# Large-scale, high-resolution comparison of the core visual object recognition behavior of humans, monkeys, and state-of-the-art deep artificial neural networks

Abbreviated title: Comparing object recognition between primates and models

Rishi Rajalingham*, Elias B. Issa*, Pouya Bashivan, Kohitij Kar, Kailyn Schmidt, and James J. DiCarlo

McGovern Institute for Brain Research and Department of Brain and Cognitive Sciences
Massachusetts Institute of Technology, Cambridge, Massachusetts 02139

*R.R. and E.B.I. contributed equally to this work.

Correspondence should be addressed to James J. DiCarlo, McGovern Institute for Brain Research, Department of Brain and Cognitive Sciences, Massachusetts Institute of Technology, 77 Massachusetts Institute of Technology, 46-6161, Cambridge, MA 02139. E-mail: dicarlo@mit.edu

E. Issa's present address: Department of Neuroscience, Zuckerman Mind Brain Behavior Institute, Columbia University, New York, NY 10027

_____

Targeting *Journal of Neuroscience*
Title 19 words
Abbreviated Title 50 characters
Abstract 242 words
Significance Statement 97 words
Introduction 649 words
Discussion 1188 words
Figures 6
*All word limits include citations
_____


AUTHOR CONTRIBUTIONS
E.B.I., R.R., and J.J.D designed the experiments. E.B.I., K.S., R.R., and K.K. carried out the experiments. R.R., E.B.I., and P.B. performed the data analysis and modeling. R.R., E.B.I., and J.J.D. wrote the manuscript.

COMPETING FINANCIAL INTERESTS
The authors declare no competing financial interests.

47  **ABSTRACT**

48

49       Primates—including humans—can typically recognize objects in visual images at a

50  glance even in the face of naturally occurring identity-preserving image transformations (e.g.

51  changes in viewpoint). A primary neuroscience goal is to uncover neuron-level mechanistic

52  models that quantitatively explain this behavior by predicting primate performance for each and

53  every image. Here, we applied this stringent behavioral prediction test to the leading mechanistic

54  models of primate vision (specifically, deep, convolutional, artificial neural networks; ANNs) by

55  directly comparing their behavioral signatures against those of humans and rhesus macaque

56  monkeys. Using high-throughput data collection systems for human and monkey psychophysics,

57  we collected over one million behavioral trials for 2400 images over 276 binary object

58  discrimination tasks. Consistent with previous work, we observed that state-of-the-art deep, feed-

59  forward convolutional ANNs trained for visual categorization (termed $DCNN_{IC}$ models)

60  accurately predicted primate patterns of object-level confusion. However, when we examined

61  behavioral performance for individual images within each object discrimination task, we found

62  that all tested $DCNN_{IC}$ models were significantly non-predictive of primate performance, and

63  that this prediction failure was not accounted for by simple image attributes, nor rescued by

64  simple model modifications. These results show that current $DCNN_{IC}$ models cannot account for

65  the image-level behavioral patterns of primates, and that new ANN models are needed to more

66  precisely capture the neural mechanisms underlying primate object vision. To this end, large-

67  scale, high-resolution primate behavioral benchmarks—such as those obtained here—could serve

68  as direct guides for discovering such models.

69

70

71   **SIGNIFICANCE STATEMENT**

72

73       Recently, specific feed-forward deep convolutional artificial neural networks (ANNs)

74   models have dramatically advanced our quantitative understanding of the neural mechanisms

75   underlying primate core object recognition. In this work, we tested the limits of those ANNs by

76   systematically comparing the behavioral responses of these models with the behavioral responses

77   of humans and monkeys, at the resolution of individual images. Using these high-resolution

78   metrics, we found that all tested ANN models significantly diverged from primate behavior.

79   Going forward, these high-resolution, large-scale primate behavioral benchmarks could serve as

80   direct guides for discovering better ANN models of the primate visual system.

81

82    **INTRODUCTION**

83

84          Primates—both human and non-human—can typically recognize objects in visual images
85    at a glance, even in the face of naturally occurring identity-preserving transformations such as
86    changes in viewpoint. This view-invariant visual object recognition ability is thought to be
87    supported primarily by the primate ventral visual stream (DiCarlo et al., 2012). A primary
88    neuroscience goal is to construct computational models that quantitatively explain the neural
89    mechanisms underlying this ability. That is, our goal is to discover artificial neural networks
90    (ANNs) that accurately predict neuronal firing rate responses at all levels of the ventral stream
91    and its behavioral output. To this end, specific models within a large family of deep,
92    convolutional neural networks (DCNNs), optimized by supervised training on large-scale
93    category-labeled image-sets (ImageNet) to match human-level categorization performance
94    (Krizhevsky et al., 2012; LeCun et al., 2015), have been put forth as the leading ANN models of
95    the ventral stream (Yamins and DiCarlo, 2016). We refer to this sub-family as $DCNN_{IC}$ models
96    (IC to denote ImageNet-categorization pre-training), so as to distinguish them from all possible
97    models in the DCNN family, and more broadly, from the super-family of all ANNs. To date, it
98    has been shown that $DCNN_{IC}$ models display internal feature representations similar to neuronal
99    representations along the primate ventral visual stream (Yamins et al., 2013; Cadieu et al., 2014;
100   Khaligh-Razavi and Kriegeskorte, 2014; Yamins et al., 2014), and they exhibit behavioral
101   patterns similar to the behavioral patterns of pairwise object confusions of primates
102   (Rajalingham et al., 2015). Thus, $DCNN_{IC}$ models may provide a quantitative account of the
103   neural mechanisms underlying primate core object recognition behavior.

104

105         However, several studies have shown that $DCNN_{IC}$ models can diverge drastically from
106   humans in object recognition behavior, especially with regards to particular images optimized to
107   be adversarial to these networks (Goodfellow et al., 2014; Nguyen et al., 2015). Related work
108   has shown that specific image distortions are disproportionately challenging to current DCNNs,
109   as compared to humans (RichardWebster et al., 2016; Dodge and Karam, 2017; Geirhos et al.,
110   2017; Hosseini et al., 2017). Such image-specific failures of the current ANN models would
111   likely not be captured by "object-level" behavioral metrics (e.g. the pattern of pairwise object
112   confusions mentioned above) that are computed by pooling over hundreds of images and thus are

113    not sensitive to variation in difficulty across images of the same object. To overcome this

114    limitation of prior work, we here aimed to use scalable behavioral testing methods to precisely

115    characterize primate behavior at the resolution of individual images and to directly compare

116    leading DCNN models to primates over the domain of core object recognition behavior at this

117    high resolution.

118

119         We focused on *core invariant object recognition*—the ability to identify objects in visual

120    images in the central visual field during a single, natural viewing fixation (DiCarlo et al., 2012).

121    We further restricted our behavioral domain to *basic-level* object discriminations, as defined

122    previously (Rosch et al., 1976). Within this domain, we collected large-scale, high-resolution

123    measurements of human and monkey behavior (over a million behavioral trials) using high-

124    throughput psychophysical techniques—including a novel home-cage behavioral system for

125    monkeys. These data enabled us to systematically compare all systems at progressively higher

126    resolution. At lower resolutions, we replicated previous findings that humans, monkeys, and

127    DCNN$_{IC}$ models all share a common pattern of object-level confusion (Rajalingham et al., 2015).

128    However, at the higher resolution of individual images, we found that the behavior of all tested

129    DCNN$_{IC}$ models was significantly different from human and monkey behavior, and this model

130    prediction failure could not be easily rescued by simple model modifications. These results show

131    that current DCNN$_{IC}$ models do not fully account for the image-level behavioral patterns of

132    primates, suggesting that new ANN models are needed to more precisely capture the neural

133    mechanisms underlying primate object vision. To this end, large-scale high-resolution behavioral

134    benchmarks, such as those obtained here, could serve as a strong top-down constraint for

135    efficiently discovering such models.

