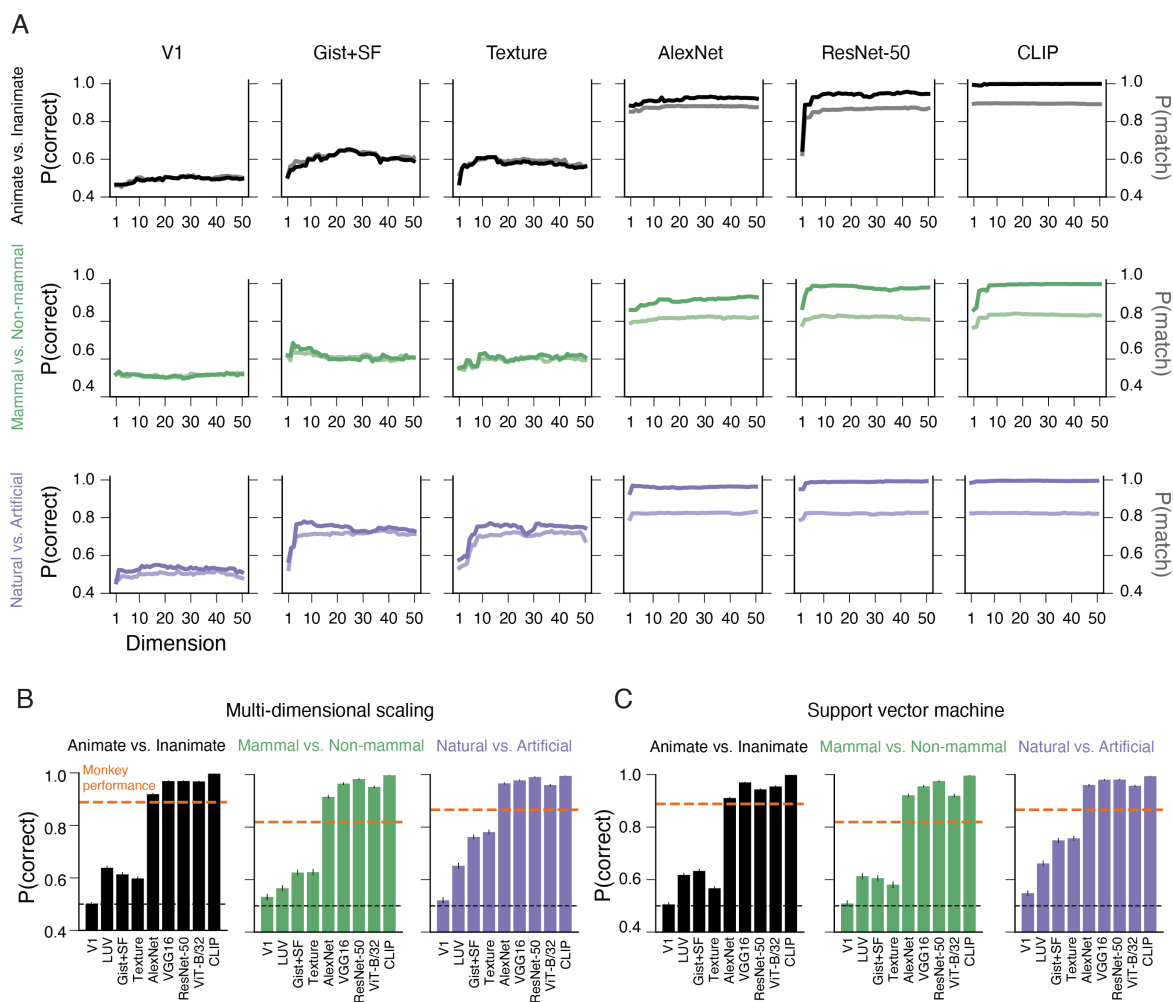
**Figure S1: Pretraining of monkeys on the task structure using shape stimuli.** (A) Behavioral testing system attached to the monkeys' home cage. Monkeys reached for a touchscreen through an arm hole in the front panel and received a juice reward from a cone-shaped reward socket (Kawaguchi et al. 2019). While performing tasks, they placed their mouth on the reward socket and looked at the screen through a viewing hole. (B) Prior to the main concept classification tasks, we trained the monkeys to perform the classification of amorphous shape stimuli. The task sequence followed the main design (Fig. 1A), but the monkeys had to categorize shape stimuli according to their shape, color, or texture. We created nine groups of stimuli as shown in this panel and chose two of them to classify for each task. In each group, there were 50 images with random bump positions. (C) The number of sessions (days) required for the monkeys to learn each shape classification task plotted in the trained order. Prior to the first task, monkeys were trained to touch a shape stimulus and move it to a target box. Then, in the first "Basic" vs. "Red" task, they were presented with two target boxes and had to learn to select one of them based on a stimulus feature. This initial stimulus-response association training took 9-11 days (see Methods). The monkeys were subsequently trained to classify different pairs of shapes in 2-5 days each. In the "Others" vs. "Red" task, the monkeys had to choose one target for red stimuli and the other target for stimuli randomly selected from the remaining eight stimulus categories.



