

1
2 **Evidence that recurrent circuits are critical to the ventral**
3 **stream's execution of core object recognition behavior**
4
5

6 Kohitij Kar^{1,2*}, Jonas Kubilius^{1,3}, Kailyn Schmidt¹, Elias B. Issa¹⁺, and James J. DiCarlo^{1,2}
7

8 1. McGovern Institute for Brain Research and Department of Brain and Cognitive Sciences,
9 Massachusetts Institute of Technology, Cambridge, MA

10 2. Center for Brains, Minds, and Machines, Massachusetts Institute of Technology, Cambridge,
11 MA

12 3. Brain and Cognition, KU Leuven, Leuven, Belgium
13
14

15 ⁺Current address: Department of Neuroscience, Zuckerman Mind Brain Behavior Institute,
16 Columbia University, NY
17
18
19
20

21 *Correspondence should be addressed to Kohitij Kar,
22

23 McGovern Institute for Brain Research
24 Massachusetts Institute of Technology,
25 77 Massachusetts Institute of Technology, 46-6161,
26 Cambridge, MA 02139.
27 E-mail: kohitij@mit.edu
28

29 **AUTHOR CONTRIBUTIONS**

30 K.K. and J.J.D designed the experiments. K.K., K.S., and E.B.I. carried out the experiments.
31 K.K. and J.K. performed the data analysis and modeling. K.K. and J.J.D. wrote the manuscript.
32

33 **ACKNOWLEDGEMENTS**

34 This research was supported by the Office of Naval Research MURI-114407 (J.J.D), and in part
35 by the US National Eye Institute grants R01-EY014970 (J.J.D.), K99-EY022671 (E.B.I.), and the
36 European Union's Horizon 2020 research and innovation programme under grant agreement No
37 705498 (J.K.). We thank Arash Afraz for his surgical assistance.
38

39 **COMPETING FINANCIAL INTERESTS**

40 The authors declare no competing financial interests.

41 Abstract

42

43

44 Non-recurrent deep convolutional neural networks (DCNNs) are currently the best models of
45 core object recognition; a behavior supported by the densely recurrent primate ventral stream,
46 culminating in the inferior temporal (IT) cortex. Are these recurrent circuits critical to the ventral
47 stream’s execution of this behavior? We reasoned that, if recurrence is critical, then primates
48 should outperform feedforward-only DCNNs for some images, and that these images should
49 require additional processing time beyond the feedforward IT response. Here we first used
50 behavioral methods to discover hundreds of these “challenge” images. Second, using large-
51 scale IT electrophysiology in animals performing core recognition tasks, we observed that
52 behaviorally-sufficient, linearly-decodable object identity solutions emerged ~30ms (on average)
53 later in IT for *challenge* images compared to DCNN and primate performance-matched “control”
54 images. We observed these same late solutions even during passive viewing. Third, consistent
55 with a failure of feedforward computations, the behaviorally-critical late-phase IT population
56 response patterns evoked by the *challenge* images were poorly predicted by DCNN activations.
57 Interestingly, very deep CNNs as well as not-so-deep but recurrent CNNs better predicted these
58 late IT responses, suggesting a functional equivalence between additional nonlinear
59 transformations and recurrence. Our results argue that automatically-evoked recurrent circuits
60 are critical even for rapid object identification. By precisely comparing current DCNNs, primate
61 behavior and IT population dynamics, we provide guidance for future recurrent model
62 development.

63

64 Introduction

65
66 In a single, natural viewing fixation (~200 ms), primates can rapidly identify objects in the central
67 visual field, despite various identity preserving image transformations, a behavior termed core
68 object recognition¹. Understanding the brain mechanisms that seamlessly solve this challenging
69 computational problem has been a key goal of visual neuroscience^{2,3}. Previous studies⁴⁻⁶ have
70 shown that object categories and identities are explicitly represented in the pattern of neural
71 activity in the primate inferior temporal (IT) cortex, and that specific IT neural population codes
72 are sufficient to explain and predict primate core object recognition. Therefore, understanding
73 how the brain solves core object recognition boils down to building a neurally-mechanistic (i.e.
74 neural network) model of the primate ventral stream that, for any image, accurately predicts the
75 neuronal firing rate responses at all levels of the ventral stream, including IT.

76
77 At present, the neural network models that best explain and predict the individual and
78 population responses (image evoked, time averaged firing rates) of primate (macaque) IT
79 neurons have been found in the architectural family of deep convolutional neural networks
80 (DCNNs) trained on object categorization⁷⁻⁹. These neural networks are also the current best
81 predictors of primate behavioral patterns over dozens of core object recognition tasks^{10,11}. All
82 neural networks in this model family are almost entirely feed-forward. Specifically, unlike the
83 ventral stream¹²⁻¹⁵, they lack cortico-cortical feedback circuits, sub-cortical feedback circuits,
84 and medium to long-range intra-area recurrent circuits (as shown in Figure 1A). The short time
85 duration (~200 ms) needed to accomplish accurate core object category and identity inferences
86 in the ventral stream^{4,16,17} suggests the possibility that recurrent-circuit driven computations are
87 not critical for these inferences. In addition, it has been argued that recurrent circuits might
88 operate at much slower time scales¹⁸, and thus may be much more relevant for processes like
89 regulating synaptic plasticity to improve future behavior (learning). Taken together, a promising
90 hypothesis is that core object recognition behavior does not require recurrent processing. The
91 primary aim of this study was to try to falsify this hypothesis, and to provide new constraints to
92 guide further neural network model development.

93
94
95 There is growing evidence that the feedforward DCNNs fall short of accurately predicting image-
96 by-image primate behavior in a variety of situations^{11,19}. We therefore hypothesized that
97 specific images for which the object identities are difficult for non-recurrent DCNNs to solve, but
98 are nevertheless easily solved by primates, might be critically benefiting from recurrent
99 computations in the primates. Furthermore, previous research (for review see²⁰) suggests that
100 the impact of recurrent computations in the ventral stream should be most relevant at later time
101 points in the image driven neural responses. Therefore we reasoned that IT neural population
102 representations of objects in images in which those object inferences critically rely on the
103 recurrent computations will require additional processing time to emerge (beyond the initial
104 evoked IT population response that begins at ~90 ms; feedforward pass).

105

106 To discover such images, we behaviorally compared primates (humans and monkeys) and a
107 particular non-recurrent DCNN (AlexNet 'fc7', ²¹) to identify two groups of images — those for
108 which object identity is easily inferred by the primate brain, but not solved by DCNNs (referred
109 to here as “*challenge images*”), and those for which both primates and models easily infer object
110 identity (referred to here as “*control images*”). To test our neural hypothesis, we simultaneously
111 measured IT population activity in response to each of 1320 images, using chronically implanted
112 multielectrode arrays across IT cortex of both the left and right hemispheres of two monkeys,
113 while monkeys performed an object discrimination task.

114
115 Our results revealed that object identity decodes from IT neural populations for the *challenge*
116 images took an average of ~30ms longer to emerge (~145 ms from stimulus onset) compared to
117 *control* images (~115 ms from stimulus onset). Consistent with previous results, we also found
118 that the top layers of DCNNs optimized for object categorization performance predicted ~50% of
119 IT image-driven neural response variance at the leading edge of the IT population response.
120 However, this fit to the IT response was significantly worse (<20% explained variance) at later
121 time points (150-200 ms post stimuli onset) — the time points where linear decoders show that
122 the IT population solutions to many of the *challenge* images emerge. Taken together, these
123 results argue against feedforward only models for the brain's execution of core object
124 recognition, and instead imply a behaviorally-critical role of recurrent computations. Notably, we
125 also found the same neural population phenomena while the monkeys passively viewed the
126 images, implying that the putative recurrent mechanisms for successful core object inference in
127 the primate are automatic and rapid, and not strongly state or task dependent. Furthermore, we
128 show that the observed image-by-image difference in DCNN and primate behavior along with
129 precisely measured IT population dynamics for each image better constrain the next generation
130 of ventral stream neural network models over previous qualitative approaches.

131

132

133

134 Results

135
136 As outlined above, we reasoned that, if recurrent circuits are critical to core object recognition
137 behavior, then current non-recurrent DCNNs should perform less accurately than the ventral
138 stream for some images. The first goal of this study was to discover many such *challenge*
139 images. Rather than making assumptions about what types of images (occluded, cluttered,
140 blurred, etc.) might most critically depend on feedback, we instead took a data driven approach
141 to identify such images.
142

143 Identification of DCNN *challenge* and *control* images

144
145 To compare the behavioral performance of primates (humans and macaques) and current
146 DCNNs image-by-image, we used a binary object discrimination task that we have previously
147 tested extensively (Figure 1C, ^{10, 11}). For each trial, monkeys used an eye movement to select
148 one of two object choices, after we briefly (100 ms) presented a test image containing one of
149 those choice objects (see Primate Behavioral Testing in Methods). Once monkeys are trained
150 in the basic task paradigm, they readily learn each new object over full viewing and background
151 transformations in just one or two days and they easily generalize to completely new images of
152 each learned object ¹⁰. This rapid learning suggests that this task taps into relatively natural
153 visual behavior, and that the object learning is unlikely to produce strong changes in the ventral
154 visual stream.
155

156 We tested a total of 1320 images (132 images of each of ten objects), in which the primary
157 visible object belonged to one of 10 different object categories (Figure 1B). To make the task
158 challenging, we included various image types (see Figure S1A): synthetic objects with high view
159 variation (scale, position and rotation) on cluttered natural backgrounds (similar to the ones
160 used in ^{6, 22}), and images with occlusion, deformation, missing object-parts, and colored
161 photographs (MS COCO dataset ²³).
162

163 Behavioral testing of all of these images was done in humans (n=88; Figure S2) and in monkeys
164 (n=2; Figure 1D). We estimated the behavioral performance of the subject pool on each image,
165 and that vector of image-wise d' performance values is referred to as I_1 (see Methods; also
166 refer ¹¹). We collected sufficient data such that the reliability of the I_1 vector was reasonably high
167 (median split half reliability $\tilde{\rho}$, humans = 0.84 and monkeys = 0.88, where 1.0 is perfect
168 reliability; see methods). To test the behavior of each DCNN model, we first extracted the
169 image evoked features of the penultimate simulated neural layer, e.g. fc7 layer of AlexNet ²¹.
170 We then trained ten linear decoders (see Methods) to derive the binary task performances, and
171 used a different set of images to test each model. Figure 1D shows an image-by-image
172 behavioral comparison between the pooled monkey population and AlexNet 'fc7'. We defined
173 and identified *control* images (blue dots; Figure 1C) as those where the absolute difference in
174 primate and DCNN performance does not exceed 0.4 (d' units), and we defined and identified

175 *challenge* images (red dots; Figure 1D) as those where the primate performance was at least
176 1.5 d' units greater than the DCNN performance. The behavioral performances for each image
177 (each object shown separately) are elaborated in the panels of Figure S3. Four examples of
178 *challenge* and *control* images are shown in Figure 1E. The *challenge* images were not
179 idiosyncratic to our choice of the AlexNet ('fc7') model. Many of them also turned out to be
180 *challenge* images for a range of other tested feedforward DCNNs with similar architectural
181 parameters, e.g., VGG-S^{24, 25}, Zeiler and Fergus²⁶; see Figure S1B. Challenge images were
182 also not specific to our synthetic image generation procedure. Figure S6A shows the challenge
183 and control image estimation from the MS COCO image-set.

184

185 Our results show that on average, both macaques and humans outperform AlexNet. Most
186 importantly, this image search procedure produced two groups of images: 1) 266 *challenge*
187 images that are accurately solved by primates but are not solved by a feedforward-only DCNN
188 (AlexNet; but see later), and 2) 149 *control* images that are solved equally well by primates and
189 the DCNN. On visual inspection, we did not observe any specific image property that
190 differentiated between these two groups of images. We also did not observe any difference in
191 performance on these two image-sets as the monkeys were repeatedly exposed to these
192 images (Figure S4). This is consistent with earlier work on monkey behavioral training¹⁰, that
193 showed — once the monkeys are trained with images of specific objects, their generalization
194 performance to new images from the same generative space is very high and consistent with
195 that of the training images. However, we observed that the reaction times of the subjects (both
196 humans and macaques) for *challenge* images were significantly higher than for the *control*
197 images (monkeys: $\Delta RT = 11.9$ ms; unpaired two-sample t-test, $t(413) = 3.4$; $p < 0.0001$;
198 humans: $\Delta RT = 25$ ms; unpaired two-sample t-test, $t(413) = 7.52$; $p < 0.0001$), suggesting that
199 additional processing time is required for the challenge images.

200

201 **Temporal evolution of image-by-image object representation in IT**

202 Previous studies^{4, 27} have shown that the identity of an object in an image is often accurately
203 conveyed in the population activity patterns of the inferior temporal cortex in the macaque.
204 Specifically, appropriately weighted linear combinations of the activities of these IT neurons can
205 approximate how neurons in downstream brain regions could integrate this information to form a
206 decision about the object identity. Such learned weighted linear combinations can accurately
207 predict the average behavioral performance in all tested core object recognition tasks⁶. That
208 previous work assumed one weighted linear combination of the neural population response
209 vector per object category (each is termed an “object decoder”) and we adopted that same
210 linear-decode assumption here as well.

211

212 In this study, we aimed to compare and contrast these linear object decodes from IT for the
213 *challenge* and *control* images. First, we wanted to know if these IT object decoders were as
214 accurate as the primates for both types of images — as predicted by the leading IT decoding
215 model⁶ — because that would demonstrate that the ventral stream successfully solves the
216 *challenge* images (images that are, by definition, not solved by current feedforward DCNNs, but

217 are somehow solved by primates). Second, we reasoned that, if recurrent computations were
218 crucial to these solutions, those computations would introduce additional processing time, and
219 therefore IT object decodes for *challenge* images should emerge later than IT object decode for
220 *control* images. Thus, we here used a sliding decoding time window (10 ms) that was narrower
221 than prior work ⁶ so that we could precisely probe the temporal dynamics of linearly-decodable
222 object category information.