136

137

138    **MATERIALS & METHODS**

139

140    *Visual images*

141         We examined basic-level, core object recognition behavior using a set of 24 broadly-

142    sampled objects that we previously found to be reliably labeled by independent human subjects,

143    based on the definition of basic-level proposed by (Rosch et al., 1976). For each object, we

144    generated 100 naturalistic synthetic images by first rendering a 3D model of the object with

145    randomly chosen viewing parameters (2D position, 3D rotation and viewing distance), and then

146    placing that foreground object view onto a randomly chosen, natural image background. To do

147    this, each object was first assigned a canonical position (center of gaze), scale (~2 degrees) and

148    pose, and then its viewing parameters were randomly sampled uniformly from the following

149    ranges for object translation ([-3,3] degrees in both h and v), rotation ([-180,180] degrees in all

150    three axes) and scale ([x0.7, x1.7]. Background images were sampled randomly from a large

151    database of high-dynamic range images of indoor and outdoor scenes obtained from Dosch

152    Design (www.doschdesign.com). This image generation procedure enforces invariant object

153    recognition, rather than image matching, as it requires the visual recognition system (human,

154    animal or model) to tackle the "invariance problem," the computational crux of object

155    recognition (Ullman and Humphreys, 1996; Pinto et al., 2008). Using this procedure, we

156    previously generated 2400 images (100 images per object) rendered at 1024x1024 pixel

157    resolution with 256-level gray scale and subsequently resized to 256x256 pixel resolution for

158    human psychophysics, monkey psychophysics and model evaluation (Rajalingham et al., 2015).

159    In the current work, we focused our analyses on a randomly subsampled, and then fixed, sub-set

160    of 240 images (10 images per object; here referred to as the "primary test images"). Figure 1A

161    shows the full list of 24 objects, with two example images of each object.

162

163    Because all of the images were generated from synthetic 3D object models, we had

164    explicit knowledge of the viewpoint parameters (position, size, and pose) for each object in each

165    image, as well as perfect segmentation masks. Taking advantage of this feature, we characterized

166    each image based on these high-level attributes, consisting of size, eccentricity, relative pose and

167    contrast of the object in the image. The size and eccentricity of the object in each image were

168    computed directly from the corresponding viewpoint parameters, under the assumption that the

169    entire image would subtend 6° at the center of visual gaze (+/-3° in both azimuth and elevation;

170    see below). For each synthetic object, we first defined its "canonical" 3D pose vector, based on

171    independent human judgments. To compute the relative pose attribute of each image, we

172    estimated the difference between the object's 3D pose and its canonical 3D pose. Pose

173    differences were computed as distances in unit quaternion representations: the 3D pose ($r_{xy}$, $r_{xz}$,

174    $r_{yz}$) was first converted into unit quaternions, and distances between quaternions $q_1$, $q_2$ were

175     estimated as $\cos^{-1}|q_1 \cdot q_2|$ (Huynh, 2009). To compute the object contrast, we measured the

176     absolute difference between the mean of the pixel intensities corresponding to the object and the

177     mean of the background pixel intensities in the vicinity of the object (specifically, within 25

178     pixels of any object pixel, analogous to computing the local foreground-background luminance

179     difference of a foreground object in an image). Figure 5C shows example images with varying

180     values for the four image attributes.

181

182     *Core object recognition behavioral paradigm*

183        Core object discrimination is defined as the ability to discriminate between two or more

184     objects in visual images presented under high view uncertainty in the central visual field (~10°),

185     for durations that approximate the typical primate, free-viewing fixation duration (~200 ms)

186     (DiCarlo and Cox, 2007; DiCarlo et al., 2012). As in our previous work (Rajalingham et al.,

187     2015), the behavioral task paradigm consisted of a interleaved set of binary discrimination tasks.

188     Each binary discrimination task is an object discrimination task between a pair of objects (e.g.

189     elephant vs. bear). Each such binary task is balanced in that the test image is equally likely

190     (50%) to be of either of the two objects. On each trial, a test image is presented, followed by a

191     choice screen showing canonical views of the two possible objects (the object that was not

192     displayed in the test image is referred to as the "distractor" object, but note that objects are

193     equally likely to be distractors and targets). Here, 24 objects were tested, which resulted in 276

194     binary object discrimination tasks. To neutralize feature attention, these 276 tasks are randomly

195     interleaved (trial by trial), and the global task is referred to as a basic-level, core object

196     recognition task paradigm.

197

198     *Testing human behavior*

199        All human behavioral data presented here were collected from 1476 human subjects on

200     Amazon Mechanical Turk (MTurk) performing the task paradigm described above. Subjects

201     were instructed to report the identity of the foreground object in each presented image from

202     among the two objects presented on the choice screen (Fig 1B). Because all 276 tasks were

203     interleaved randomly (trial-by-trial), subjects could not deploy feature attentional strategies

204     specific to each object or specific to each binary task to process each test image.

205

206    Figure 1B illustrates the time course of each behavioral trial, for a particular object

207    discrimination task (zebra versus dog). Each trial initiated with a central black point for 500 ms,

208    followed by 100 ms presentation of a test image containing one foreground object presented

209    under high variation in viewing parameters and overlaid on a random background, as described

210    above (see *Visual images* above). Immediately after extinction of the test image, two choice

211    images, each displaying a single object in a canonical view with no background, were shown to

212    the left and right. One of these two objects was always the same as the object that generated the

213    test image (i.e., the correct object choice), and the location of the correct object (left or right) was

214    randomly chosen on each trial. After clicking on one of the choice images, the subject was

215    queued with another fixation point before the next test image appeared. No feedback was given;

216    human subjects were never explicitly trained on the tasks. Under assumptions of typical

217    computer ergonomics, we estimate that images were presented at 6–8° of visual angle at the

218    center of gaze, and the choice object images were presented at ±6–8° of eccentricity along the

219    horizontal meridian.

220

221    We measured human behavior using the online Amazon MTurk platform (see Figure 1C),

222    which enables efficient collection of large-scale psychophysical data from crowd-sourced

223    "human intelligence tasks" (HITs). The reliability of the online MTurk platform has been

224    validated by comparing results obtained from online and in-lab psychophysical experiments

225    (Majaj et al., 2015; Rajalingham et al., 2015). We pooled 927,296 trials from 1472 human

226    subjects to characterize the aggregate human behavior, which we refer to as the "pooled" human

227    (or "archetypal" human). Each human subject performed only a small number of trials (~150) on

228    a subset of the images and binary tasks. All 2400 images were used for behavioral testing, but in

229    some of the HITs, we biased the image selection towards the 240 primary test images (1424±70

230    trials/image on this subsampled set, versus 271±93 trials/image on the remaining images, mean ±

231    SD) to efficiently characterize behavior at image level resolution. Images were randomly drawn

232    such that each human subject was exposed to each image a relatively small number of times

233    (1.5±2.0 trials/image per subject, mean ± SD), in order to mitigate potential alternative

234    behavioral strategies (e.g. "memorization" of images) that could arise from a finite image set.

235    Behavioral signatures at the object-level (B.O1, B.O2, see *Behavioral metrics and signatures*)

236    were measured using all 2400 test images, while image-level behavioral signatures (B.I1n, B.I2n,

237    see *Behavioral metrics and signatures*) were measured using the 240 primary test images. (We

238    observed qualitatively similar results using those metrics on the full 2400 test images, but we

239    here focus on the primary test images as the larger number of trials leads to lower noise levels).

240

241         Five other human subjects were separately recruited on MTurk to each perform a large

242    number of trials on the same images and tasks ($53,097 \pm 15,278$ trials/subject, mean $\pm$ SD).

243    Behavioral data from these five subjects was not included in the characterization of the pooled

244    human described above, but instead aggregated together to characterize a distinct held-out

245    human pool. For the scope of the current work, this held-out human pool—which largely

246    replicated all behavioral signatures of the larger archetypal human (see Figures 2 and 3)—served

247    as an independent validation of our human behavioral measurements.

248

249    *Testing monkey behavior*

250         Five adult male rhesus macaque monkeys (*Macaca mulatta, subjects M, Z, N, P, B*) were

251    tested on the same basic-level, core object recognition task paradigm described above, with

252    minor modification as described below. All procedures were performed in compliance with

253    National Institutes of Health guidelines and the standards of the Massachusetts Institute of

254    Technology Committee on Animal Care and the American Physiological Society. To efficiently

255    characterize monkey behavior, we used a novel home-cage behavioral system developed in our

256    lab (termed MonkeyTurk, see Fig. 1C). This system leveraged a tablet touchscreen (9" Google

257    Nexus or 10.5" Samsung Galaxy Tab S) and used a web application to wirelessly load the task

258    and collect the data (code available at https://github.com/dicarlolab/mkturk). Analogous to the

259    online Amazon Mechanical Turk, which allows for efficient psychophysical assays of a large

260    number (hundreds) of human users in their native environments, MonkeyTurk allowed us to test

261    many monkey subjects simultaneously in their home environment. Each monkey voluntarily

262    initiated trials, and each readily performed the task a few hours each day that the task apparatus

263    was made available to it. At an average rate of ~2,000 trials per day per monkey, we collected a

264    total of 836,117 trials from the five monkey subjects over a period of ~3 months.