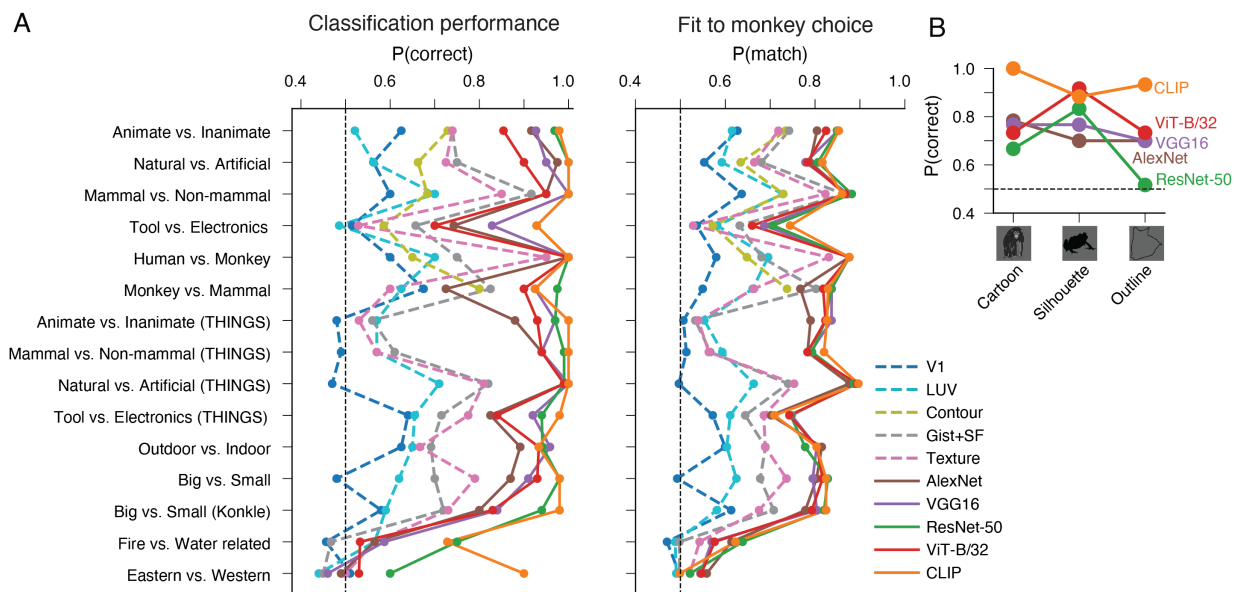
**Figure S2: Monkey reach trajectories reflected their decision-making process.** (A) Monkeys moved a stimulus to a target with a straight path in most trials, but they occasionally made detours as if they changed their mind (red lines; Kaufman et al. 2015). These trials were detected by looking for trajectories whose average position was the opposite to the position of the chosen target with respect to the midline. The right target box corresponded to the animate category in this example. (B) Monkeys exhibited these detour trajectories much more often for difficult images. For each task, we selected the 10 easiest and most difficult images based on correct rates and plotted the probability of trials in which the monkeys exhibited detour trajectories. The plots aggregate the data from the three monkeys.