223
224 To estimate the temporal evolution of the IT object decode for each image, we used large scale
225 multi-electrode array recordings (Figure 2A) to sample and record hundreds of neural sites
226 across IT cortex in two awake, behaving macaques. In each monkey, we implanted multiple
227 chronic 96-channel microelectrode arrays, inferior to the superior temporal sulcus (STS) and
228 anterior to the posterior middle temporal sulcus (pMTS); each array sampled from ~25 mm² of
229 the posterior, central and anterior part of IT. Recording sites that yielded a significant visual
230 drive (d'_{visual}), high selectivity and high image rank order response reliability (ρ_{site}^{IRO}) across
231 trials were considered for further analyses (see Neural recording quality metrics in Methods;
232 Figure S5 shows the average neural reliability across all neurons over time). In total, we
233 recorded from 424 valid IT sites which included 159 and 139 sites in the right hemisphere and
234 32 and 94 sites in the left hemisphere of monkey M (shown as inset in Figure 2A) and monkey
235 N respectively.

236
237
238 To determine the time at which explicit object identity representations are sufficiently formed in
239 the IT population activity, we plotted the temporal trajectory of the IT object decode accuracy for
240 each image. The IT decodes were obtained by training 10 linear (SVM) classifiers to predict the
241 respective object categories from the IT population vector using 10 ms non-overlapping time
242 bins. We computed the neural decoding accuracies (NDA) per time-bin by training and testing
243 independently at each time bin. Consistent with prior work ²⁷, this reveals that the the linearly
244 available information is not the same at each time — for example decoders trained at early time
245 bins (~100-130) do not generalize to late time bins with respect to decoding accuracies(Figure
246 S16). Thus, we determined the time at which the NDA measured for each image reached the
247 level of the subject's (pooled monkey) behavioral accuracy. We termed this time, the *Object*
248 *Solution Time* (OST), and we emphasize that each image has a potentially unique solution time
249 (OST_{image}). Briefly, OST for each image, was defined as the time (relative to image onset) when
250 the linear IT population decode (see Methods; Figure 2A, top panel) first rose to within the error
251 margins of the pooled monkey behavioral score for that image (see examples in Figure. 2B).
252 Because we recorded many repetitions of each image, we were able to measure OST_{image} very
253 accurately (standard error of ~9ms on average, as determined via bootstrapping across
254 repetitions). We also observed that the OSTs estimated by randomly subsampling half (n=212)
255 the total number of sites were significantly correlated (Spearman R was 0.77 and 0.76 for
256 control and challenge images respectively; $p < 0.00001$; and ΔOST was maintained ~30 ms) with
257 the OSTs from the total number of sites (n=424).

258
259

260 Figure 2B shows the temporal evolution of the IT object decode and the OST estimates for two
261 *control* images and two *challenge* images. For all four images, the correct (ground truth)
262 answer is the object 'bear' (insets in Figure 2B). Two observations are apparent in these
263 examples. First, for both the *control* and the *challenge* images, the IT decodes achieve the
264 behavioral accuracy of the monkey (note, behavioral accuracy is similar for all four images, by
265 design). Second, the IT decode solutions for *challenge* images emerge slightly later than the
266 solutions for the *control* images.

267
268 Both of these observations were also found on average in the full sets of *challenge* and *control*
269 images. First, the IT decodes achieved the primate behavioral level of accuracy on average for
270 the challenge and *control* image-sets (~91 % of *challenge* images and ~97 % of *control*
271 images), which meant that we could determine an OST for essentially all of these images.
272 Second, and consistent with our hypothesis, we observed that IT object solution times
273 (OST_{image}) for the *challenge* images were, on average, ~30 ms later compared to the *control*
274 images. Specifically, the median OST for the *challenge* images was 145 ± 1.4 ms (median \pm
275 SE) from stimulus onset and the median OST for the *control* images was 115 ± 1.4 ms (median
276 \pm SE) (Figure 2C). The average difference (~30 ms) between the OSTs of *challenge* and *control*
277 images did not depend on our choice of behavioral accuracy levels (Figure S7A) or image-set
278 type (Figure S6B). We also observed that there is a significant correlation between OSTs
279 estimated using a random half of the total number of sites (20 random splits) with that of the
280 entire dataset (Spearman R was 0.77 and 0.76 for control and challenge images respectively;
281 $p < 0.00001$; and ΔOST was maintained ~30 ms; $OST_{control} = 122 \pm 2.4$ ms, $OST_{challenge} = 151 \pm$
282 3.1 ms; estimated as median \pm SE of OST across control and challenge images, which were
283 estimated by averaging across 20 random split halves of the full neural population).

284
285
286

287 These results are consistent with the hypothesis that recurrent circuit computations are critical
288 to core object recognition (see Introduction). Thus, we next carried out a series of *controls* to
289 rule out alternative explanations for these results.

290 **Comparison of initial visual drive in IT evoked by *control* and *challenge* images**

291
292 We considered the possibility that the observed *OST* lag for the *challenge* images might have
293 been due to the IT neurons taking longer to start responding to these images. For example, if
294 the information in those images took longer to be transmitted by the retina. However, the data
295 do not support this possibility. First, we observed that *control* and *challenge* images share the
296 same population neural onset response latencies — the difference in IT response onset latency
297 was only 0.17 ms (median; ± 0.21 ms, SE; paired t-test; $t(423) = 0.3896$, $p = 0.69$; see Figure
298 3A, Figure S7B), suggesting that the initial visual drive for the images in both sets arrive at
299 approximately the same time in IT.

300
301

302
303 We considered the possibility that the difference in the OST between control and challenge
304 images for each object category is primarily driven by neurons that specifically prefer that
305 category (*object relevant neurons*). To address this, we first asked whether the object relevant
306 neurons show a significant difference in response latency (i.e. Δt_{onset} (challenge - control image)
307 > 0) when measured for their preferred object category. Our results (as demonstrated in Figure
308 S8 A-C) show that Δt_{onset} was not significant for any object category. In fact a closer inspection
309 (top panel of Figure S8C) reveals that for some objects (e.g. bear, elephant, dog) Δt_{onset} was
310 actually negative — that is, a trend for slightly *shorter* response latency for challenge images.
311 Finally, to test the possibility that there was an overall trend for the most selective neurons to
312 show a significant Δt_{onset} , we computed the correlation between the Δt_{onset} and the individual
313 object selectivity per neuron, per object category. We observed (bottom panel: Figure S8C) that
314 there was no dependence of object selectivity per neuron on the response latency differences.
315 In sum, the later mean OST for challenge images cannot be simply explained by longer
316 response latencies in the IT neurons that “care” about the object categories.

317
318 Interestingly however, we found that firing rates (R) were significantly higher ($\% \Delta R = 17.3\%$,
319 paired t-test; $t(423) = 6.8848$, $p < 0.0001$) for *challenge* images compared to *control* images,
320 tested on a 30 ms window centered at 150 ms post stimuli onset (see Figure 3A). We do not
321 yet know how to interpret this higher firing rate, but one possible explanation of this difference in
322 IT mean firing rate is the effect of additional inputs from activated recurrent circuits into the IT
323 neural sites at later time points (see Discussion). Regardless, these observations show that the
324 *challenge* images drive IT neurons just as quickly and at least as strongly as the *control* images.

325
326 When we closely examined the neural population response latencies for each image, we found
327 that the time at which the IT population firing rates started to increase from baseline (onset
328 latency; t_{onset}) and when the population firing rate reached its peak (t_{peak}) were on average earlier
329 than the OST for the images (Figure 3B and 3C). We also found no correlation (Pearson $r =$
330 0.009 ; $p = 0.8$) between the population response onset latency for each image (see Methods)
331 and the OST for that image (see Figure 3D). For example, inspection of Figure 3D reveals that
332 some of the *challenge* images evoke faster-than-average latency responses in IT, yet have slow
333 OSTs (~ 200 ms). Conversely, some of the *control* images evoke slower-than-average IT
334 responses, yet have relatively fast OSTs (~ 110 ms). In sum, these results show that visual drive
335 rapidly reaches IT for nearly all of these images, but that, for some images (mostly the
336 *challenge* images), that visually driven population activity takes longer to evolve to an accurate,
337 linearly-decodable format (OST).

338

339 **Controls for low level image properties**

340
341 We next considered the possibility that the average time lag for the *challenge* image OSTs
342 might have been due to low level image property differences between the two image-sets. From
343 previous research, we know that temporal properties of IT neurons depend critically on low level

344 image features like total image contrast energy²⁸, spatial frequency power distribution²⁹, and
345 spatial location of the visual objects³⁰. So we asked if these low level explanations might
346 explain the lag of the *challenge* image *OSTs*. First, we did not find any significant differences
347 ($\Delta t_{\text{onset}} = 0.17$ ms, paired t-test; $t(423) = 0.3896$; $p=0.697$) in neural firing rate onset latencies
348 (Figure 3A, Figure S4B) between *control* and *challenge* images across the recorded neural
349 sites. We also observed that solution times were not significantly correlated with image contrast
350 (Spearman $\rho=-0.04$; $p=0.47$). Second, we used the SHINE (spectrum, histogram, and intensity
351 normalization and equalization; Figure S7C) technique³¹ to equate low level image properties
352 across the *control* and *challenge* image-sets, and re-ran the recording experiment (subsampling
353 118 images each from the *control* and *challenge* image-sets; no. of repetitions per image = 44;
354 see Methods). The average estimated difference in *OST* values between “SHINED” *challenge*
355 and *control* images was still ~24 ms (Figure S7D). Third, we tested whether the overall
356 difference in *OST* between the challenge and control images, was specific to certain low or high
357 values of various image based properties (image clutter, blur, contrast, object size and object
358 eccentricity; for definition — see Methods). We observed that although certain image properties
359 were significantly correlated with the absolute *OST* values, the ΔOST was consistently ~30 ms at
360 different levels of these factors (Figure S18).

361
362 To test whether ΔOST (challenge - control) depends on neurons with higher or lower absolute
363 latencies, we divided the neural population into two groups — low latencies (<25 percentile of
364 the neural latencies; $n = 67$) and high latencies (>75 percentile of all neural latencies; $n = 67$).
365 We found that both neural groups conveyed similar information about the two types of images.
366 Specifically, we observed that there was no significant difference between control and challenge
367 image decoding accuracies estimated at the *OST* of each image, for both the low and high
368 latency populations (median $d'_{\text{high-latency}}^{\text{control}} = 1.23$, $d'_{\text{high-latency}}^{\text{challenge}} = 1.3$, $d'_{\text{low-latency}}^{\text{control}} = 1.05$,
369 $d'_{\text{low-latency}}^{\text{challenge}} = 1.04$; unpaired t-test for high latency group, $t(388)=0.17$, $p=0.86$; unpaired t-test
370 for low latency group, $t(388)=1.2$, $p = 0.2$). Consistent with our main result, we also found that
371 the low latency group of neurons and the the high latency group of neurons each showed a
372 positive lag for decoding of challenge images relative to control images ($\Delta DecodeLatency_{th=1.0}^{\text{low}}$
373 $= \sim 22$ ms, $\Delta DecodeLatency_{th=1.0}^{\text{high}} = \sim 18$ ms; note that we here set a decoding threshold of 1.0 to
374 compensate for the smaller number of neurons relative to the ~400 needed to achieve monkey
375 behavioral d').

376
377 To test whether the response latencies of an earlier area in the ventral stream hierarchy (area
378 V4) to the control and challenge images are different, we also simultaneously recorded from
379 area V4 in the left (95 sites) and right (56 sites) hemispheres of monkey M and N respectively
380 (see Methods). We found no significant difference in the response latencies (both onset and
381 peak) between control and challenge images across the V4 sites (Figure S9; paired t-test;
382 $t(150)=0.2$; $p=0.8$). These results further support the hypothesis that the ΔOST between the
383 challenge and the control images in IT is not driven by image properties that evoke shorter
384 latencies for control images at lower levels of the visual system.

385
386

387 Object solution estimates and timing during passive viewing

388
389

390 To test whether the late-emerging object solutions in IT only emerge when the animal is
391 performing the task (“active” condition), we also recorded IT population activity during “passive”
392 viewing of all the *challenge* and *control* images. Monkeys fixated a dot, while images were each
393 presented for 100 ms (same duration as the active task viewing of the image, see Figure. 1),
394 followed by 100 ms of no image, followed by the next image for 100 ms, etc. (typically 5 images
395 were presented per fixation trial; see Methods). *A priori*, several outcomes of switching from
396 active to passive viewing seemed likely: a decreased goodness of both the early-emerging and
397 the late-emerging IT decoded solutions, a decreased goodness of the late-emerging solutions, a
398 further delay in the late-emerging solutions, or no effect.

399

400 First, similar to the active condition (% $\Delta R = 17.3\%$), we observed that *challenge* images evoked
401 a significant higher firing rate (% $\Delta R = 13.2\%$, paired t-test; $t(423) = 8.27$, $p < 0.0001$) at later
402 time points (tested on a 30 ms window centered at 150 ms post stimuli onset) compared to the
403 *control* images (Figure S10A). Second, similar to the active viewing, we observed that we could
404 successfully estimate the object solution times for 92% of *challenge* and 98% of *control* images.
405 The object solution times estimated during the active and passive conditions were also strongly
406 correlated (Spearman $\rho = 0.76$; $p < 0.0001$). Similar to the active condition, *challenge* image
407 solutions required an additional time of ~ 28 ms (on average) to achieve full solution compared
408 to the *control* images (Figure S10B). In sum, we observe that the solutions in IT emerge with a
409 similar lag and overall accuracy (goodness) during passive viewing. This suggest that the
410 putative recurrent computations that underlie the late-emerging IT object information are not
411 task dependent, but are instead reflexive and automatically triggered by the images. This is
412 consistent with previous findings of McKee et al.³², where they reported that macaque IT cortex
413 predominantly shows task-independent visual feature representation. Similar results have also
414 been reported in humans³³.