265

266         Monkey training is described in detail elsewhere (Rajalingham et al., 2015). Briefly, all

267    monkeys were initially trained on the match-test-image-to-object rule using other images and

268    were also trained on discriminating the particular set of 24 objects tested here using a separate set

269    of training images rendered from these objects, in the same manner as the main testing images.

270    Two of the monkeys subjects (Z and M) were previously trained in the lab setting, and the

271    remaining three subjects were trained using MonkeyTurk directly in their home cages and did

272    not have significant prior lab exposure. Once monkeys reached saturation performance on

273    training images, we began the behavioral testing phase to collect behavior on test images.

274    Monkeys did improve throughout the testing phase, exhibiting an increase in performance

275    between the first and second half of trials of 4%±0.9% (mean ± SEM over five monkey subjects).

276    However, the image-level behavioral signatures obtained from the first and the second halves of

277    trials were highly correlated to each other (B.I1 noise-adjusted correlation of 0.85±0.06, mean ±

278    SEM over five monkey subjects, see *Behavioral metrics and signatures* below), suggesting that

279    monkeys did not significantly alter strategies (e.g. did not "memorize" images) throughout the

280    behavioral testing phase.

281

282        The monkey task paradigm was nearly identical to the human paradigm (see Figure 1B),

283    with the exception that trials were initiated by touching a white "fixation" circle horizontally

284    centered on the bottom third of the screen (to avoid occluding centrally-presented test images

285    with the hand). This triggered a 100ms central presentation of a test image, followed

286    immediately by the presentation of the two choice images (Fig. 1B, location of correct choice

287    randomly assigned on each trial, identical to the human task). Unlike the main human task,

288    monkeys responded by directly touching the screen at the location of one of the two choice

289    images. Touching the choice image corresponding to the object shown in the test image resulted

290    in the delivery of a drop of juice through a tube positioned at mouth height (but not obstructing

291    view), while touching the distractor choice image resulted in a three second timeout. Because

292    gaze direction typically follows the hand during reaching movements, we assumed that the

293    monkeys were looking at the screen during touch interactions with the fixation or choice targets.

294    In both the lab and in the home cage, we maintained total test image size at ~6 degrees of visual

295    angle at the center of gaze, and we took advantage of the retina-like display qualities of the tablet

296    by presenting images pixel matched to the display (256 x 256 pixel image displayed using 256 x

297    256 pixels on the tablet at a distance of 8 inches) to avoid filtering or aliasing effects.

298

299    As with Mechanical Turk testing in humans, MonkeyTurk head-free home-cage testing
300    enables efficient collection of reliable, large-scale psychophysical data but it likely does not yet
301    achieve the level of experimental control that is possible in the head-fixed laboratory setting.
302    However, we note that when subjects were engaged in home-cage testing, they reliably had their
303    mouth on the juice tube and their arm positioned through an armhole. These spatial constraints
304    led to a high level of head position trial-by-trial reproducibility during performance of the task
305    paradigm. Furthermore, when subjects were in this position, they could not see other animals as
306    the behavior box was opaque, and subjects performed the task at a rapid pace 40 trials/minute
307    suggesting that they were not frequently distracted or interrupted. The location of the upcoming
308    test image (but not the location of the object within that test image) was perfectly predictable at
309    the start of each behavioral trial, which likely resulted in a reliable, reproduced gaze direction at
310    the moment that each test image was presented. The relatively short—but natural and high
311    performing (Cadieu et al., 2014)—test image duration (100 ms) ensured that saccadic eye
312    movements were unlike to influence test image performance (as they generally take ~200 ms to
313    initiate in response to the test image, and thus well after the test image has been extinguished).

314

315    *Testing model behavior*

316    We tested a number of different deep convolutional neural network (DCNN) models on
317    the exact same images and tasks as those presented to humans and monkeys. Importantly, our
318    core object recognition task paradigm is closely analogous to the large-scale ImageNet 1000-way
319    object categorization task for which these networks were optimized and thus expected to perform
320    well. We focused on publicly available DCNN model architectures that have proven highly
321    successful with respect to this computer vision benchmark over the past five years: AlexNet
322    (Krizhevsky et al., 2012), NYU (Zeiler and Fergus, 2014), VGG (Simonyan and Zisserman,
323    2014), GoogleNet (Szegedy et al., 2013), Resnet (He et al., 2016), and Inception-v3 (Szegedy et
324    al., 2013). As this is only a subset of possible DCNN models, we refer to these as the $DCNN_{IC}$
325    (to denote ImageNet-Categorization) visual system model sub-family. For each of the publicly
326    available model architectures, we first used ImageNet-categorization-trained model instances,
327    either using publicly available trained model instances or training them to saturation on the 1000-
328    way classification task in-house. Training took several days on 1-2 GPUs.

329

330    We then performed psychophysical experiments on each ImageNet-trained DCNN model

331    to characterize their behavior on the exact same images and tasks as humans and monkeys. We

332    first adapted these ImageNet-trained models to our 24-way object recognition task by re-training

333    the final class probability layer (initially corresponding to the probability output of the 1000-way

334    ImageNet classification task) while holding all other layers fixed. In practice, this was done by

335    extracting features from the penultimate layer of each DCNN$_{IC}$ (i.e. top-most prior to class

336    probability layer), on the same images that were presented to humans and monkeys, and training

337    back-end multi-class logistic regression classifiers to determine the cross-validated output class

338    probability for each image. This procedure is illustrated in Figure 1C. To estimate the hit rate of

339    a given image in a given binary classification task, we renormalized the 24-way class

340    probabilities of that image, considering only the two relevant classes, to sum to one. Object-level

341    and image-level behavioral metrics were computed based on these hit rate estimates (as

342    described in *Behavioral metrics and signatures* below). Importantly, this procedure assumes that

343    the model "retina" layer processes the central 6 degrees of the visual field. It also assumes that

344    linear discriminants ("readouts") of the model's top feature layer are its behavioral output (as

345    intended by the model designers). Manipulating either of these choices (e.g. resizing the input

346    images such that they span only part of the input layer, or building linear discriminates for

347    behavior using a different model feature layer) would result in completely new, testable ANN

348    models that we do not test here.

349

350    From these analyses, we selected the most *human-consistent* DCNN$_{IC}$ architecture

351    (Inception-v3, see *Behavioral consistency* below), fixed that architecture, and then performed

352    post-hoc analyses in which we varied: the input image sampling, the initial parameter settings

353    prior to training, the filter training images, the type of classifiers used to generate the behavior

354    from the model features, and the classifier training images. To examine input image sampling,

355    we re-trained the Inception-v3 architecture on images from ImageNet that were first spatially

356    filtered to match the spatial sampling of the primate retina (i.e. an approximately exponential

357    decrease in cone density away from the fovea) by effectively simulating a fish-eye

358    transformation on each image. These images were at highest resolution at the "fovea" (i.e. center

359    of the image) with gradual decrease in resolution with increasing eccentricity. To examine the

360    analog of "inter-subject variability", we constructed multiple trained model instances

361 ("subjects"), where the architecture and training images were held fixed (Inception-v3 and
362 ImageNet, respectively) but the model filter weights initial condition and order of training
363 images were randomly varied for each model instance. Importantly, this procedure is only one
364 possible choice for simulating inter-subject variability for DCNN models, a choice that is an
365 important open research direction that we do not address here. To examine the effect of model
366 training, we fine-tuned an ImageNet-trained Inception-v3 model on a synthetic image set
367 consisting of ~6.9 million images of 1049 objects (holding out 50,000 images for model
368 validation). These images were generated using the same rendering pipeline as our test images,
369 but the objects were non-overlapping with the 24 test objects presented here. As expected, fine-
370 tuning on synthetic images led to an overall increase in performance of ~5%. We tested the effect
371 of different classifiers to generate model behavior by testing both multi-class logistic regression
372 and support vector machine classifiers. Additionally, we tested the effect of varying the number
373 of training images used to train those classifiers (20 versus 50 images per class).