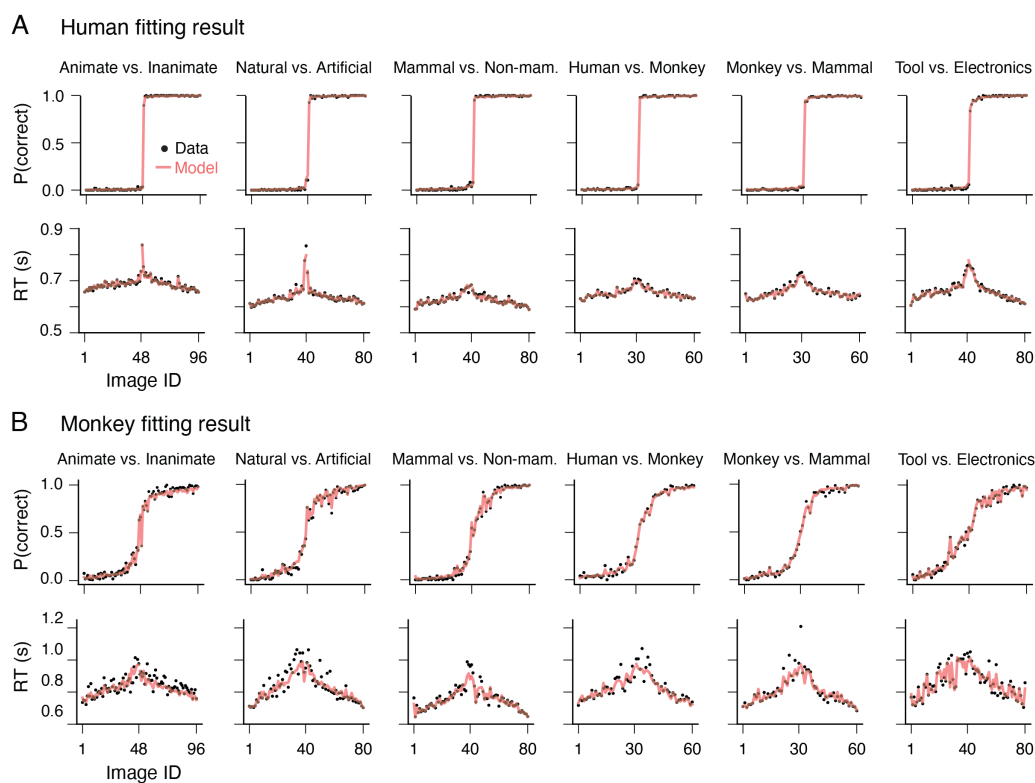
**Figure S3: Monkeys failed to generalize a rule to texform images.** It has been shown that animacy information is partially preserved even in images that retain only mid-level textural visual features (texform; Long et al. 2018). We found that the monkeys failed to generalize the learned animate vs. inanimate rule to texform images. Thus, it is unlikely that the monkeys relied primarily on mid-level textural features in our main task. Error bars are S.E.M. across images. **(A)** We first confirmed that monkeys could generalize the rule to the original, source images (object photographs) used to generate texform images by Long et al. (2018). **(B)** However, when we showed texform images, the monkey performance decreased to the chance level ( $p = 0.26 - 0.82$ ,  $BF_{01} = 4.7 - 6.3$ , binomial test). We used 60 texform images (set 1). **(C)** We then attempted to train the monkeys to discriminate animate vs. inanimate using the same set of texform images (Set 1). Although the monkeys gradually improved their performance was consistently lower than that of the main concept tasks for all monkeys ( $p = < 9.5 \times 10^{-4}$  for the first 4 days). In particular, monkey El presented a much slower learning curve. **(D)** Monkeys Ju and Ol were able to generalize the learned classification to other texform images ( $p < 2.1 \times 10^{-5}$ ), but the performance of monkey El was statistically indistinguishable from chance level ( $p = 0.078$ ,  $BF_{01} = 1.9$ , binomial test).



**Figure S4: Model performance did not depend on classification or dimension reduction methods (A)** We confirmed that the classification and behavioral fitting performances of the models (Fig. 4) did not strongly depend on the number of dimensions of the model parameters reduced by principal component analysis (PCA). In each panel, the darker line represents the model performance in correctly classifying object concepts, and the lighter line represents the model performance in fitting monkey choices. In both cases, the fitting was performed using the training image sets, and the performance was evaluated using the generalization image sets used in the monkey experiments shown in Fig. 2. **(B, C)** Model performance was not strongly dependent on the methods used for dimension reduction and classification. We tested multi-dimensional scaling for dimension reduction (B) and used a support vector machine instead of logistic regression for classification (C) and found that the results did not change qualitatively.



**Figure S5: Model classification and behavioral fitting performance for all tasks. (A)** Summary of model performance for all tasks shown in Figs. 1, 2, and 6. For the tasks using large-scale THINGS images (Fig. 2), we used only the first 100 generalization stimuli to compute the performance for proper comparison; thus, the performances do not exactly match those shown in Fig. 4. The classification performances of ResNet-50 and CLIP are the same as those shown in Fig. 6A. **(B)** DNNs could also classify modified stimuli outside the range of natural images (cartoon, silhouette, and outline images), which we tested with monkeys (Fig. 5G). The performance tended to be lower for the outline images, which is consistent with the monkey and human behavioral results.



**Figure S6: Drift diffusion model (DDM) could accurately fit choices and reaction times. (A, B) Fitting of the DDM (Fig. 5C) to the behavioral data of individual tasks for humans (A) and monkeys (B). The conventions are the same as those in Fig. 5D. Model fit was performed for the tasks using grayscale image sets (Fig. 1), in which the monkeys repeated many trials for the same stimuli. See Methods for the details of the model fit. The image IDs were sorted by the sensitivity parameter ( $dr_i$  in Eq. 2) of the model fit averaged across participants for each task. The fit lines are not smooth because the plots are the average across participants who had different rank orders of sensitivity across image IDs.**

## **Movies 1-3**

Example videos showing the monkeys' performance on the tasks using the large-scale THINGS database (Fig. 2; Movie 1: monkey Ju performing the animate vs. inanimate task, Movie 2: monkey El performing the natural vs. artificial task, Movie 3: monkey Ol performing the mammal vs. non-mammal task). The videos were taken during the generalization session, where half of the images were from the generalization set, which was shown only once in each task. The other half were from the training set. The upper parts of the videos are screenshots of the monitor the monkey was looking at.