415

416 However, because these animals were trained on the object discrimination task prior to any
417 neural data collection, it might be that the OST difference is due to internal processes that are
418 only activated in trained monkeys (e.g. mental task performance?) or somehow due to the
419 training history itself. To test this, we carried out the same analyses on smaller sets of data
420 from two untrained animals. Specifically, given our behavioral work in the current study to sort
421 images into two types (challenge and control images; Figure. 1D), we were able to sort control
422 and challenge images out of a pool of images that we had previously collected in two untrained
423 monkeys during a passive fixation task (previously reported in^{6,35,7,8}). To appropriately compare
424 with the results from the trained monkeys in this study, we matched the set of common images
425 (640 images), array implant locations (left hemispheres; posterior and central IT), number of
426 neural sites (168), and number of image repetitions (43). We observed a small, but significant
427 overall decrease in decoding accuracy across all images in the untrained monkey IT decodes
428 (paired t-test; median $\Delta d' = 0.23$, $t(639) = 7.78$; $p < 0.0001$). Most importantly however, we found
429 that the IT cortex of untrained monkeys demonstrated lagged decode solutions for the challenge
430 images (relative to the control images) that are very comparable to those observed from the

431 trained monkey IT populations (estimated at a decoding accuracy threshold of 1.8;
432 $\Delta DecodeLatency_{th=1.8}^{untrained} = \sim 34$ ms; $\Delta DecodeLatency_{th=1.8}^{trained} = \sim 30$ ms; see Figure S11). In
433 sum, the main experimental observation we report here (lagged OST for challenge images)
434 appears to be largely automatic, and it does not require, and is not the results of, laboratory
435 training.

436

437 **IT predictivity across time using current feedforward deep neural network models** 438 **of the ventral stream**

439

440 We reasoned that, if the late-emerging IT population solutions are indeed dependent on
441 recurrent computations that are lacking in current DCNN models, then perhaps the previously
442 demonstrated ability of those models to (partially) explain and predict individual IT neurons⁸
443 was due mostly to the similarity of the DCNN population response to the feedforward portion of
444 the IT population response. To test this idea, we asked how well the DCNN “IT” population
445 response pattern (which is not temporally evolving) could predict the time-evolving IT neural
446 population response pattern up to and including the *OST* of each image. To do this, we used
447 previously described methods (similar to⁸). Specifically, we quantified the IT population
448 goodness of fit as the median (over neurons) of the noise corrected explained response
449 variance score (IT predictivity; Figure S12; also see Methods).

450

451 First, we observed that the top layers (penultimate) of the DCNN (AlexNet ‘fc7’) predicted $44.3 \pm$
452 0.7% of the potentially explainable (i.e. image-driven) IT neural response variance during the
453 early response phase (90-110 ms; Figure 4A) for all images. This result further confirms that
454 feedforward DCNNs indeed approximate the initial (putative largely feedforward) IT population
455 response pattern. However, we observed that the ability of this DCNN’s “IT” population
456 response to predict the IT population pattern significantly worsened (<20% explained variance)
457 as that response pattern evolved over time (Figure 4A). This drop in IT predictivity was not due
458 to low signal to noise ratio of the neural responses during those time points because our
459 explained variance measure already compensates for any changes in SNR, and also because
460 SNR remains relatively high in the late part of the IT responses (Figure S5). In sum, the gradual
461 drop in IT predictivity by these feedforward DCNN models is consistent with the hypothesis that
462 late-phase IT population responses are modified by the action of recurrent circuits that are not
463 contained in those DCNN models. Consistent with our hypothesis that *challenge* images rely
464 more strongly on those recurrent circuits than *control* images, we observed that the drop in IT
465 predictivity coincided with the solution times of the *challenge* images (refer top panel histograms
466 for OST distributions of *challenge* and *control* images).

467 **Evaluation of deeper CNNs as models of ventral visual stream processing**

468 Although, the above results suggest the likely importance of recurrent computations in the
469 primate ventral stream for some images, we are still left with the open question: what specific
470 computational function do recurrent circuits provide beyond the feedforward representation

471 during core object recognition behavior? It is understood in the artificial neural network
472 community that finite-time recurrent neural networks can be constructed as very deep,
473 feedforward-only neural networks with weight sharing across layers that are recurrently
474 connected in the original recurrent network³⁴. Thus, we reasoned that the actions of recurrent
475 circuits in the ventral stream might be computationally equivalent to stacking further non-linear
476 transformations onto the initially evoked (~feedforward) IT population response pattern. In
477 particular, perhaps neural populations from newer DCNNs for visual categorization that have an
478 even higher number of stacked nonlinear transformations might better approximate the
479 recurrent computations of the ventral stream, even though they were not specifically designed to
480 emulate the anatomical recurrent circuits of the ventral stream. To test this idea, we asked if
481 existing very deep CNNs provide a better neural match to the IT response at its late phase and
482 to the image-by-image patterns of behavioral performance. Currently there are many deeper
483 CNNs available that outperform the baseline DCNN (AlexNet) used here, such as inception-v3
484 ³⁵, inception-v4 ³⁶ and ResNet-50, ResNet-101³⁷. Based on the number of layers (non-linear
485 transformations), we divided the tested DCNN models into two groups, deep (8 layers; AlexNet,
486 Zeiler and Fergus model, VGG-S) and deeper (>20 layers, inception-v3, inception-v4, ResNet-
487 50, ResNet-101) CNNs. We made three observations, that corroborate our speculation.

488
489
490 Given that decodes out of IT neural populations, typically have the highest behavioral
491 consistency to that of primates⁶, compared to any other area in the ventral stream, we first
492 searched all the above mentioned neural networks to determine which layer of the models has
493 the highest I_1 behavioral consistency on our image-set. We referred to this layer as the model-
494 IT layer. Interestingly, we observed that the model-IT layers of very deep CNNs predict IT neural
495 responses at the late phases (150-250 ms) significantly higher (Δ Predictivity = 5.8%, paired t-
496 test; $t(423) = 14.26$, $p < 0.0001$) than “regular-deep” models like AlexNet (Figure 4B; scatter plot
497 comparisons with AlexNet shown separately in Figure S13). This observation suggests that very
498 deep CNNs might indeed be approximating “unrolled” versions of the ventral stream’s recurrent
499 circuits. Second, as expected from the ImageNet challenge results ³⁸, we observed an
500 increased performance and therefore reduced number of *challenge* images for very deep CNNs.
501 Third, we found that the images that remain unsolved by these very deep CNNs (i.e. *challenge*
502 images for these models) showed even longer *OST*s in IT cortex than the original full set of
503 *challenge* images (Figure 4C). Assuming that longer *OST* is a signature of more recurrent
504 computations, this suggests that the newer, very deep CNNs have implicitly, but only partially,
505 approximated — in a feedforward network — some of the computations that the ventral stream
506 implements recurrently to solve some of the *challenge* images.

508 **Evaluation of CORnet (a regular-deep-recurrent CNN) as a model of the ventral** 509 **visual stream**

510
511 To more directly ask if the experimental observations above might indeed be the result of
512 recurrent computations, we implemented an ANN model that does recurrent computations.

513 Specifically we tested a regular-deep (i.e. less than 10 layers) recurrent neural network model,
514 termed CORnet³⁹ We chose this particular network given its very high performance on brain-
515 score (<http://brain-score.org/>⁴⁰), an online platform that hosts the neural and behavioral
516 benchmarks for core object recognition models. The IT-layer of CORnet has within-area
517 recurrent connections (with shared weights). The model currently implements five time-steps
518 (pass1- pass5 in Figure 4B). Therefore, the activity arising at the first time-step in the model-IT
519 layer is nonlinearly transformed to arrive at the output of the second time step and so on.
520 Indeed, we observed that CORnet had higher predictivity (Figure 4C) for the late-phase IT
521 responses (for images that had late OSTs; >145 ms). In addition, pass-1 and pass-2
522 (corresponding to time-step 1) of the network had a significant (multiple-comparison corrected-
523 paired t-test; $t(423)=12.78$; $p<0.00001$) lower IT predictivity compared to pass-3 and 4 for later
524 time-steps, whereas the opposite was true for earlier time-steps (Figure S14). Taken together,
525 these results further argue for recurrent computations in the ventral stream.

526 **Comparison of backward visual masking between *challenge* and *control* images**

527 So far we have observed that feedforward DCNNs poorly predict the IT neural responses at
528 later times beyond the putative feedforward response (90-110 ms post image onset), during
529 which a majority of the *challenge* images (~82%) evoke their object solutions in IT. Based on
530 these results, we hypothesized that these later IT population responses are critical for
531 successful core object recognition behavior for many of the *challenge* images (~57% of
532 *challenge* images have OST>140 ms). To further test this idea, we performed an additional
533 behavioral experiment that aimed to corroborate the neurophysiology results. We modified the
534 original object discrimination paradigm by adding a visual mask (phase scrambled image,⁴¹) for
535 500 ms (Figure 5A), immediately following the test image presentation: a manipulation
536 commonly known as backward visual masking. Such backward masking has been previously
537 associated with selective disruption of the recurrent inputs to an area from other areas^{42, 43},
538 limiting the visual processing to the initial feedforward response⁴⁴. Given that solutions for the
539 *challenge* images can arise in IT cortex only at later time points compared to the *control* images,
540 we reasoned that if disruption in processing produced by a visual mask affects IT at earlier
541 times, it will produce larger behavioral deficits for *challenge* images compared to *control* images.
542 However, we predicted that these differences should subside at longer presentation times when
543 enough time is provided for the recurrent processes to build a sufficient object representation for
544 both *control* and *challenge* images in IT. Therefore, during this experiment, we tested a range of
545 masking disruption times by randomly interleaving the sample image duration (and thus the
546 mask onset). Specifically, we tested 34, 67, 100, 167 and 267 ms (see Methods). Our results
547 (Figure 5B) show that visual masking indeed had a significantly stronger effect on the *challenge*
548 images at smaller presentation durations compared to the *control* images. Consistent with our
549 hypothesis, we did not observe any measurable masking differences between the two image-
550 sets at longer presentation times (~267 ms). Median $\Delta d'$ (difference between *control* and
551 *challenge* images grouped by objects) averaged across all 10 objects were 0.5, 0.81, 0.33, 0.40,
552 and -0.02 for 34, 67, 100, 167 and 267 ms presentation duration respectively. The difference in
553 performance was statistically significant at the .05 significance level (Bonferroni adjusted) for
554 all presentation durations except 267 ms. Together with the neurophysiology results, these

555 observations provide converging evidence that rapid, automatic, recurrent ventral stream
556 computations are critical to the brain's ability to infer object identity in the *challenge* images,
557 even at the rapid time scale of natural vision (~200 ms per fixation).

558 **Model-driven versus image-property driven approaches to study recurrence**

559 Previous research has suggested that recurrent computations in the ventral stream might be
560 necessary to achieve pattern completion when exposed to occluded images⁴⁵⁻⁴⁷, object based
561 attention in cluttered scenes^{45, 48}, etc. Indeed, we observe that several image properties like
562 object size, presence of occlusion, and object eccentricity, as well as a combination of all these
563 factors (Figure 6) are significant, but very weak predictors of our putative recurrence signal (the
564 OST vector; see Methods: Estimation of the OST prediction strength). In comparison, the
565 performance gap between AlexNet and the monkey behavior ($\Delta d'$) is a significantly stronger
566 predictor of OST. Therefore, our results suggest another possible image-wise predictor of
567 ventral stream recurrence — the difference in performance between feed-forward DCNNs and
568 primates, $\Delta d'$. This vector is likely itself dependent on a complex combination of image
569 properties, such as those mentioned above. However, it is directly computable and our results
570 show that it can serve as a much better model guide. In particular, we find that $\Delta d'$ is
571 significantly predictive of the OST for each image (Spearman $\rho = 0.44$; $p < 0.001$), and, in this
572 sense, is a much better predictor of the engagement of ventral stream recurrence than any of
573 the individual image properties.

574

575 Discussion

576

577 The overall goal of this study was to ask if recurrent circuits are critical to the ventral stream's
578 execution of core recognition behavior — the ability to report object category in the central 10 °
579 with less than 200 ms of image viewing duration. We reasoned that, if computations mediated
580 by recurrent circuits are critical for some images, then one way to find such images is by finding
581 images that are difficult for non-recurrent DCNNs to solve, but are nevertheless easily solved by
582 primates. Thus we first used extensive behavioral testing to find such *challenge* images along
583 with behaviorally matched *control* images. With these in hand, we then aimed to look for a likely
584 empirical signature of recurrence — the requirement of additional time to complete successful
585 processing. To ask this question, we first had to confirm that the *challenge* images that are
586 behaviorally solved (by definition) were, in fact, solved by the ventral stream — as predicted by
587 current models of the neural mechanisms underlying core recognition ⁶. Using large-scale IT
588 population neurophysiology, we confirmed part of this prediction: behaviorally-sufficient linearly
589 decodable object solutions emerged in the IT population activity for essentially all of the
590 *challenge* images (assessed with the same number of neurons and training exemplars as for the
591 *control* images). But looking at the temporal evolution of these IT population solutions
592 simultaneously revealed a key observation not revealed in prior work ⁶ — the IT solutions were
593 lagged by an average of ~30ms later for *challenge* images compared to the *control* images. In
594 addition, we also found that the temporally lagged IT population response patterns that
595 contained the linearly-decodable object identity solutions were poorly predicted by DCNN model
596 “neural” population responses to the same *challenge* images. This stands in contrast to the
597 early IT population responses, which were much better predicted by the DCNN model,
598 consistent with prior work ⁸. Notably, we observed both of these findings during active task
599 performance (when the animals had to report the identity of the dominant object in the image),
600 but we found all of these results to be almost identical during passive viewing. Taken together,
601 these results imply that automatically-evoked recurrent circuits are critical for object
602 identification behavior even at the fast timescales of core object recognition.