374

375 *Behavioral metrics and signatures*

376      To characterize the behavior of any visual system, we here introduce four behavioral ($B$)
377 metrics of increasing richness, requiring increasing amounts of data to measure reliably. Each
378 behavioral metric computes a pattern of unbiased behavioral performance, using a sensitivity
379 index: $d' = Z(HitRate) - Z(FalseAlarmRate)$, where Z is the inverse of the cumulative
380 Gaussian distribution. The various metrics differ in the resolution at which hit rates and false
381 alarm rates are computed. Table 1 summarizes the four behavioral metrics, varying the hit-rate
382 resolution (object-level or image-level) and the false-alarm resolution (one-versus-all or one-
383 versus-other). When each metric is applied to the behavioral data of a visual system—biological
384 or artificial—we refer to the result as one behavioral "signature" of that system. Note that we do
385 not consider the signatures obtained here to be the final say on the behavior of these biological or
386 artificial systems—in the terms defined here, new experiments using new objects/images but the
387 same metrics would produce additional behavioral signatures.

388

389      The four behavioral metrics we chose are as follows: First, the one-versus-all object-level
390 performance metric (termed B.O1) estimates the discriminability of each object from all other
391 objects, pooling across all distractor object choices. Since we here tested 24 objects, the resulting

392    B.O1 signature has 24 independent values. Second, the one-versus-other object-level

393    performance metric (termed B.O2) estimates the discriminability of each specific pair of objects,

394    or the pattern of pairwise object confusions. Since we here tested 276 interleaved binary object

395    discrimination tasks, the resulting B.O2 signature has 276 independent values (the off-diagonal

396    elements on one half of the 24x24 symmetric matrix). Third, the one-versus-all image-level

397    performance metric (termed B.I1) estimates the discriminability of each image from all other

398    objects, pooling across all possible distractor choices. Since we here focused on the primary

399    image test set of 240 images (10 per object, see above), the resulting B.I1 signature has 240

400    independent values. Fourth, the one-versus-other image-level performance metric (termed B.I2)

401    estimates the discriminability of each image from each distractor object. Since we here focused

402    on the primary image test set of 240 images (10 per object, see above) with 23 distractors, the

403    resulting B.I2 signature has 5520 independent values.

404

405          Naturally, object-level and image-level behavioral signatures are tightly linked. For

406    example, images of a particularly difficult-to-discriminate object would inherit lower

407    performance values on average as compared to images from a less difficult-to-discriminate

408    object. To isolate the behavioral variance that is specifically driven by image variation and not

409    simply predicted by the objects (and thus already captured by B.O1 and B.O2), we defined

410    normalized image-level behavioral metrics (termed B.I1n, B.I2n) by subtracting the mean

411    performance values over all images of the same object and task. This process is schematically

412    illustrated in Figure 3A. We note that the resulting normalized image-level behavioral signatures

413    capture a significant proportion of the total image-level behavioral variance in our data (e.g.

414    52%, 58% of human B.I1 and B.I2 variance is driven by image variation, independent of object

415    identity). In this study, we use these normalized metrics for image-level behavioral comparisons

416    between models and primates (see Results).

417

418    *Behavioral Consistency*

419          To quantify the similarity between a model visual system and the human visual system

420    with respect to a given behavioral metric, we used a measure called the "*human-consistency*" as

421    previously defined (Johnson et al., 2002). *Human-consistency* ($\tilde{\rho}$) is computed, for each of the

422    four behavioral metrics, as a noise-adjusted correlation of behavioral signatures (DiCarlo and

423    Johnson, 1999). For each visual system, we randomly split all behavioral trials into two equal

424    halves and applied each behavioral metric to each half, resulting in two independent estimates of

425    the system's behavioral signature with respect to that metric. We took the Pearson correlation

426    between these two estimates of the behavioral signature as a measure of the reliability of that

427    behavioral signature given the amount of data collected, i.e. the split-half internal reliability. To

428    estimate the *human-consistency*, we computed the Pearson correlation over all the independent

429    estimates of the behavioral signature from the model (**m**) and the human (**h**), and we then divide

430    that raw Pearson correlation by the geometric mean of the split-half internal reliability of the

431    same behavioral signature measured for each system: $\tilde{\rho}(\boldsymbol{m}, \boldsymbol{h}) = \frac{\rho(\boldsymbol{m},\boldsymbol{h})}{\sqrt{\rho(\boldsymbol{m},\boldsymbol{m})\rho(\boldsymbol{h},\boldsymbol{h})}}.$

432

433        Since all correlations in the numerator and denominator were computed using the same

434    amount of trial data (exactly half of the trial data), we did not need to make use of any prediction

435    formulas (e.g. extrapolation to larger number of trials using Spearman-Brown prediction

436    formula). This procedure was repeated 10 times with different random split-halves of trials. Our

437    rationale for using a reliability-adjusted correlation measure for *human-consistency* was to

438    account for variance in the behavioral signatures that arises from "noise," i.e., variability that is

439    not replicable by the experimental condition (image and task) and thus that no model can be

440    expected to predict (DiCarlo and Johnson, 1999; Johnson et al., 2002). In sum, if the model (m)

441    is a replica of the archetypal human (h), then its expected human-consistency is 1.0, regardless of

442    the finite amount of data that are collected.

443

444    *Characterization of Residuals*

445        In addition to measuring the similarity between the behavioral signatures of primates and

446    models (using *human-consistency* analyses, as described above), we examined the corresponding

447    differences, termed "residual signatures." Each candidate visual system model's residual

448    signature was estimated as the residual of a linear least squares regression of the model's

449    signature on the corresponding human signature and a free intercept parameter. This procedure

450    effectively captures the differences between human and model signatures after accounting for

451    overall performance differences. Residual signatures were estimated on disjoint split-halves of

452    trials, repeating 10 times with random trial permutations. Residuals were computed with respect

453    to the normalized one-versus-all image-level performance metric (B.I1n) as this metric showed a

454    clear difference between $DCNN_{IC}$ models and primates, and the behavioral residual can be
455    interpreted based only the test images (i.e. we can assign a residual per image).

456

457        To examine the extent to which the difference between each model and the archetypal
458    human is reliably shared across different models, we measured the Pearson correlation between
459    the residual signatures of pairs of models. Residual similarity was quantified as the proportion of
460    shared variance, defined as the square of the noise-adjusted correlation between residual
461    signatures (the noise-adjustment was done as defined in equation above). Correlations of residual
462    signatures were always computed across distinct split-halves of data, to avoid introducing
463    spurious correlations from subtracting common noise in the human data. We measured the
464    residual similarity between all pairs of tested models, holding both architecture and optimization
465    procedure fixed (between instances of the ImageNet-categorization trained Inception-v3 model,
466    varying in filter initial conditions), varying the architecture while holding the optimization
467    procedure fixed (between all tested ImageNet-categorization trained DCNN architectures), and
468    holding the architecture fixed while varying the optimization procedure (between ImageNet-
469    categorization trained Inception-v3 and synthetic-categorization fine-tuned Inception-v3
470    models). This analysis addresses not only the reliability of the failure of $DCNN_{IC}$ models to
471    predict human behavior (deviations from humans), but also the relative importance of the
472    characteristics defining similarities within the model sub-family (namely, the architecture and the
473    optimization procedure). We first performed this analysis for residual signatures over the 240
474    primary test images, and subsequently zoomed in on subsets of images that humans found to be
475    particularly difficult. This image selection was made relative to the distribution of image-level
476    performance of held-out human subjects (B.I1 metric from five subjects); difficult images were
477    defined as ones with performance below the 25th percentile of this distribution.

478

479        To examine whether the difference between each model and humans can be explained by
480    simple human-interpretable stimulus attributes, we regressed each $DCNN_{IC}$ model's residual
481    signature on image attributes (object size, eccentricity, pose, and contrast). Briefly, we
482    constructed a design matrix from the image attributes (using individual attributes, or all
483    attributes), and used multiple linear least squares regression to predict the image-level residual
484    signature. The multiple linear regression was tested using two-fold cross-validation over trials.

485    The relative importance of each attribute (or groups of attributes) was quantified using the
486    proportion of explainable variance (i.e. variance remaining after accounting for noise variance)
487    explained from the residual signature.

488

489    *Primate behavior zone*

490         In this work, we are primarily concerned with the behavior of an "archetypal human",
491    rather than the behavior of any given individual human subject. We operationally defined this
492    concept as the common behavior over many humans, obtained by pooling together trials from a
493    large number of individual human subjects and treating this human pool as if it were acquired
494    from a single behaving agent. Due to inter-subject variability, we do not expect any given human
495    or monkey subject to be perfectly consistent with this archetypal human (i.e. we do not expect it
496    to have a *human-consistency* of 1.0). Given current limitations of monkey psychophysics, we are
497    not yet able to measure the behavior of very large number of monkey subjects at high resolution
498    and consequently cannot directly estimate the *human-consistency* of the corresponding
499    "archetypal monkey" to the human pool. Rather, we indirectly estimated this value by first
500    measuring *human-consistency* as a function of number of individual monkey subjects pooled
501    together (n), and extrapolating the *human-consistency* estimate for pools of very large number of
502    subjects (as n approaches infinity). Extrapolations were done using least squares fitting of an
503    exponential function $\tilde{\rho}(n) = a + b \cdot e^{-cn}$ (see Figure 4).