603

604 The idea that “feedback”, broadly construed, is important to vision and to object recognition is
605 not new (see ⁴⁹ for review). Previous reports ⁵⁰ demonstrated that different forms of information
606 can be decoded from early and late responses in IT, suggesting a potential role of intra-areal
607 recurrent inputs during the late-phase IT responses. While such broad concepts about the
608 potential role of feedback in vision have been previously suggested and partly explored, we
609 believe that this is the first work to examine these questions at such large scale, at the fast time
610 scales of core object recognition; the first to do so using image computable models of the neural
611 processing to guide the choice of experiments (i.e. the images and discrimination tasks), and
612 the first to do so with an implemented linking model (decoder) of how IT supports recognition
613 behavior.

614 **Late object identity solution times in IT imply recurrent computations underlie**
615 **core recognition**

616 The most parsimonious interpretation of the results reported here is that the late phases of the
617 stimulus evoked responses in IT depend on some type (or types) of recurrent computations that
618 are not present in today's non-recurrent DCNN ventral stream models. And our comparisons
619 with behavior suggest that these IT dynamics are not epiphenomenal, but are critical to core
620 object recognition behavior. But what kind(s) of additional computations are taking place and
621 where in the brain do those recurrent circuit elements live? We do not yet know the answers to
622 these questions, but we can speculate to generate a testable set of hypotheses. Based on the
623 number of synapses between V1 and IT, Tovee⁵¹ proposed that the ventral stream comprises
624 of stages that are approximately 10-15 ms away from each other. Our observation of an
625 additional processing time of ~30 ms for *challenge* images is therefore equivalent to at least two
626 additional processing stages. Thus, one possible hypothesis is a cortico-cortical recurrent
627 pathway between the ventral stream cortical areas including IT and lower areas like V4, V2 and
628 V1 (similar to suggestions of⁵²⁻⁵⁴). This possibility is consistent with observations of temporally-
629 specific effects in the response dynamics of V4 neurons⁵⁵ for images with occlusion.
630 Alternatively, the temporal lag signature we report here is also consistent with the possibility that
631 IT is receiving important recurrent flow from downstream areas like the prefrontal and perirhinal
632 cortices (e.g. as suggested by^{56, 57}). We also cannot rule out the possibility that all of the
633 additional computations are due to recurrence within IT itself (e.g. consistent with recent models
634 such as⁴⁷), or due to subcortical circuits (e.g. basal ganglia loops,⁵⁸). These hypotheses are
635 not mutually exclusive. Given all that prior work, the main contribution of our work is to take the
636 very broad notion of "feedback" and pin down a narrower case that is both experimentally
637 tractable (i.e. the neural phenomena is observable in IT for a prescribed set of images) and is
638 guaranteed to have high behavioral relevance. The present results now motivate the need for
639 direct perturbation studies that aim to independently suppress each of those circuit motifs to
640 assess the relative importance of each of these circuit motifs. Such perturbations should be
641 paired with IT electrophysiological recordings and behavior. The results of the present study
642 also provide sets of images and predictions of exactly how and when IT will be disrupted when
643 the critical circuit motif(s) is/are suppressed. Specifically, our measurements of both the $\Delta d'$
644 and the OST_{image} vectors provide observable signatures of recurrent computations that make
645 clear predictions for such direct neural suppression studies. Based on our results here, we
646 predict that a specific disruption of the relevant recurrent circuits will prevent the emergence of
647 the object solutions to the *challenge* images in IT. This will in turn result in larger behavioral
648 deficits in the *challenge* images (relative to the *control* images). Note however, that the results
649 reported here provide more specific predictions for future perturbation experiments — beyond
650 control and challenge image differences. The estimated OST vector (putative "recurrence"
651 signal) predicts exactly which individual images will be most affected (i.e. the images showing
652 longer solution times). This knowledge can be used to optimize the image-sets and behavioral
653 tasks for these next experiments.
654

655 **Temporally specific failures of current ventral stream encoding models imply that**
656 **recurrent circuits are needed to improve those models**

657
658 Prior to this study, the best models of the ventral visual stream belonged to a class of
659 feedforward DCNNs, e.g. HMO⁸, AlexNet²¹ and VGG^{25, 59}. These studies^{7, 8} have
660 demonstrated that feedforward DCNNs can explain ~50% of the within-animal explainable
661 response variance in stimulus evoked V4 and IT responses (averaged responses from 70 - 170
662 ms post-stimulus onset). Our results here confirm that feedforward DCNNs indeed approximate
663 ~50% of the first 30 ms (~90-120 ms) of the stimulus evoked, within-animal explainable IT
664 response variance, thus establishing DCNNs as a good functional approximation of the
665 feedforward pass of the primate ventral stream. However, in addition, we observed that the
666 ability of DCNN neural populations to predict IT neural responses drops significantly at later
667 phases of the stimulus evoked IT responses (>150 ms after image onset, see Figure 4A). This
668 is consistent with our inference that the late object solution times for *challenge* images are
669 primarily caused by the additional processing time required by recurrent processes in the ventral
670 stream. Recruitment of recurrent circuits in the form of both intra and inter-cortical feedback
671 during these times might explain why the feedforward-only DCNN activations poorly predict the
672 late IT responses. In addition, other forms of dynamics coding, for instance, short-lived
673 spatiotemporal patterns of spiking⁶⁰ might also be relevant, and currently are missing from
674 DCNNs.

675 **Unique object solution times per image motivate the search for better models of**
676 **the link between IT neural population patterns and core object recognition**
677 **behavior**

678
679 Majaj et al.⁶ experimentally rejected a large number of alternative models that link ventral
680 stream population activity to core object recognition behavior (“decoding models”). The authors
681 showed that a simple linear decoding model, formed by linearly weighting the population activity
682 of IT neurons (integrated from 70-170 ms post image onset) was sufficient to explain and predict
683 the average performance of human subjects in each of a set of 64 tested core object recognition
684 tasks. However, in the Majaj et al.⁶ study, the key predictor variable (behavioral performance)
685 was computed as an average over all test images for any given task. The authors (one of us
686 among them) speculated that a much finer-grain predictor variable, e.g. image-level behavioral
687 performance, could provide a stronger test of these decoding models. Here we observe that,
688 even for images that have statistically non-distinguishable levels of behavioral performance, the
689 linearly-decodable information in the IT population pattern varies quite substantially over the IT
690 response time window used by the decoding models proposed by Majaj et al. (specifically — 70-
691 170 ms post stimulus onset). Taken together, this argues that future work in this direction might
692 successfully reject most or even all of the LaWS of RAD IT decoding models, and thus drive the
693 field to create better mechanistic neuronal-to-behavioral linking hypotheses.

694
695

696 **Role of recurrent computations: deliverables from these data and insights from**
697 **deeper CNNs**

698
699 Prior studies have strongly associated the role of recurrent computations during visual object
700 recognition tasks with overcoming certain specific challenging image properties that might be
701 boiled down to a single word or phrase such as “occlusion”⁴³, high levels of “clutter”⁴⁵,
702 “grouping” of behaviorally relevant image regions⁶¹ or the need for visual “pattern completion”⁴⁷.
703^{61, 62}. While we agree that such ideas or task conditions might recruit recurrent processes in the
704 ventral stream, the present work argues that picking any one of these single ideas is not the
705 most efficient approach to constrain future models of the mechanisms of object recognition.
706 Specifically, we have here found that a very good way to expose which images rely most heavily
707 on recurrent computations in the ventral stream is model-based. That is, we use the shallower
708 models to find images for which the difference between feedforward-only DCNN and primate
709 behavior ($\Delta d'$) is the largest, and this difference is a better predictor of the neural phenomena of
710 recurrence than any of the image-based properties (see Figure. 6). We interpret this to mean
711 that the models effectively embed knowledge about multiple interacting image properties that
712 cannot be described by single words or phrases, but that this knowledge better accounts for the
713 what happens in the feedforward part of the response than those other types of explanations.

714
715 While this is a good way to focus experimental efforts, it does not yet expose the computational
716 role of recurrence, i.e., the exact nature of the computational problem solved by recurrent
717 circuits during core object recognition. Interestingly, we found that deeper CNNs like inception-
718 v3, v4³⁶, ResNet-50,101³⁷, that introduce more nonlinear transformations to the image pixels,
719 compared to shallower networks like AlexNet or VGG, are better models of the late phase of IT
720 responses (the phase that is most behaviorally relevant for DCNN-*challenge* images). This is
721 also consistent with a previous study³⁴ where it was shown that a shallow recurrent neural
722 network (RNN) is equivalent to a very deep CNN (e.g. ResNet) with weight sharing among the
723 layers. Therefore, we speculate that what the computer vision community has achieved by
724 stacking more layers into the CNNs, is a partial approximation of something that is more
725 efficiently built into the primate brain architecture in the form of recurrent circuits. That is, during
726 core (~200 ms) object recognition, recurrent computations act as additional non-linear
727 transformations of the initial feedforward IT response, to produce more explicit (linearly
728 separable) solutions. This provides a qualitative explanation for what recurrent circuits provide
729 in a variety of challenging image conditions, the purpose of which is to achieve a more explicit
730 object representation at the level of IT. What is now needed are new recurrent artificial neural
731 networks (here we provided results from one such model: CORnet³⁹) that successfully
732 incorporate these ideas. While the data presented here cannot fully specify the form of those
733 ANNs, they will provide a strong check (see below) on any model that aims to succeed in these
734 more advanced vision challenges where primates still exceed machines, as well as behavioral
735 tasks that deal with more dynamic visual input (i.e. movies) and associated tasks such as action
736 recognition, etc.

737

738 **Constraints for future models provided by our data**

739

740

741 Our results motivate a change in the architecture of artificial neural networks that aim to model
742 the ventral visual stream (i.e. addition of recurrent circuits) — motivating a switch from largely
743 feedforward DCNNs to recurrent DCNNs. However, a primary goal of experiments is not simply
744 to provide motivation, but to also provide validation and strong constraints for guiding the
745 construction of those new models. The results obtained here provide three precisely measured
746 constraints for next generation neural network models. First, we provide a behavioral vector, $\Delta d'$
747 that quantifies the performance gap between current feedforward DCNNs (e.g. AlexNet) and the
748 image-by-image primate core object recognition behavior (I_1). Second, for each of these
749 images, we have estimated the time at which object solutions are sufficiently represented in the
750 macaque IT cortex (i.e. the OST_{image} vector). Third, we have reliably measured the neural
751 responses to each of the tested images at their respective object solution times (potential target
752 features for models). Next generation dynamic models of the ventral stream should be
753 constrained to produce the target features (object solutions) at these times. We will also host
754 the images, primate behavioral scores, estimated object solution times, and the modeling
755 results at <http://brain-score.org>⁴⁰.

756

757

758

759 Methods

760

761 Subjects

762 The nonhuman subjects in our experiments were two adult male rhesus monkeys (*Macaca*
763 *mulatta*). All human studies were done in accordance with the Massachusetts Institute of
764 Technology Committee on the Use of Humans as Experimental Subjects. A total of 88
765 observers participated in the binary object discrimination task. Observers completed these 20-
766 25 min tasks through Amazon's Mechanical Turk, an online platform in which subjects can
767 complete experiments for a small payment.

768 Visual stimuli: generation

769 Generation of synthetic (“naturalistic”) images

770

771 High-quality images of single objects were generated using free ray-tracing software
772 (<http://www.povray.org>), similar to Majaj et al. ⁶. Each image consisted of a 2D projection of a
773 3D model (purchased from Dosch Design and TurboSquid) added to a random background. The
774 ten objects chosen were **bear, elephant, face, apple, car, dog, chair, plane, bird** and **zebra**
775 (Figure 1B). By varying six viewing parameters, we explored three types of identity while
776 preserving object variation, position (x and y), rotation (x , y , and z), and size. All images were
777 achromatic with a native resolution of 256×256 pixels (see Figure S1A for example images). A
778 total of 1120 naturalistic images (112 per object category) were used.

779 Generation of natural images (photographs)

780

781 Images pertaining to the 10 nouns, were download from <http://cocodataset.org>. Each image was
782 resized to $256 \times 256 \times 3$ pixel size and presented within the central 8° . We used the same
783 images while testing the feedforward DCNNs. A total of 200 COCO images (20 per object
784 category) was used.

785

786 Quantification of image properties

787

788 We have compared the ability of different image properties to predict the putative recurrence
789 signal, inferred from our results. These image properties were either pre-defined during the
790 image generation process (e.g. object size, object eccentricity, and the object rotation vectors,
791 presence of an object occluder) or computed after the image generation procedure. The post
792 image generation properties are listed below:

793

794 *Image contrast.* This was defined as the variance of the luminance distribution per image
795 (grayscale images only).

796

797 *Image blur.* The image processing literature contains multiple measures of image focus
798 based on first order differentiation or smoothing followed by differentiation. We have used
799 a technique from Santos et al. ⁶³ to define the focus of an image.

800

801 *Image clutter.* This measure (Feature Congestion) of visual clutter is related to the local
802 variability in certain key features, e.g., color, contrast, and orientation ⁶⁴.

803 **Primate behavioral testing**

804 Humans tested on amazon mechanical turk

805

806 We measured human behavior (88 subjects) using the online Amazon MTurk platform which
807 enables efficient collection of large-scale psychophysical data from crowd-sourced “human
808 intelligence tasks” (HITs). The reliability of the online MTurk platform has been validated by
809 comparing results obtained from online and in-lab psychophysical experiments ^{6, 10}. Each trial
810 started with a 100 ms presentation of the sample image (one out of 1360 images). This was
811 followed by a blank gray screen for 100 ms; followed by a choice screen with the target and
812 distractor objects (similar to ¹¹). The subjects indicated their choice by touching the screen or
813 clicking the mouse over the target object. Each subject saw an image only once. We collected
814 the data such that, there were 80 unique subject responses per image, with varied distractor
815 objects.

816

817 Monkeys tested during simultaneous electrophysiology

818

819 Active binary object discrimination task

820

821 We measured monkey behavior from two male rhesus macaques. Images were presented on a
822 24-inch LCD monitor (1920 × 1080 at 60 Hz) positioned 42.5 cm in front of the animal. Monkeys
823 were head fixed. Monkeys fixated a white square dot (0.2°) for 300 ms to initiate a trial. The trial
824 started with the presentation of a sample image (from a set of 1360 images) for 100 ms. This
825 was followed by a blank gray screen for 100 ms, after which the choice screen was shown
826 containing a standard image of the target object (the correct choice) and a standard image of
827 the distractor object. The monkey was allowed to view freely the choice objects for up to 1500
828 ms and indicated its final choice by holding fixation over the selected object for 400 ms. Trials
829 were aborted if gaze was not held within $\pm 2^\circ$ of the central fixation dot during any point until the

830 choice screen was shown. Prior to the final behavioral testing, both monkeys were trained in
831 their home-cages on a touchscreen (for details see¹¹; details of the code and hardware available
832 at <https://github.com/dicarlolab/mkturk>) to perform the binary object discrimination tasks. We
833 used a separate set of images that were synthesized using the same image generation protocol
834 to train the monkeys on the binary object discrimination task. Once the behavioral performance
835 stabilized during the training, we then tested the monkeys on the image-set described in the
836 manuscript along with simultaneous electrophysiology.
837

838 Passive Viewing

839
840 During the passive viewing task, monkeys fixated a white square dot (0.2°) for 300 ms to initiate
841 a trial. We then presented a sequence of 5 to 10 images, each ON for 100 ms followed by a 100
842 ms gray (background) blank screen. This was followed by fluid reward and an inter trial interval
843 of 500 ms, followed by the next sequence. Trials were aborted if gaze was not held within $\pm 2^\circ$ of
844 the central fixation dot during any point.

845 Behavioral Metrics

846 We have used the same one-vs-all image level behavioral performance metric (I_1) to quantify
847 the performance of the humans, monkeys, deep HCNNs and neural based decoding models for
848 the binary match sample tasks. This metric estimates the overall discriminability of each image
849 containing a specific target object from all other objects (pooling across all 9 possible distractor
850 choices).

851 For example, given an image of object 'i', and all nine distractor objects ($j \neq i$) we first compute
852 the average hit rate,

$$853 \quad HitRate_{image}^i = \frac{\sum_{j=1}^{10} Pc_{image}^{i,j \neq i}}{9}, \text{ where } Pc \text{ refers to the fraction of correct responses}$$

854 for the binary task between objects 'i' and 'j'. We then compute the false alarm rate for the
855 object 'i' as

$$FalseAlarm^i = 1 - avg(HitRate_{image}^{j \neq i})$$

856 The unbiased behavioral performance, per image, was then computed using a sensitivity index
857 d' ,

$$858 \quad d'_{image} = z(HitRate_{image}^i) - z(FalseAlarm^i),$$

859 where z is the inverse of the cumulative Gaussian distribution. The values of d' were bounded
860 between -5 and 5. Given the size of our image-set, the I_1 vector contains 1320 independent d'
861 values. The estimated median false alarm rate across objects were 0.11 and 0.18 for the
862 monkey behavior and neural decoding performance respectively.
863

864 To compute the reliability of the estimated I_1 vector, we split the trials per image into two equal
865 halves by resampling without substitution. The Spearman-Brown corrected correlation of the two
866 corresponding I_1 vectors (one from each split half) was used as the reliability score (i.e. internal
867 consistency) of our I_1 estimation.

868 **Large scale multielectrode recordings and simultaneous behavioral recording**

869 Surgical implant of chronic micro-electrode arrays

870 Before training, we surgically implanted each monkey with a head post under aseptic conditions.
871 After behavioral training, we recorded neural activity using 10 × 10 micro-electrode arrays (Utah
872 arrays; Blackrock Microsystems). A total of 96 electrodes were connected per array. Each
873 electrode was 1.5 mm long and the distance between adjacent electrodes was 400 μm. Before
874 recording, we implanted each monkey multiple Utah arrays in the IT and V4 cortex. In monkey
875 M, we implanted 3 arrays in right hemisphere (all 3 in IT) and 3 arrays in the left hemisphere (2
876 in IT and 1 in V4). In monkey N, we implanted 3 arrays in the left hemisphere (all 3 in IT) and 3
877 arrays in the right hemisphere (2 in IT and 1 in V4). The left and right hemisphere arrays were
878 not implanted simultaneously. We recorded for ~6-8 months from implants in one hemisphere
879 before explanting the arrays and implanting new arrays in the opposite hemisphere. Array
880 placements were guided by the sulcus pattern, which was visible during surgery. The electrodes
881 were accessed through a percutaneous connector that allowed simultaneous recording from all
882 96 electrodes from each array. Behavioral testing was performed using standard operant
883 conditioning (fluid reward), head stabilization, and real-time video eye tracking. All surgical and
884 animal procedures were performed in accordance with National Institutes of Health guidelines
885 and the Massachusetts Institute of Technology Committee on Animal Care.
886

887 Eye Tracking

888
889 We monitored eye movements using video eye tracking (SR Research EyeLink 1000). Using
890 operant conditioning and water reward, our 2 subjects were trained to fixate a central white
891 square (0.2°) within a square fixation window that ranged from ±2°. At the start of each
892 behavioral session, monkeys performed an eye-tracking calibration task by making a saccade to
893 a range of spatial targets and maintaining fixation for 500 ms. Calibration was repeated if drift
894 was noticed over the course of the session.

895 Electrophysiological Recording

896 During each recording session, band-pass filtered (0.1 Hz to 10 kHz) neural activity was
897 recorded continuously at a sampling rate of 20 kHz using Intan Recording Controller (Intan
898 Technologies, LLC). The majority of the data presented here were based on multiunit activity.
899 We detected the multiunit spikes after the raw data was collected. A multiunit spike event was
900 defined as the threshold crossing when voltage (falling edge) deviated by more than three times
901 the standard deviation of the raw voltage values. Of 960 implanted electrodes, five arrays
902 (combined across the two hemispheres) × 96 electrodes × two monkeys, we focused on the 424
903 most visually driven, selective and reliable neural sites. Our array placements allowed us to
904 sample neural sites from different parts of IT, along the posterior to anterior axis. However, for
905 all the analyses, we did not consider the specific spatial location of the site, and treated each
906 site as a random sample from a pooled IT population.
907

908 Neural recording quality metrics per site

909

910 *Visual drive per neuron (d'_{visual}):* We estimated the overall visual drive for each electrode. This
 911 metric was estimated by comparing the COCO image responses of each site to a blank (gray
 912 screen) response.
 913

$$d'_{visual} = \frac{avg(R_{coco}) - avg(R_{gray})}{\sqrt{\frac{1}{2}(\sigma_{R_{coco}}^2 + \sigma_{R_{gray}}^2)}}$$

914
 915
 916
 917 *Image rank-order response reliability per neural site (ρ_{site}^{IRO}):* To estimate the reliability of the
 918 responses per site, we computed a Spearman-Brown corrected, split half (trial-based)
 919 correlation between the rank order of the image responses (all images).
 920

921 *Selectivity per neural site:* For each site, we measured selectivity as the d' for separating that
 922 site's best (highest response-driving) stimulus from its worst (lowest response-driving) stimulus.
 923 d' was computed by comparing the response mean of the site over all trials on the best stimulus
 924 as compared to the response mean of the site over all trials on the worst stimulus, and
 925 normalized by the square-root of the mean of the variances of the sites on the two stimuli:

$$selectivity_i = \frac{mean(\vec{b}_i) - mean(\vec{w}_i)}{\sqrt{\frac{var(\vec{b}_i) + var(\vec{w}_i)}{2}}}$$

926
 927 where \vec{b}_i is the vector of responses of site i to its best stimulus over all trials and \vec{w}_i is the vector
 928 of responses of site i to its worst stimulus. We computed this number in a cross-validated
 929 fashion, picking the best and worst stimulus on a subset of trials and then computing the
 930 selectivity measure on a separate set of trials, and averaging the selectivity value of 50 trial
 931 splits.
 932

933
 934
 935 *Inclusion criterion for neural sites:* For our analyses, we only included the neural recording sites
 936 that had an overall significant visual drive (d'_{visual}), an image rank order response reliability
 937 (ρ_{site}^{IRO}) that was greater than 0.6 and a selectivity score that was greater than 1. Given that most
 938 of our neural metrics are corrected by the estimated noise at each neural site, the criterion for
 939 selection of neural sites is not that critical. It was mostly done to reduce computation time and
 940 eliminate noisy recordings.
 941

942 Population Neural response latency estimation

943
 944 Onset latencies (t_{onset}) were determined as the earliest time from sample image onset when the
 945 firing rates of neurons were higher than one-tenth of the peak of its response. We averaged the
 946 latencies estimated across individual neural sites to compute the population latency.
 947

948 Peak latencies (t_{peak}) were estimated as the time of maximum response (firing rate) of a neural
949 site in response to an image. We averaged the peak latencies estimated across individual
950 neural sites to compute the population peak latency per image.

951
952 Both of these latency measures were computed across different sets of images (*control* and
953 *challenge*) as mentioned in the article.

954
955
956

957 **Estimation of solution for object identity per image**

958

959 IT cortex

960 To estimate what information downstream neurons could easily “read” from a given IT neural
961 population, we used a simple, biologically plausible linear decoder (i.e., linear classifiers), that
962 has been previously shown to link IT population activity and primate behavior⁶. Such decoders
963 are simple in that they can perform binary classifications by computing weighted sums (each
964 weight is analogous to the strength of synapse) of input features and separate the outputs
965 based on a decision boundary (analogous to a neuron’s spiking threshold). Here we have used
966 a support vector machine (SVM) algorithm with linear kernels. The SVM learning model
967 generates a decoder with a decision boundary that is optimized to best separate images of the
968 target object from images of the distractor objects. The optimization is done under a
969 regularization constraint that limits the complexity of the boundary. We used L2 (ridge)
970 regularization, where the objective function for the minimization comprises of an additional term
971 (to reduce model complexity),

972

$$973 \quad \text{L2 (penalty)} = \frac{\lambda}{2} \sum_{j=1}^p \beta_j^2$$

974

975 where β and p are the classifier weights associated with ‘p’ predictors (e.g. 424 neurons). The
976 strength of regularization, λ was optimized for each train-set and a stochastic gradient descent
977 solver was used to estimate 10 (one for each object) one-vs-all classifiers. After training each of
978 these classifiers with a set of 100 training images per object, we generated a class score (sc)
979 per classifier for all held out test images given by,

980

$$981 \quad sc = R\beta + bias,$$

982

983 where R is the population response vector and the bias is estimated by the SVM solver.

984

985 The train and test sets were pseudo-randomly chosen multiple times until every image of our
986 image set was part of the held-out test set. We then converted the class scores into probabilities
987 by passing them through a *softmax* (normalized exponential) function.

$$988 \quad P_{image}^i = \frac{e^{sc_i}}{\sum_{i=1}^{10} e^{sc_i}}$$

989 Our behavioral I_1 scores are all trial-averaged metrics. Therefore, in order to generate a
990 comparable trial-averaged performance per image — a probability for each classifier output,
991 given any image (P_{image}^i) was generated. The decoders are therefore trained and tested with
992 trial-averaged data.

993
994 We then computed the binary task performances, by calculating the percent correct score for
995 each pair of possible binary task given an image. For instance, if an image was from object i ,
996 then the percent correct score for the binary task between object i and object j , $Pr^{i,j}$ was
997 computed as,
998

$$Pr_{image}^{i,j} = \frac{P_{image}^i}{P_{image}^i + P_{image}^j}$$

999 From each percent correct score, we then estimated a neural I_1 score (per image), following the
1000 same procedures as the behavioral metric.

1001 Object solution time per image in IT (OST_{image})

1002

1003 Object solution time per image, OST_{image} was defined as the time it takes for linear IT
1004 population decodes to reach within the error margins of the pooled monkey behavioral I_1 score
1005 for that image. In order to estimate this time, we first computed a neural I_1 vector for non-
1006 overlapping 10 ms time bins post the sample image onset. We then used linear interpolation to
1007 predict the value of the I_1 vector per image at any given time between 0 and 250 ms. We then
1008 used the Levenberg-Marquardt algorithm to estimate the time at which the neural I_1 vector
1009 reached the error margins of the pooled monkey behavioral I_1 .
1010

1011 We balanced the control and challenge image populations at each level of
1012 the monkeys' performance. Therefore, we discarded challenge images that
1013 showed a d' of 5 or higher since there were no equivalent control images at
1014 that behavioral-accuracy level. However, we estimated the average OST
1015 for the challenge images at $d' \geq 5$ to be 150.2 ms (well within the range of
1016 other challenge image OSTs). Deep Convolutional Neural Networks
1017 (DCNN)

1018 **Binary object discrimination tasks with DCNNs**

1019 We have used two different techniques to train and test the DCNN features on the binary object
1020 discrimination task.