504

505         For each behavioral metric, we defined a "primate zone" as the range of *human-*
506    *consistency* values delimited by estimates $\tilde{\rho}_{M\infty}$ and $\tilde{\rho}_{H\infty}$ as lower and upper bounds respectively.
507    $\tilde{\rho}_{M\infty}$ corresponds to the extrapolated estimate of *human-consistency* of a large (i.e. infinitely
508    many) pool of rhesus macaque monkeys; $\tilde{\rho}_{H\infty}$ is by definition equal to 1.0. Thus, the primate
509    zone defines a range of *human-consistency* values that correspond to models that accurately
510    capture the behavior of the human pool, at least as well as an extrapolation of our monkey
511    sample. In this work, we defined this range of *human-consistency* values as the criterion for
512    success for computational models of primate visual object recognition behavior.

513

514         To make a global statistical inference about whether models sampled from the DCNN$_{IC}$
515    sub-family meet or fall short of this criterion for success, we attempted to reject the hypothesis

516    that, for a given behavioral metric, the *human-consistency* of DCNN$_{IC}$ models is within the

517    primate zone. To test this hypothesis, we estimated the empirical probability that the distribution

518    of *human-consistency* values, estimated over different model instances within this family, could

519    produce *human-consistency* values within the primate zone. Specifically, we estimated a p-value

520    for each behavioral metric using the following procedure: We first estimated an empirical

521    distribution of Fisher-transformed *human-consistency* values for this model family (i.e. over all

522    tested DCNN$_{IC}$ models and over all trial-resampling of each DCNN$_{IC}$ model). From this

523    empirical distribution, we fit a Gaussian kernel density function, optimizing the bandwidth

524    parameter to minimize the mean squared error to the empirical distribution. This kernel density

525    function was evaluated to compute a p-value, by computing the cumulative probability of

526    observing a *human-consistency* value greater than or equal to the criterion of success (i.e. the

527    Fisher transformed $\tilde{\rho}_{M\infty}$ value). This p-value indicates the probability that *human-consistency*

528    values sampled from the observed distribution would fall into the primate zone, with smaller p-

529    values indicating stronger evidence against the hypothesis that the *human-consistency* of DCNN

530    models is within the primate zone.

531

532    **RESULTS**

533

534        In the present work, we systematically compared the basic level core object recognition

535    behavior of primates and state-of-the-art artificial neural network models using a series of

536    behavioral metrics ranging from low to high resolution within a two-alternative forced choice

537    match-to-sample paradigm. The behavior of each visual system, whether biological or artificial,

538    was tested on the same 2400 images (24 objects, 100 images/object) in the same 276 interleaved

539    binary object recognition tasks. Each system's behavior was characterized at multiple resolutions

540    (see *Behavioral metrics and signatures* in Methods) and directly compared to the corresponding

541    behavioral metric applied on the archetypal human (defined as the average behavior of a large

542    pool of human subjects tested; see Methods). The overarching logic of this study was that, if two

543    visual systems are equivalent, they should produce statistically indistinguishable behavioral

544    signatures with respect to these metrics. Specifically, our goal was to compare the behavioral

545    signatures of visual system models with the corresponding behavioral signatures of primates.

546

547      *Object-level behavioral comparison*

548            We first examined the pattern of one-versus-all object-level behavior (termed "B.O1

549      metric") computed across all images and possible distractors. Since we tested 24 objects here, the

550      B.O1 signature was 24 dimensional. Figure 2A shows the B.O1 signatures for the pooled human

551      (pooling n=1472 human subjects), pooled monkey (pooling n=5 monkey subjects), and several

552      DCNN$_{IC}$ models as 24-dimensional vectors using a color scale. Each element of the vector

553      corresponds to the system's discriminability of one object against all others that were tested (i.e.

554      all other 23 objects). The color scales span each signature's full performance range, and warm

555      colors indicate lower discriminability. For example, red indicates that the tested visual system

556      found the object corresponding to that element of the vector to be very challenging to

557      discriminate from other objects (on average over all 23 discrimination tests, and on average over

558      all images). Figure 2B directly compares the B.O1 signatures computed from the behavioral

559      output of two visual system models—a pixel model (top panel) and a DCNN$_{IC}$ model (Inception-

560      v3, bottom panel)—against that of the human B.O1 signature. We observe a tighter

561      correspondence to the human behavioral signature for the DCNN$_{IC}$ model visual system than for

562      the baseline pixel model visual system. We quantified that similarity using a noise-adjusted

563      correlation between each pair of B.O1 signatures (termed *human-consistency,* following

564      (Johnson et al., 2002)); the noise adjustment means that a visual system that is identical to the

565      human pool will have an expected *human-consistency* score of 1.0, even if it has irreducible trial-

566      by-trial stochasticity; see Methods). Figure 2C shows the B.O1 *human-consistency* for each of

567      the tested model visual systems. We additionally tested the behavior of a held-out pool of five

568      human subjects (black dot) and a pool of five macaque monkey subjects (gray dot), and we

569      observed that both yielded B.O1 signatures that were highly human-consistent (*human-*

570      *consistency* $\tilde{\rho}$ = 0.90, 0.97 for monkey pool and held-out human pool, respectively). We defined

571      a range of *human-consistency* values, termed the "primate zone" (shaded gray area), delimited by

572      extrapolated *human-consistency* estimates of large pools of macaques (see Methods, Figure 4).

573      We found that the baseline pixel visual system model and the low-level V1 visual system model

574      were not within this zone ($\tilde{\rho}$ = 0.40, 0.67 for pixels and V1 models, respectively), while all tested

575      DCNN$_{IC}$ visual system models were either within or very close to this zone. Indeed, we could not

576      reject the hypothesis that DCNN$_{IC}$ models are primate-like (p = 0.54, exact test, see Methods).

577

578    Next, we compared the behavior of the visual systems at a slightly higher level of

579    resolution. Specifically, instead of pooling over all discrimination tasks for each object, we

580    computed the mean discriminability of each of the 276 pairwise discrimination tasks (still

581    pooling over images within each of those tasks). This yielded a symmetric matrix that is referred

582    to here as the B.O2 signature. Figure 2D shows the B.O2 signatures of the pooled human, pooled

583    monkey, and several DCNN$_{IC}$ visual system models as 24x24 symmetric matrices. Each bin $(i,j)$

584    corresponds to the system's discriminability of objects $i$ and $j$, where warmer colors indicate

585    lower performance; color scales are not shown but span each signature's full range. We observed

586    strong qualitative similarities between the pairwise object confusion patterns of all of the high

587    level visual systems (e.g. camel and dog are often confused with each other by all three systems).

588    This similarity is quantified in Figure 2E, which shows the *human-consistency* of all examined

589    visual system models with respect to this metric. Similar to the B.O1 metric, we observed that

590    both a pool of macaque monkeys and a held-out pool of humans are highly *human-consistent*

591    with respect to this metric ($\tilde{\rho}$ = 0.77, 0.94 for monkeys, humans respectively). Also similar to the

592    B.O1 metric, we found that all DCNN$_{IC}$ visual system models are highly *human-consistent* ($\tilde{\rho}$ >

593    0.8) while the baseline pixel visual system model and the low-level V1 visual system model were

594    not ($\tilde{\rho}$ = 0.41, 0.57 for pixels, V1 models respectively). Indeed, all DCNN$_{IC}$ visual system

595    models are within the defined "primate zone" of *human-consistency*, and we could not falsify the

596    hypothesis that DCNN$_{IC}$ models are primate-like (p = 0.99, exact test).

597

598    Taken together, humans, monkeys, and current DCNN$_{IC}$ models all share similar patterns

599    of object-level behavioral performances (B.O1 and B.O2 signatures) that are not shared with

600    lower-level visual representations (pixels and V1). However, object-level performance patterns

601    do not capture the fact that some images of an object are more challenging than other images of

602    the same object because of interactions of the variation in the object's pose and position with the

603    object's class. To overcome this limitation, we next examined the patterns of behavior at the

604    resolution of individual images on a subsampled set of images where we specifically obtained a

605    large number of behavioral trials to accurately estimate behavioral performance on each image.

606    Note that, from the point of view of the subjects, the behavioral tasks are identical to those

607    already described. We simply aimed to measure and compare their patterns of performance at

608    much higher resolution.

609

610    *Image-level behavioral comparison*

611         To isolate purely image-level behavioral variance, i.e. variance that is not predicted by

612    the object and thus already captured by the B.O1 signature, we computed the normalized image-

613    level signature. This normalization procedure is schematically illustrated in Figure 3A which

614    shows that the one-versus-all image-level signature (240-dimensional, 10 images/object) is used

615    to obtain the normalized one-versus-all image-level signature (termed B.I1n, see *Behavioral*

616    *metrics and signatures*). Figure 3B shows the B.I1n signatures for the pooled human, pooled

617    monkey, and several $DCNN_{IC}$ models as 240 dimensional vectors. Each bin's color corresponds

618    to the discriminability of a single image against all distractor options (after subtraction of object-

619    level discriminability, see Figure 3A), where warmer colors indicate lower values; color scales

620    are not shown but span each signature's full range. Figure 3D shows the *human-consistency* with

621    respect to the B.I1n signature for all tested models. Unlike with object-level behavioral metrics,

622    we now observe a divergence between $DCNN_{IC}$ models and primates. Both the monkey pool and

623    the held-out human pool remain highly *human-consistent* ($\tilde{\rho} = 0.77$, $0.96$ for monkeys, humans

624    respectively), but all $DCNN_{IC}$ models were significantly less *human-consistent* (Inception-

625    v3: $\tilde{\rho} = 0.62$) and well outside of the defined "primate zone" of B.I1n *human-consistency*.

626    Indeed, the hypothesis that the *human-consistency* of $DCNN_{IC}$ models is within the primate zone

627    is strongly rejected ($p = 6.16e-8$, exact test, see Methods).

628

629         We can zoom in further by examining not only the overall performance for a given image

630    but also the object confusions for each image, i.e. the additional behavioral variation that is due

631    not only to the test image but to the interaction of that test image with the alternative (incorrect)

632    object choice that is provided after the test image (see Fig. 1B). This is the highest level of

633    behavioral accuracy resolution that our task design allows. In raw form, it corresponds to one-

634    versus-other image-level confusion matrix, where the size of that matrix is the total number of

635    images by the total number of objects (here, 240x24). Each bin ($i,j$) corresponds to the behavioral

636    discriminability of a single image $i$ against distractor object $j$. Again, we isolate variance that is

637    not predicted by object-level performance by subtracting the average performance on this binary

638    task (mean over all images) to convert the raw matrix B.I2 above into the normalized matrix,

639    referred to as B.I2n. Figure 3D shows the B.I2n signatures as 240x24 matrices for the pooled

640 human, pooled monkey and top $DCNN_{IC}$ visual system models. Color scales are not shown but

641 span each signature's full range; warmer colors correspond to images with lower performance in

642 a given binary task, relative to all images of that object in the same task. Figure 3E shows the

643 *human-consistency* with respect to the B.I2n metric for all tested visual system models.

644 Extending our observations using B.I1n, we observe a similar divergence between primates and

645 $DCNN_{IC}$ visual system models on the matrix pattern of image-by-distractor difficulties (B.I2n).

646 Specifically, both the monkey pool and held-out human pool remain highly *human-consistent*

647 ($\tilde{\rho} = 0.75$, $0.77$ for monkeys, humans respectively), while all tested $DCNN_{IC}$ models are

648 significantly less *human-consistent* (Inception-v3: $\tilde{\rho} = 0.53$) falling well outside of the defined

649 "primate zone" of B.I2n *human-consistency* values. Once again, the hypothesis that the *human-*

650 *consistency* of $DCNN_{IC}$ models is within the primate zone is strongly rejected (p = 3.17e-18,

651 exact test, see Methods).

652

653 *Natural subject-to-subject variation*

654 For each behavioral metric (B.O1, BO2, B.I1n, BI2n), we defined a "primate zone" as the

655 range of consistency values delimited by *human-consistency* estimates $\tilde{\rho}_{M\infty}$ and $\tilde{\rho}_{H\infty}$ as lower

656 and upper bounds respectively. $\tilde{\rho}_{M\infty}$ corresponds to the extrapolated estimate of the *human-*

657 *consistency* of a large (i.e. infinitely many subjects) pool of rhesus macaque monkeys. Thus, the

658 fact that a particular tested visual system model falls outside of the primate zone can be

659 interpreted as a failure of that visual system model to accurately predict the behavior of the

660 archetypal human at least as well as the archetypal monkey.

661

662 However, from the above analyses, it is not yet clear whether a visual system model that

663 fails to predict the archetypal human might nonetheless accurately correspond to one or more

664 individual human subjects found within the natural variation of the human population. Given the

665 difficulty of measuring individual subject behavior at the resolution of single images for large

666 numbers of human and monkey subjects, we could not yet directly test this hypothesis. Instead,

667 we examined it indirectly by asking whether an archetypal model—that is a pool that includes an

668 increasing number of model "subjects"—would approach the human pool. We simulated model

669 inter-subject variability by retraining a fixed DCNN architecture with a fixed training image set

670 with random variation in the initial conditions and order of training images. This procedure

671 results in models that can still perform the task but with slightly different learned weight values.

672 We note that this procedure is only one possible choice of generating inter-subject variability

673 within each visual system model type, a choice that is an important open research direction that

674 we do not address here. From this procedure, we constructed multiple trained model instances

675 ("subjects") for a fixed DCNN architecture, and asked whether an increasingly large pool of

676 model "subjects" better captures the behavior of the human pool, at least as well as a monkey

677 pool. This post-hoc analysis was conducted for the most *human-consistent* DCNN architecture

678 (Inception-v3).

679

680  Figure 4A shows, for each of the four behavioral metrics, the measured *human-*

681 *consistency* of subject pools of varying size (number of subjects *n*) of rhesus macaque monkeys

682 (black) and ImageNet-trained Inception-v3 models (blue). The *human-consistency* increases with

683 growing number of subjects for both visual systems across all behavioral metrics. To estimate

684 the expected *human-consistency* for a pool of infinitely many monkey or model subjects, we fit

685 an exponential function mapping *n* to the mean *human-consistency* values and obtained a

686 parameter estimate for the asymptotic value (see Methods). We note that estimated asymptotic

687 values are not significantly beyond the range of the measured data—the *human-consistency* of a

688 pool of five monkey subjects reaches within 97% of the *human-consistency* of an estimated

689 infinite pool of monkeys for all metrics—giving credence to the extrapolated *human-consistency*

690 values. This analysis suggests that under this model of inter-subject variability, a pool of

691 Inception-v3 subjects accurately capture archetypal human behavior at the resolution of objects

692 (B.O1, B.O2) by our primate zone criterion (see Figure 4A, first two panels). In contrast, even a

693 large pool of Inception-v3 subjects still fails at its final asymptote to accurately capture human

694 behavior at the image-level (B.I1n, B.I2n) (Figure 4A, last two panels).

695

696 *Modification of visual system models to try to rescue their human-consistency*

697  Next, we wondered if some relatively simple changes to the $DCNN_{IC}$ visual system

698 models tested here could bring them into better correspondence with the primate visual system

699 behavior (with respect to B.I1n and B.I2n metrics). Specifically, we considered and tested the

700 following modifications to the most *human-consistent* $DCNN_{IC}$ model visual system (Inception-

701 v3): we (1) changed the input to the model to be more primate-like in its retinal sampling

702   (Inception-v3 + retina-like), (2) changed the transformation (aka "decoder") from the internal
703   model feature representation into the behavioral output by augmenting the number of decoder
704   training images or changing the decoder type (Inception-v3 + SVM, Inception-v3 +
705   classifier_train), and (3) modified all of the internal filter weights of the model (aka "fine
706   tuning") by augmenting its ImageNet training with additional images drawn from the same
707   distribution as our test images (Inception-v3 + synthetic-fine-tune). While some of these
708   modifications (e.g. fine-tuning on synthetic images and increasing the number of classifier
709   training images) had the expected effect of increasing mean overall performance (not shown, see
710   Methods), we found that none of these modifications led to a significant improvement in its
711   *human-consistency* on the behavioral metrics (Figure 4B). Thus, the failure of current $DCNN_{IC}$
712   models to accurately capture the image-level signatures of primates cannot be rescued by simple
713   modifications on a fixed architecture.