1021

1022 1. Back-end training (transfer learning): Here we have used the same linear decoding scheme
1023 mentioned above (for the IT neurons) to estimate the object solution strengths per image for the
1024 DCNNs. Briefly, we first obtained an ImageNet pre-trained DCNN (e.g AlexNet). We then

1025 replaced the last three layers (i.e. anything beyond 'fc7') of this network with a fully connected
1026 layer containing 10 nodes (each representing one of the 10 objects we have used in this study).
1027 We then trained this last layer with a back-end classifier (L2 regularized linear SVM; similar to
1028 the one mentioned for IT) on a subset of images from our image-set (containing both control
1029 and challenge images). These images were selected randomly from our imageset and used as
1030 the train-set. The remaining images were then used for the testing (such that there is no overlap
1031 between the train and test images). Repeating this procedure multiple times allowed us to use
1032 all images as test images providing us with the performance of the model for each image. The
1033 features extracted from each of the DCNN models were projected onto the first 1000 principle
1034 components (ranked in the order of variance explained) to construct the final feature set used.
1035 This was done to maintain consistency while comparing different layers across various DCNNs
1036 (some include ~20000 features) and control for the total number of features used in the
1037 analyses.

1038
1039 2. Fine-tuning: Although the steps mentioned above (transfer learning) is more similar to how
1040 we think the monkey implements the learning of the task in his brain, we cannot completely rule
1041 out the possibility that the representations of the images in IT do not change after training with
1042 our image-set. Prior work suggests that such IT population response changes are modest at
1043 best⁶⁵. Therefore, we also fine-tune (end-to-end) the ImageNet pre-trained AlexNet with images
1044 (randomly selected from our own image-set) and test on the remaining held out images. This
1045 technique also involves first obtaining an imagenet pertained DCNN, and replacing the final 3
1046 layers (e.g. beyond AlexNet 'fc7') with a fully connected layer of 10 nodes. However, the key
1047 difference of this technique with the transfer learning technique is that the new network is now
1048 trained end-to-end with stochastic gradient decent on separate training images from our own
1049 image-set used to test the monkeys. Figure S15 shows that the three main findings of our article
1050 (discovery of challenge images; lagged solutions for challenge images and lower IT predictivity
1051 for late-phase IT responses) are well replicated even with a fine-tuned ImageNet pre-trained
1052 AlexNet.

1053
1054

1055 **Prediction of neural response from DCNN features**

1056
1057 We modeled each IT neural site as a linear combination of the DCNN model features (illustrated
1058 in Figure S12). We first extracted the features per image, from the DCNNs' layers. The features
1059 extracted were then projected onto its first 1000 principle components (ranked in the order of
1060 variance explained) to construct the final feature set used. For example, we used the features
1061 from AlexNet's ²¹ 'fc7' layer to generate Figure 4A. Using a 50%/50% train/test split of the
1062 images, we then estimated the regression weights (i.e how we can linearly combine the model
1063 features to predict the neural site's responses) using a partial least squares (MATLAB
1064 command: *plsregress*) regression procedure, using 20 retained components. The neural
1065 responses used for training (R^{TRAIN}) and testing (R^{TEST}) the encoding models were averaged
1066 firing rates (measured at the specific sites) within the time window considered. We treated each
1067 time window (10 ms bins) independently for training and testing. The training images used for
1068 regressing the model features onto a neuron, at each time point, were sampled randomly
1069 (repeats included random subsampling) from the entire image set. For each set of regression

1070 weights (w) estimated on the training image responses (R^{TRAIN}), we generated the output of that
1071 ‘synthetic neuron’ for the held out test set (M^{PRED}) as

1072
1073
$$M^{PRED} = (w * F^{TEST}) + \beta,$$

1074 where w and β are estimated via the PLS regression and F^{TEST} are the model activation
1075 features for the test image-set.
1076

1077 The percentage of explained variance, *IT predictivity* (for details refer ⁸) for that neural site, was
1078 then computed by normalizing the r^2 prediction value for that site by the self-consistency of the
1079 test image responses ($\rho^{R^{TEST}}$) for that site and the self-consistency of the regression model
1080 predictions ($\rho^{M^{PRED}}$) for that site (estimated by a Spearman Brown corrected trial-split
1081 correlation score).
1082

1083
$$IT\ predictivity = \left(\frac{corr(R^{TEST}, M^{PRED})}{\sqrt{\rho^{R^{TEST}} * \rho^{M^{PRED}}}} \right)^2$$

1084
1085 To achieve accurate cross-validation results, we had to test the prediction of the model on held
1086 out image responses. But to make sure we have exposed the mapping procedure (mapping the
1087 model features on to individual IT neural sites) to images from the same full generative space
1088 and especially from both the control and challenge image categories, for each time step — we
1089 randomly sub-sampled image responses from the entire image set (measured at that specific
1090 time step). This ensured that the mapping step was exposed to exemplars from both the control
1091 and the challenge images groups.
1092

1093 Estimation of the OST prediction strength

1094 We compared how well different factors and $\Delta d'$ between monkey behavior and AlexNet ‘fc7’,
1095 predicted the differences in the object solution time (OST) estimates. Each image has an
1096 associated value for different image properties, either categorical e.g. occluded/non-occluded
1097 or continuous e.g. object size etc. We first divided the image-sets into two groups, *high* and *low*,
1098 for each factor. The *high* group for each factor contained images with values higher than 95th
1099 percentile of the factor distribution, and the *low* group contained the ones with values less than
1100 5th percentile of the distribution. For the categorical factor like occlusion, the *high* group
1101 contained images with occlusion and the *low* group contained images without occlusion. Then,
1102 for each factor we performed a one-way ANOVA with object solution time as the dependent
1103 variable. The rationale behind this test was if the experimenter(s) were to create image-sets
1104 based on any one of these factor, how likely is it expose a large difference between the OST
1105 values. Therefore, we used the F-value of the test (y-axis in Figure 6) to quantify the OST
1106 prediction strength.

1107

1108 Data and code availability

1109

1110 At the time of publishing, the images used in this study and the data associated with all the
1111 figures will be publicly available at our github repository (<https://github.com/kohitij-kar>). We will
1112 also host the images, primate behavioral scores, estimated object solution times, and the
1113 modeling results at <http://brain-score.org>⁴⁰.

1114

1115

1116

1117 Figure Caption

1118

1119 **Figure 1.** *Behavioral screening and identification of control and challenge images.* A) We task both
1120 primates (humans and macaques; top row) and feedforward DCNNs (bottom row) to identify which object
1121 is present in each Test image (1320 images). The top row shows the stages in the ventral visual pathway
1122 in primates (retina, LGN: lateral geniculate nucleus, areas V1, V2, V4, and IT), which is implicated in core
1123 object recognition. We can conceptualize each stage as rapidly transforming the representation of the
1124 image ultimately yielding to the primates' behavior (i.e. producing a behavioral report of which object was
1125 present). The blue arrows indicate the known anatomical feedforward projections from one area to the
1126 other. The red arrows indicate the known lateral and top down recurrent connections. The bottom row
1127 demonstrate a schematic of a similar pathway commonly present in the DCNNs. These networks contain
1128 a series of convolutional and pooling layers with nonlinear transforms at each stage, followed by fully
1129 connected layers (which approximates macaque IT neural responses) that ultimately gives rise to the
1130 models' "behavior." Note that the DCNNs only have feedforward (blue) connections. B) Object categories.
1131 We used ten different object types; bear, elephant, face, plane, dog, car, apple, chair, bird and zebra. C)
1132 Binary object discrimination task. Here we show the timeline of events on each trial. Subjects fixate a dot.
1133 The test image (8 °) containing one of ten possible objects was shown for 100 ms. After a 100 ms delay,
1134 a canonical view of the target object (the same noun as that present in the test image) and a distractor
1135 object (from the other 9 objects) appeared, and the human or monkey indicated which object was present
1136 in the test image by clicking on or making a saccade to one of the two choices respectively. D)
1137 Comparison of monkey performance (pooled across 2 monkeys) and DCNN performance (AlexNet; 'fc7'
1138 ²¹). Each dot represents the behavioral task performance (I_1 ; refer Methods) for a single image. We
1139 reliably identified *challenge* images (red dots) and *control* images (blue dots). Error bars are bootstrapped
1140 s.e.m. E) Examples of four *challenge* and four *control* images.

1141

1142 **Figure 2.** Large scale multiunit array recordings in the macaque inferior temporal cortex. A) Schematic of
1143 array placement, neural data recording and object solution time estimation. We recorded extracellular
1144 voltage in IT from two monkeys, each hemisphere implanted with 2 or 3 Utah arrays. For each image
1145 presentation (100 ms), we counted multiunit spike events (see Methods for details), per site, in non
1146 overlapping 10 ms windows, post stimulus onset to construct a single population activity vector per time
1147 bin. These population vectors (image evoked neural features) were then used to train and test cross-
1148 validated linear support vector machine decoders (d) separately per time bin. The decoder outputs per
1149 image (over time) were then used to perform a binary match to sample task, and obtain neural decode
1150 accuracies (NDA) at each time bin. An example of the neural decode accuracy over time is shown in the
1151 top panel. The time at which the neural decodes equal the primate (monkey) performance, is then
1152 recorded as the object solution time (OST) for that specific image. B) Examples of IT population decodes
1153 over time, with the estimated object solution times for four images; two *control* (top panel: blue curves)
1154 and two *challenge* images (bottom panel: red curves). The red and blue dots are the estimated neural
1155 decode accuracies at each time bins. The solid lines are nonlinear fits of the decoder accuracies over
1156 time (see Methods). The gray lines indicate the I_1 performance of the primates (pooled monkey) for the
1157 specific images. Error bar indicates bootstrapped s.e.m. C) Distribution of object solution times for both
1158 *control* (blue) and *challenge* (red) images. The median OST for *control* (blue) and *challenge* (red) images
1159 are shown in the plot with dashed lines. The inset in the top shows the median evolution of IT decodes
1160 over time until the OST for control (blue) and challenge (red) images.

1161

1162

1163

1164 **Figure 3.** Relationship between object solution times and neural response latencies. A) Comparison of
1165 neural responses evoked by *control* (blue) and *challenge* (red) images. We estimated two measures of
1166 population response latency: Population onset latency (t_{onset}) and Population peak latency (t_{peak}). B)
1167 Distributions of the population onset latencies (median across 424 sites), population peak response
1168 latencies (median across 424 sites) and object solution times for *control* images ($n=149$). C) Same as in
1169 B) but for *challenge* images ($n = 266$). D) Comparison of population onset latencies and object solution
1170 times for both *control* (blue) and *challenge* images (red). Vertical error bars show s.e.m across neurons
1171 and horizontal error bars show bootstrap (across trial repetition) standard deviation of *OST* estimates.

1172

1173

1174

1175

1176 **Figure 4.** Predicting IT neural responses with DCNN features. A) IT predictivity of AlexNet's 'fc7' layer as
1177 a function of object solution time (ms). For each time bin, we consider IT predictivity only for images that
1178 have a solution time equal to or higher than that time bin. Error bars indicate the standard error of mean
1179 across neurons. Top panel shows the distribution of object solution times for *control* (blue) and *challenge*
1180 (red) images. B) IT predictivity computed separately for late OST images (OST>150 ms; total of 349
1181 images) at the corresponding object solution times, as function of deep (AlexNet, Zeiler and Fergus,
1182 VGG-S), deeper (Inception, ResNet) CNNs and deep-recurrent CNNs (CORnet). * indicates a statistically
1183 significant difference between two groups. The inset to the right shows a schematic representation of
1184 CORnet that has recurrent connections (shown in red) at each layer (V1, V2, V4 and IT) C) Comparison
1185 of median OST for different sets of *challenge* images: the set of *challenge* images is defined with respect
1186 to each DCNN model (thus, the exact set of images is different for each bar, and the number of images is
1187 indicated on top of the bars). In each case, the *challenge* images are defined as the set of images that
1188 remain unsolved by each model (using the fixed definitions of this study, see text). Note that the use of
1189 deeper CNNs and the deep-recurrent CNN, resulted in the discovery of *challenge* images that required
1190 even longer *OST*s in IT cortex than the original set *challenge* images (defined for AlexNet 'fc7'). * indicates
1191 a statistically significant difference between two groups.

1192

1193

1194 **Figure 5.** A) Binary object discrimination with backward visual masking. The test image (presented for 34,
1195 67, 100, 134 or 267 ms) was followed immediately by a visual mask (phase scrambled image) for 500 ms,
1196 followed by a blank gray screen for 100 ms and then the object choice screen. Monkeys reported the
1197 target object by fixating it on the choice screen. B) Difference in behavioral performance between *control*
1198 and *challenge* image after backward visual masking. Each bar on the plot (y-axis) is the difference in the
1199 pooled monkey performance during the visual masking task (A) between *control* and *challenge* images at
1200 the respective sample image presentation durations (x-axis). The top panel inset shows the raw
1201 performance (d') for the two groups of images (blue: *control* images, red: *challenge* images). Error bars
1202 denote the standard error of mean across all objects.

1203

1204 **Figure 6.** Comparison of *OST* prediction strength between different image properties, a combination of all
1205 estimated image properties, and the $\Delta d'$ vector (deviation of model behavior from pooled monkey
1206 behavior). The red dashed line denotes the significance threshold of the F-statistic. Image properties like
1207 object size, eccentricity, presence of an occluder, as well as a combination of these properties (referred to
1208 as “all-factors”) significantly predict *OST*. However, the $\Delta d'$ vector provides the strongest *OST*
1209 predictions. Error bars denote the bootstrap standard deviation over images. * denotes a significant
1210 difference between the two groups — image properties vs $\Delta d'$, estimated with repeated measures
1211 ANOVA ($F(1,10) > 100$; $p < 0.0001$; multiple-comparison using Turkey test showed a significant difference
1212 between $\Delta d'$ and all other image properties).

1213

1214

1215

1216

1217

Bibliography

1218

1219 1. DiCarlo, J.J., Zoccolan, D. & Rust, N.C. How does the brain solve visual object
1220 recognition? *Neuron* **73**, 415-434 (2012).

1221 2. Riesenhuber, M. & Poggio, T. Models of object recognition. *Nat Neurosci* **3 Suppl**, 1199-
1222 1204 (2000).