714

715   *Looking for clues: Image-level comparisons of models and primates*

716       Taken together, Figures 2, 3 and 4 suggest that current $DCNN_{IC}$ visual system models fail
717   to accurately capture the image-level signatures of humans and monkeys. To further examine this
718   failure in the hopes of providing clues for model improvement, we examined the image-level
719   residual signatures of all the visual system models, relative to the pooled human. For each model,
720   we computed its residual signature as the difference (positive or negative) of a linear least
721   squares regression of the model signature on the corresponding human signature. For this
722   analysis, we focused on the B.I1n metric as it showed a clear divergence of $DCNN_{IC}$ models and
723   primates, and the behavioral residual can be interpreted based only on the test images (whereas
724   B.I2n depends on the interaction between test images and distractor choice).

725

726       We first asked to what extent the residual signatures are shared between different visual
727   system models. Figure 5A shows the similarity between the residual signatures of all pairs of
728   models; the color of bin (*i,j*) indicates the proportion of explainable variance that is shared
729   between the residual signatures of visual systems *i* and *j*. For ease of interpretation, we ordered
730   visual system models based on their architecture and optimization procedure and partitioned this
731   matrix into four distinct regions. Each region compares the residuals of a "source" model group
732   with fixed architecture and optimization procedure (five Inception-v3 models optimized for

733    categorization on ImageNet, varying only in initial conditions and training image order) to a

734    "target" model group. The target groups of models for each of the four regions are: 1) the pooled

735    monkey, 2) other $DCNN_{IC}$ models from the source group, 3) $DCNN_{IC}$ models that differ in

736    architecture but share the optimization procedure of the source group models and 4) $DCNN_{IC}$

737    models that differ slightly using an augmented optimization procedure but share the architecture

738    of the source group models. Figure 5B shows the mean (±SD) variance shared in the residuals

739    averaged within these four regions for all images (black dots), as well as for images that humans

740    found to be particularly difficult (gray dots, selected based on held-out human data, see

741    Methods). First, consistent with the results shown in Figure 3, we note that the residual

742    signatures of this particular $DCNN_{IC}$ model are not well shared with the pooled monkey ($r^2=0.39$

743    in region 1), and this phenomenon is more pronounced for the images that humans found most

744    difficult ($r^2=0.17$ in region 1). However, this relatively low correlation between model and

745    primate residuals is not indicative of spurious model residuals, as the model residual signatures

746    were highly reliable between different instances of this fixed $DCNN_{IC}$ model, across random

747    training initializations (region 2: $r^2=0.79$, 0.77 for all and most difficult images, respectively).

748    Interestingly, residual signatures were still largely shared with other $DCNN_{IC}$ models with vastly

749    different architectures (region 3: $r^2=0.70$, 0.65 for all and most difficult images, respectively).

750    However, residual signatures were more strongly altered when the visual training diet of the

751    same architecture was altered (region 4: $r^2=0.57$, 0.46 for all and most difficult images

752    respectively, cf. region 3). Taken together, these results indicate that the images where $DCNN_{IC}$

753    visual system models diverged from humans (and monkeys) were not spurious but were rather

754    highly reliable across different model architectures, demonstrating that current $DCNN_{IC}$ models

755    systematically and similarly diverge from primates.

756

757          To look for clues for model improvement, we asked what, if any, characteristics of

758    images might account for this divergence of models and primates. We regressed the residual

759    signatures of $DCNN_{IC}$ models on four different image attributes (corresponding to the size,

760    eccentricity, pose, and contrast of the object in each image). We used multiple linear regressions

761    to predict the model residual signatures from all of these image attributes, and also considered

762    each attribute individually using simple linear regressions. Figure 6A shows example images

763    (sampled from the full set of 2400 images) with increasing attribute value for each of these four

764     image attributes. While the $DCNN_{IC}$ models were not directly optimized to display primate-like

765     performance dependence on such attributes, we observed that the Inception-v3 visual system

766     model nonetheless exhibited qualitatively similar performance dependencies as primates (see

767     Figure 6B). For example, humans (black), monkeys (gray) and the Inception-v3 model (blue) all

768     performed better, on average, for images in which the object is in the center of gaze (low

769     eccentricity) and large in size. Furthermore, all three systems performed better, on average, for

770     images when the pose of the object was closer to the canonical pose (see Figure 6B); this

771     sensitivity to object pose manifested itself as a non-linear dependence due to the fact that all

772     tested objects exhibited symmetry in at least one axis. The similarity of the patterns in Figure 6B

773     between primates and the $DCNN_{IC}$ visual system models is not perfect but is striking,

774     particularly in light of the fact that these models were not optimized to produce these patterns.

775     However, this similarity is analogous to the similarity in the B.O1 and B.O2 metrics in that it

776     only holds on average over many images. Looking more closely at the image-by-image

777     comparison, we again found that the $DCNN_{IC}$ models failed to capture a large portion of the

778     image-by-image variation (Figure 3). In particular, Figure 6C shows the proportion of variance

779     explained by specific image attributes for the residual signatures of monkeys (black) and

780     $DCNN_{IC}$ models (blue). We found that, taken together, all four of these image attributes

781     explained only ~10% of the variance in $DCNN_{IC}$ residual signatures, and each individual

782     attribute could explain at most a small amount of residual variance (<5% of the explainable

783     variance). In sum, these analyses show that some behavioral effects that might provide intuitive

784     clues to modify the $DCNN_{IC}$ models are already in place in those models (e.g. a dependence on

785     eccentricity). But the quantitative image-by-image analyses of the remaining unexplained

786     variance (Figure 6C) argue that the $DCNN_{IC}$ visual system models' failure to capture primate

787     image-level signatures cannot be further accounted for by these simple image attributes and

788     likely stem from other factors.

789

790     **DISCUSSION**

791

792        The current work was motivated by the broad scientific goal of discovering models that

793     quantitatively explain the neuronal mechanisms underlying primate invariant object recognition

794     behavior. To this end, previous work had shown that specific artificial neural network models

795   (ANNs), drawn from a large family of deep convolutional neural networks (DCNNs) and

796   optimized to achieve high levels of object categorization performance on large-scale image-sets,

797   capture a large fraction of the variance in primate visual recognition behaviors (Rajalingham et

798   al., 2015; Jozwik et al., 2016; Kheradpisheh et al., 2016; Kubilius et al., 2016; Peterson et al.,

799   2016; Wallis et al., 2017), and the internal hidden neurons of those same models also predict a

800   large fraction of the image-driven response variance of brain activity at multiple stages of the

801   primate ventral visual stream (Yamins et al., 2013; Cadieu et al., 2014; Khaligh-Razavi and

802   Kriegeskorte, 2014; Yamins et al., 2014; Güçlü and van Gerven, 2015; Cichy et al., 2016; Hong

803   et al., 2016; Seibert et al., 2016; Cadena et al., 2017; Wen et al., 2017). For clarity, we here

804   referred to this sub-family of models as $DCNN_{IC}$ (to denote ImageNet-Categorization training),

805   so as to distinguish them from all possible models in the DCNN family, and more broadly, from

806   the super-family of all ANNs. In this work, we directly compared leading $DCNN_{IC}$ models to

807   primates (humans and monkeys) with respect to their behavioral signatures at both object and

808   image level resolution in the domain of core object recognition. In order to do so, we measured

809   and characterized primate behavior at larger scale and higher resolution than previously possible.

810   We first replicate prior work (Rajalingham et al., 2015) showing that, at the object level,

811   $DCNN_{IC}$ models produce statistically indistinguishable behavior from primates, and we extend

812   that work by showing that these models also match the *average* primate sensitivities to object

813   contrast, eccentricity, size, and pose, a noteworthy similarity in light of the fact that these models

814   were not optimized to produce these performance patterns. However, our primary novel result is

815   that, examining behavior at the higher resolution of individual images, all leading $DCNN_{IC}$

816   models failed to replicate the image-level behavioral signatures of primates. An important related

817   claim is that rhesus monkeys are more consistent with the archetypal human than any of the

818   tested $DCNN_{IC}$ models (at the image-level).

819

820        While it had previously been shown that $DCNN_{IC}$ models can diverge from human

821   behavior on specifically chosen adversarial images (Szegedy et al., 2013), a strength of our work

822   is that we did not optimize images to induce failure but instead randomly sampled the image

823   generative parameter space broadly. As such, our results highlight a *general*, rather than

824   adversarial-induced, failure of $DCNN_{IC}$ models to fully capture the neural mechanisms

825   underlying primate core object recognition behavior. Furthermore, we showed that this failure of

826  current $DCNN_{IC}$ models cannot be explained by simple image attributes and cannot be rescued

827  by simple model modifications (input image sampling, model training, and classifier variations).