1223 3. Yamins, D.L. & DiCarlo, J.J. Eight open questions in the computational modeling of
1224 higher sensory cortex. *Curr Opin Neurobiol* **37**, 114-120 (2016).

1225 4. Hung, C.P., Kreiman, G., Poggio, T. & DiCarlo, J.J. Fast readout of object identity from
1226 macaque inferior temporal cortex. *Science* **310**, 863-866 (2005).

1227 5. Freiwald, W.A., Tsao, D.Y. & Livingstone, M.S. A face feature space in the macaque
1228 temporal lobe. *Nat Neurosci* **12**, 1187-1196 (2009).

1229 6. Majaj, N.J., Hong, H., Solomon, E.A. & DiCarlo, J.J. Simple Learned Weighted Sums of
1230 Inferior Temporal Neuronal Firing Rates Accurately Predict Human Core Object Recognition
1231 Performance. *J Neurosci* **35**, 13402-13418 (2015).

1232 7. Cadieu, C.F., *et al.* Deep neural networks rival the representation of primate IT cortex for
1233 core visual object recognition. *PLoS Comput Biol* **10**, e1003963 (2014).

1234 8. Yamins, D.L., *et al.* Performance-optimized hierarchical models predict neural responses
1235 in higher visual cortex. *Proc Natl Acad Sci U S A* **111**, 8619-8624 (2014).

1236 9. Guclu, U. & van Gerven, M.A. Deep Neural Networks Reveal a Gradient in the
1237 Complexity of Neural Representations across the Ventral Stream. *J Neurosci* **35**, 10005-10014
1238 (2015).

1239 10. Rajalingham, R., Schmidt, K. & DiCarlo, J.J. Comparison of Object Recognition Behavior
1240 in Human and Monkey. *J Neurosci* **35**, 12127-12136 (2015).

1241 11. Rajalingham, R., *et al.* Large-scale, high-resolution comparison of the core visual object
1242 recognition behavior of humans, monkeys, and state-of-the-art deep artificial neural networks.
1243 *bioRxiv*, 240614 (2018).

1244 12. Rockland, K.S. & Virga, A. Terminal arbors of individual "feedback" axons projecting
1245 from area V2 to V1 in the macaque monkey: a study using immunohistochemistry of
1246 anterogradely transported Phaseolus vulgaris-leucoagglutinin. *J Comp Neurol* **285**, 54-72
1247 (1989).

1248 13. Felleman, D.J. & Van Essen, D.C. Distributed hierarchical processing in the primate
1249 cerebral cortex. *Cereb Cortex* **1**, 1-47 (1991).

1250 14. Rockland, K.S., Saleem, K.S. & Tanaka, K. Divergent feedback connections from areas
1251 V4 and TEO in the macaque. *Vis Neurosci* **11**, 579-600 (1994).

1252 15. Rockland, K.S. & Van Hoesen, G.W. Direct temporal-occipital feedback connections to
1253 striate cortex (V1) in the macaque monkey. *Cereb Cortex* **4**, 300-313 (1994).

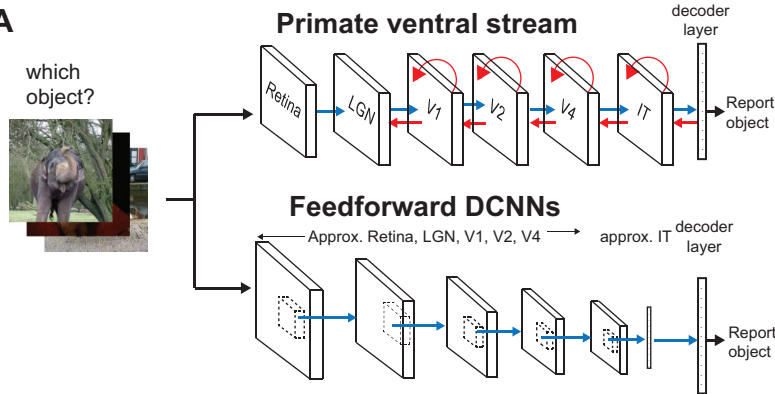
1254 16. Thorpe, S., Fize, D. & Marlot, C. Speed of processing in the human visual system.
1255 *Nature* **381**, 520-522 (1996).

1256 17. Liu, H., Agam, Y., Madsen, J.R. & Kreiman, G. Timing, timing, timing: fast decoding of
1257 object information from intracranial field potentials in human visual cortex. *Neuron* **62**, 281-290
1258 (2009).

- 1259 18. Hinton, G.E., Dayan, P., Frey, B.J. & Neal, R.M. The "wake-sleep" algorithm for
1260 unsupervised neural networks. *Science* **268**, 1158-1161 (1995).
- 1261 19. Geirhos, R., *et al.* Comparing deep neural networks against humans: object recognition
1262 when the signal gets weaker. *arXiv preprint arXiv:1706.06969* (2017).
- 1263 20. Lamme, V.A. & Roelfsema, P.R. The distinct modes of vision offered by feedforward and
1264 recurrent processing. *Trends Neurosci* **23**, 571-579 (2000).
- 1265 21. Krizhevsky, A., Sutskever, I. & Hinton, G.E. ImageNet classification with deep
1266 convolutional neural networks. in *Proceedings of the 25th International Conference on Neural*
1267 *Information Processing Systems - Volume 1* 1097-1105 (Curran Associates Inc., Lake Tahoe,
1268 Nevada, 2012).
- 1269 22. Pinto, N., Cox, D.D. & DiCarlo, J.J. Why is real-world visual object recognition hard?
1270 *PLoS Comput Biol* **4**, e27 (2008).
- 1271 23. Lin, T.-Y., *et al.* Microsoft coco: Common objects in context. in *European conference on*
1272 *computer vision* 740-755 (Springer, 2014).
- 1273 24. Sermanet, P., *et al.* Overfeat: Integrated recognition, localization and detection using
1274 convolutional networks. *arXiv preprint arXiv:1312.6229* (2013).
- 1275 25. Chatfield, K., Simonyan, K., Vedaldi, A. & Zisserman, A. Return of the devil in the
1276 details: Delving deep into convolutional nets. *arXiv preprint arXiv:1405.3531* (2014).
- 1277 26. Zeiler, M.D. & Fergus, R. Visualizing and understanding convolutional networks. in
1278 *European conference on computer vision* 818-833 (Springer, 2014).
- 1279 27. Meyers, E.M., Freedman, D.J., Kreiman, G., Miller, E.K. & Poggio, T. Dynamic
1280 population coding of category information in inferior temporal and prefrontal cortex. *J*
1281 *Neurophysiol* **100**, 1407-1419 (2008).
- 1282 28. Oram, M.W. Contrast induced changes in response latency depend on stimulus
1283 specificity. *J Physiol Paris* **104**, 167-175 (2010).
- 1284 29. Rolls, E.T., Baylis, G.C. & Leonard, C.M. Role of low and high spatial frequencies in the
1285 face-selective responses of neurons in the cortex in the superior temporal sulcus in the monkey.
1286 *Vision Res* **25**, 1021-1035 (1985).
- 1287 30. Op De Beeck, H. & Vogels, R. Spatial sensitivity of macaque inferior temporal neurons. *J*
1288 *Comp Neurol* **426**, 505-518 (2000).
- 1289 31. Willenbockel, V., *et al.* Controlling low-level image properties: the SHINE toolbox. *Behav*
1290 *Res Methods* **42**, 671-684 (2010).
- 1291 32. McKee, J.L., Riesenhuber, M., Miller, E.K. & Freedman, D.J. Task dependence of visual
1292 and category representations in prefrontal and inferior temporal cortices. *J Neurosci* **34**, 16065-
1293 16075 (2014).
- 1294 33. Bugatus, L., Weiner, K.S. & Grill-Spector, K. Task alters category representations in
1295 prefrontal but not high-level visual cortex. *Neuroimage* **155**, 437-449 (2017).
- 1296 34. Liao, Q. & Poggio, T. Bridging the gaps between residual learning, recurrent neural
1297 networks and visual cortex. *arXiv preprint arXiv:1604.03640* (2016).
- 1298 35. Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J. & Wojna, Z. Rethinking the inception
1299 architecture for computer vision. in *Proceedings of the IEEE Conference on Computer Vision*
1300 *and Pattern Recognition* 2818-2826 (2016).
- 1301 36. Szegedy, C., Ioffe, S., Vanhoucke, V. & Alemi, A.A. Inception-v4, inception-resnet and
1302 the impact of residual connections on learning. in *AAAI* 12 (2017).

- 1303 37. He, K., Zhang, X., Ren, S. & Sun, J. Deep residual learning for image recognition. in
1304 *Proceedings of the IEEE conference on computer vision and pattern recognition* 770-778
1305 (2016).
- 1306 38. Russakovsky, O., *et al.* Imagenet large scale visual recognition challenge. *International*
1307 *Journal of Computer Vision* **115**, 211-252 (2015).
- 1308 39. Kubilius, J., *et al.* CORnet: Modeling the Neural Mechanisms of Core Object
1309 Recognition. *bioRxiv*, 408385 (2018).
- 1310 40. Schrimpf, M., *et al.* Brain-Score: Which Artificial Neural Network for Object Recognition
1311 is most Brain-Like? *bioRxiv*, 407007 (2018).
- 1312 41. Stojanoski, B. & Cusack, R. Time to wave good-bye to phase scrambling: creating
1313 controlled scrambled images using diffeomorphic transformations. *J Vis* **14** (2014).
- 1314 42. Lamme, V.A., Zipser, K. & Spekreijse, H. Masking interrupts figure-ground signals in V1.
1315 *J Cogn Neurosci* **14**, 1044-1053 (2002).
- 1316 43. Fahrenfort, J.J., Scholte, H.S. & Lamme, V.A. Masking disrupts reentrant processing in
1317 human visual cortex. *J Cogn Neurosci* **19**, 1488-1497 (2007).
- 1318 44. Elsayed, G.F., *et al.* Adversarial examples that fool both human and computer vision.
1319 *arXiv preprint arXiv:1802.08195* (2018).
- 1320 45. Walther, D., Rutishauser, U., Koch, C. & Perona, P. Selective visual attention enables
1321 learning and recognition of multiple objects in cluttered scenes. *Computer Vision and Image*
1322 *Understanding* **100**, 41-63 (2005).
- 1323 46. Spoerer, C.J., McClure, P. & Kriegeskorte, N. Recurrent Convolutional Neural Networks:
1324 A Better Model of Biological Object Recognition. *Front Psychol* **8**, 1551 (2017).
- 1325 47. Tang, H., *et al.* Recurrent computations for visual pattern completion. *arXiv preprint*
1326 *arXiv:1706.02240* (2017).
- 1327 48. Bichot, N.P., Heard, M.T., DeGennaro, E.M. & Desimone, R. A Source for Feature-
1328 Based Attention in the Prefrontal Cortex. *Neuron* **88**, 832-844 (2015).
- 1329 49. Lehky, S.R. & Tanaka, K. Neural representation for object recognition in inferotemporal
1330 cortex. *Curr Opin Neurobiol* **37**, 23-35 (2016).
- 1331 50. Sugase, Y., Yamane, S., Ueno, S. & Kawano, K. Global and fine information coded by
1332 single neurons in the temporal visual cortex. *Nature* **400**, 869-873 (1999).
- 1333 51. Tovee, M.J. Neuronal processing. How fast is the speed of thought? *Curr Biol* **4**, 1125-
1334 1127 (1994).
- 1335 52. Ullman, S. Sequence seeking and counter streams: a computational model for
1336 bidirectional information flow in the visual cortex. *Cereb Cortex* **5**, 1-11 (1995).
- 1337 53. van Kerkoerle, T., *et al.* Alpha and gamma oscillations characterize feedback and
1338 feedforward processing in monkey visual cortex. *Proc Natl Acad Sci U S A* **111**, 14332-14341
1339 (2014).
- 1340 54. Nurminen, L., Merlin, S., Bijanzadeh, M., Federer, F. & Angelucci, A. Top-down
1341 feedback controls spatial summation and response amplitude in primate visual cortex. *Nat*
1342 *Commun* **9**, 2281 (2018).
- 1343 55. Fyall, A.M., El-Shamayleh, Y., Choi, H., Shea-Brown, E. & Pasupathy, A. Dynamic
1344 representation of partially occluded objects in primate prefrontal and visual cortex. *Elife* **6**
1345 (2017).

- 1346 56. Tomita, H., Ohbayashi, M., Nakahara, K., Hasegawa, I. & Miyashita, Y. Top-down signal
1347 from prefrontal cortex in executive control of memory retrieval. *Nature* **401**, 699-703 (1999).
- 1348 57. Bar, M., *et al.* Top-down facilitation of visual recognition. *Proc Natl Acad Sci U S A* **103**,
1349 449-454 (2006).
- 1350 58. Seger, C.A. How do the basal ganglia contribute to categorization? Their roles in
1351 generalization, response selection, and learning via feedback. *Neurosci Biobehav Rev* **32**, 265-
1352 278 (2008).
- 1353 59. Simonyan, K. & Zisserman, A. Very deep convolutional networks for large-scale image
1354 recognition. *arXiv preprint arXiv:1409.1556* (2014).
- 1355 60. Havenith, M.N., *et al.* Synchrony makes neurons fire in sequence, and stimulus
1356 properties determine who is ahead. *J Neurosci* **31**, 8570-8584 (2011).
- 1357 61. Jeurissen, D., Self, M.W. & Roelfsema, P.R. Serial grouping of 2D-image regions with
1358 object-based attention in humans. *Elife* **5** (2016).
- 1359 62. Roelfsema, P.R. Cortical algorithms for perceptual grouping. *Annu Rev Neurosci* **29**,
1360 203-227 (2006).
- 1361 63. Santos, A., *et al.* Evaluation of autofocus functions in molecular cytogenetic analysis. *J*
1362 *Microsc* **188**, 264-272 (1997).
- 1363 64. Rosenholtz, R., Li, Y. & Nakano, L. Measuring visual clutter. *J Vis* **7**, 17 11-22 (2007).
- 1364 65. Baker, C.I., Behrmann, M. & Olson, C.R. Impact of learning on representation of parts
1365 and wholes in monkey inferotemporal cortex. *Nat Neurosci* **5**, 1210-1216 (2002).
- 1366 66. Buffalo, E.A., Fries, P., Landman, R., Liang, H. & Desimone, R. A backward progression
1367 of attentional effects in the ventral stream. *Proc Natl Acad Sci U S A* **107**, 361-365 (2010).
- 1368 67. Hong, H., Yamins, D.L., Majaj, N.J. & DiCarlo, J.J. Explicit information for category-
1369 orthogonal object properties increases along the ventral stream. *Nat Neurosci* **19**, 613-622
1370 (2016).
- 1371 68. Lin, J.Y., Hubert-Wallander, B., Murray, S.O. & Boynton, G.M. Rapid and reflexive
1372 feature-based attention. *J Vis* **11** (2011).
- 1373 69. Maunsell, J.H. & Treue, S. Feature-based attention in visual cortex. *Trends Neurosci* **29**,
1374 317-322 (2006).
- 1375 70. Poort, J., *et al.* The role of attention in figure-ground segregation in areas V1 and V4 of
1376 the visual cortex. *Neuron* **75**, 143-156 (2012).
- 1377 71. Salti, M., *et al.* Distinct cortical codes and temporal dynamics for conscious and
1378 unconscious percepts. *Elife* **4** (2015).

A**B**

Object categories

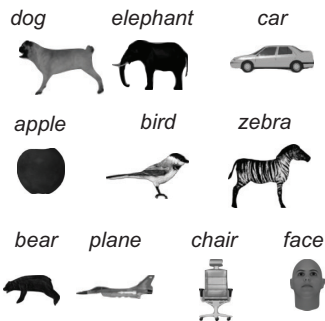
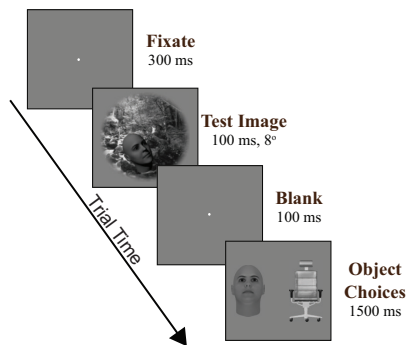
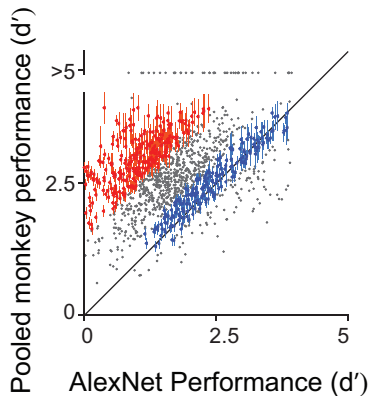
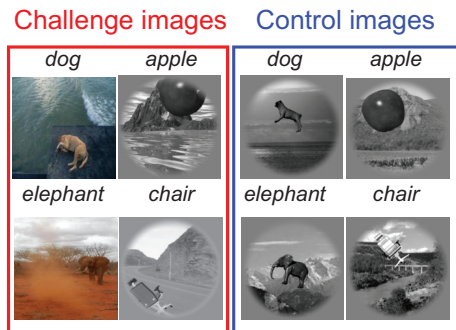
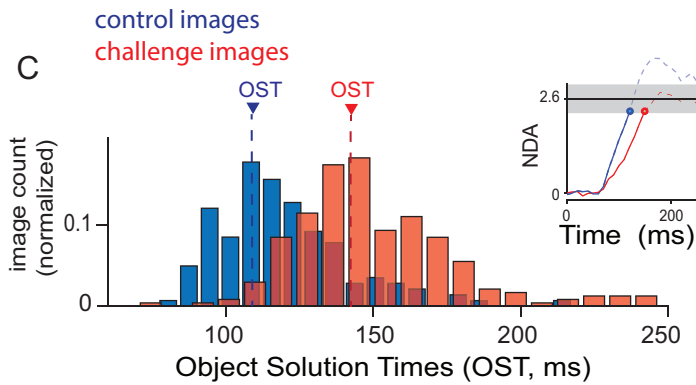
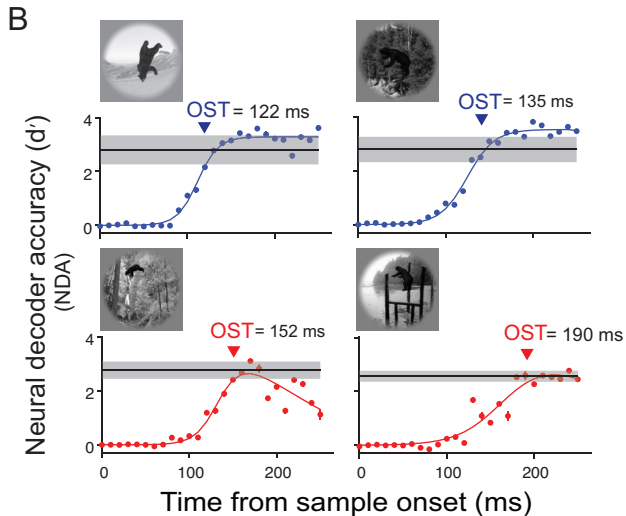
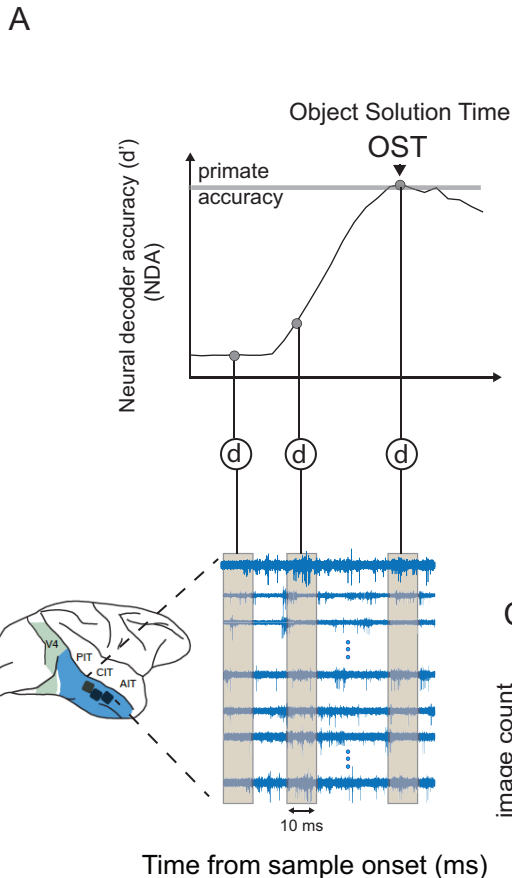
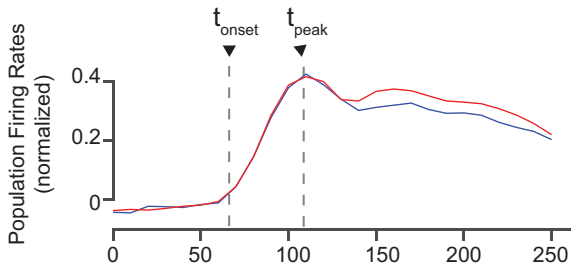
**C** Object Discrimination Task**D** Behavioral comparison
monkeys vs DCNN**E**

Image examples

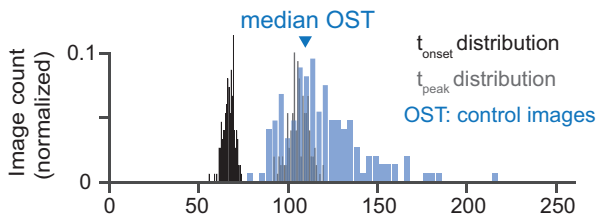




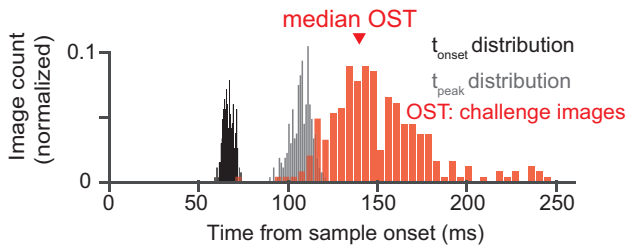
A



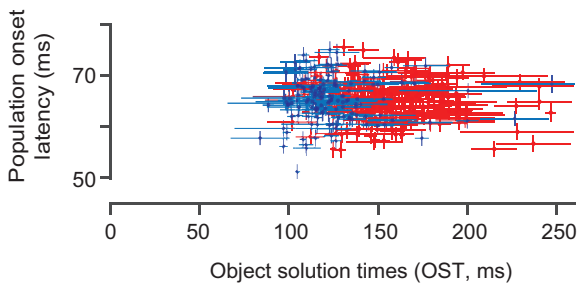
B



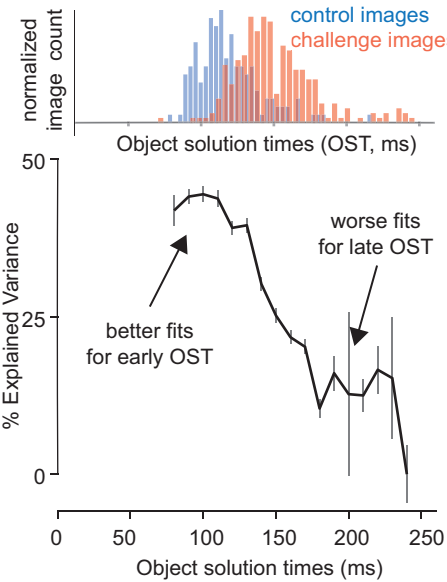
C



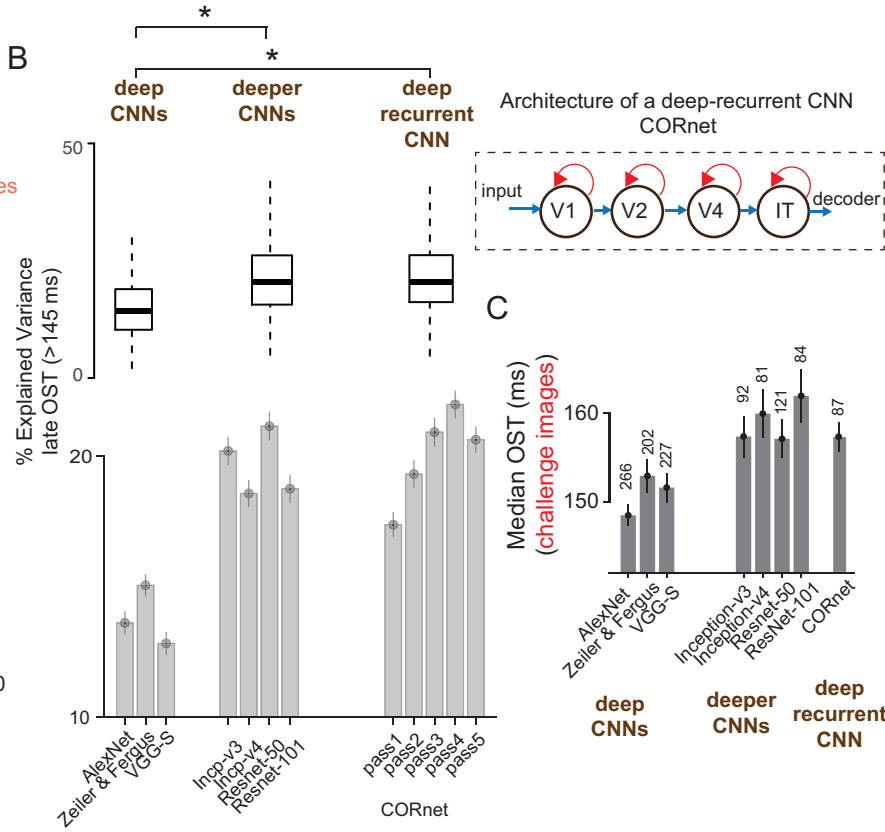
D



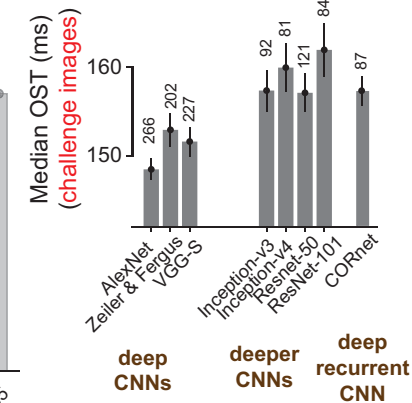
A



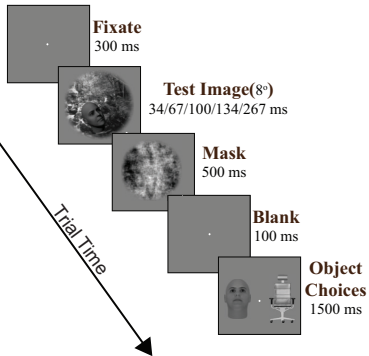
B



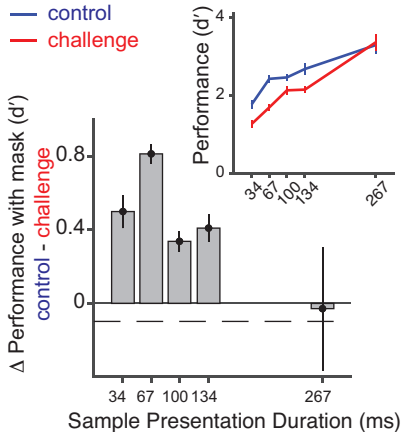
C



A



B



Ability to predict "recurrence"