828  Taken together, these results suggest that new ANN models are needed to more precisely capture

829  the neural mechanisms underlying primate object vision.

830

831  With regards to new ANN models, we can attempt to make prospective inferences about

832  future possible $DCNN_{IC}$ models from the data presented here. Based on the observed distribution

833  of image-level *human-consistency* values for the $DCNN_{IC}$ models tested here, we infer that yet

834  untested model instances sampled identically (i.e. from the $DCNN_{IC}$ model sub-family) are very

835  likely to have similarly inadequate image-level *human-consistency*. While we cannot rule out the

836  possibility that at least one model instance within the $DCNN_{IC}$ sub-family would fully match the

837  image-level behavioral signatures, the probability of sampling such a model is vanishingly small

838  ($p<10^{-17}$ for B.I2n *human-consistency*, estimated using exact test using Gaussian kernel density

839  estimation, see Methods, Results). An important caveat of this inference is that we may have a

840  biased estimate of the *human-consistency* distribution of this model sub-family, as we did not

841  exhaustively sample the sub-family. In particular, if the model sampling process is non-

842  stationary over time (e.g. increases in computational power over time allows larger models to be

843  successfully trained), the *human-consistency* of new (i.e. yet to be sampled) models may lie

844  outside the currently estimated distribution. Consistent with the latter, we observed that current

845  $DCNN_{IC}$ cluster into two distinct "generations" separated in time (before/after the year 2015; e.g.

846  Inception-v3 improves over AlexNet though both lie outside the primate zone in Figure 3). Thus,

847  following this trend, it is possible that the evolution of "next-generation" models within the

848  $DCNN_{IC}$ sub-family could meet our criteria for successful matching primate-like behavior.

849

850  Alternatively, it is possible—and we think likely—that future $DCNN_{IC}$ models will also

851  fail to capture primate-like image-level behavior, suggesting that either the architectural

852  limitations (e.g. convolutional, feed-forward) and/or the optimization procedure (including the

853  diet of visual images) that define this model sub-family are fundamentally limiting. Thus, ANN

854  model sub-families utilizing different architectures (e.g. recurrent neural networks) and/or

855  optimized for different behavioral goals (e.g. loss functions other than object classification

856  performance, and/or images other than category-labeled ImageNet images) may be necessary to

857    accurately capture primate behavior. To this end, we propose that testing even individual
858    changes to the $DCNN_{IC}$ models—each creating a new ANN model sub-family—may be the best
859    way forward, because $DCNN_{IC}$ models currently offer the best explanations (in a predictive
860    sense) of both the behavioral and neural phenomena of core object recognition.

861

862           To reach that goal of finding a new ANN model sub-family that is a better mechanistic
863    model of the primate ventral visual stream, we propose that even larger-scale, high-resolution
864    behavioral measurements, such as expanded versions of the patterns of image-level performance
865    presented here, could serve as a useful top-down optimization guides. Not only do these high-
866    resolution behavioral signatures have the statistical power to reject the currently leading ANN
867    models, but they can also be efficiently collected at very large scale, in contrast to other guide
868    data (e.g. large-scale neuronal measurements). Indeed, current technological tools for high-
869    throughput psychophysics in humans and monkeys (e.g. Amazon Mechanical Turk for humans,
870    Monkey Turk for rhesus monkeys) enable time- and cost-efficient collection of large-scale
871    behavioral datasets, such as the ~1 million behavioral trials obtained for the current work. These
872    systems trade off an increase in efficiency with a decrease in experimental control. For example,
873    we did not impose experimental constraints on subjects' acuity and we can only infer likely head
874    and gaze position. Previous work has shown that patterns of behavioral performance on object
875    recognition tasks from in-lab and online subjects were equally reliable and virtually identical
876    (Majaj et al., 2015), but it is not yet clear to what extent this holds at the resolution of individual
877    images, as one might expect that variance in performance across images is more sensitive to
878    precise head and gaze location. For this reason, we here refrain from making strong inferences
879    from small behavioral differences, such as the small difference between humans and monkeys.
880    Nevertheless, we argue that this sacrifice in exact experimental control while retaining sufficient
881    power for model comparison is a good tradeoff for efficiently collecting large behavioral datasets
882    toward the goal of constraining future models of the primate ventral visual stream.

883

**REFERENCES**

Cadena SA, Denfield GH, Walker EY, Gatys LA, Tolias AS, Bethge M, Ecker AS (2017) Deep convolutional models improve predictions of macaque V1 responses to natural images. bioRxiv:201764.

Cadieu CF, Hong H, Yamins DL, Pinto N, Ardila D, Solomon EA, Majaj NJ, DiCarlo JJ (2014) Deep neural networks rival the representation of primate IT cortex for core visual object recognition. PLoS computational biology 10:e1003963.

Cichy RM, Khosla A, Pantazis D, Torralba A, Oliva A (2016) Comparison of deep neural networks to spatio-temporal cortical dynamics of human visual object recognition reveals hierarchical correspondence. Scientific reports 6:27755.

DiCarlo JJ, Johnson KO (1999) Velocity invariance of receptive field structure in somatosensory cortical area 3b of the alert monkey. Journal of Neuroscience 19:401-419.

DiCarlo JJ, Cox DD (2007) Untangling invariant object recognition. Trends in cognitive sciences 11:333-341.

DiCarlo JJ, Zoccolan D, Rust NC (2012) How does the brain solve visual object recognition? Neuron 73:415-434.

Dodge S, Karam L (2017) A Study and Comparison of Human and Deep Learning Recognition Performance Under Visual Distortions. arXiv preprint arXiv:170502498.

Geirhos R, Janssen DH, Schütt HH, Rauber J, Bethge M, Wichmann FA (2017) Comparing deep neural networks against humans: object recognition when the signal gets weaker. arXiv preprint arXiv:170606969.

Goodfellow IJ, Shlens J, Szegedy C (2014) Explaining and harnessing adversarial examples. arXiv preprint arXiv:14126572.

Güçlü U, van Gerven MA (2015) Deep neural networks reveal a gradient in the complexity of neural representations across the ventral stream. Journal of Neuroscience 35:10005-10014.

He K, Zhang X, Ren S, Sun J (2016) Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition, pp 770-778.

Hong H, Yamins DL, Majaj NJ, DiCarlo JJ (2016) Explicit information for category-orthogonal object properties increases along the ventral stream. Nature neuroscience 19:613.

Hosseini H, Xiao B, Jaiswal M, Poovendran R (2017) On the Limitation of Convolutional Neural Networks in Recognizing Negative Images. human performance 4:6.

Huynh DQ (2009) Metrics for 3D rotations: Comparison and analysis. Journal of Mathematical Imaging and Vision 35:155-164.

Johnson KO, Hsiao SS, Yoshioka T (2002) Neural coding and the basic law of psychophysics. The Neuroscientist 8:111-121.

Jozwik KM, Kriegeskorte N, Mur M (2016) Visual features as stepping stones toward semantics: Explaining object similarity in IT and perception with non-negative least squares. Neuropsychologia 83:201-226.

Khaligh-Razavi S-M, Kriegeskorte N (2014) Deep supervised, but not unsupervised, models may explain IT cortical representation. PLoS computational biology 10:e1003915.

Kheradpisheh SR, Ghodrati M, Ganjtabesh M, Masquelier T (2016) Deep networks can resemble human feed-forward vision in invariant object recognition. Scientific reports 6:32672.

Krizhevsky A, Sutskever I, Hinton GE (2012) Imagenet classification with deep convolutional neural networks. In: Advances in neural information processing systems, pp 1097-1105.

Kubilius J, Bracci S, de Beeck HPO (2016) Deep neural networks as a computational model for human shape sensitivity. PLoS computational biology 12:e1004896.

LeCun Y, Bengio Y, Hinton G (2015) Deep learning. Nature 521:436-444.

932  Majaj NJ, Hong H, Solomon EA, DiCarlo JJ (2015) Simple Learned Weighted Sums of Inferior
933  Temporal Neuronal Firing Rates Accurately Predict Human Core Object Recognition
934  Performance. The Journal of Neuroscience 35:13402-13418.
935  Nguyen A, Yosinski J, Clune J (2015) Deep neural networks are easily fooled: High confidence
936  predictions for unrecognizable images. In: Proceedings of the IEEE Conference on Computer
937  Vision and Pattern Recognition, pp 427-436.
938  Peterson JC, Abbott JT, Griffiths TL (2016) Adapting deep network features to capture
939  psychological representations. arXiv preprint arXiv:160802164.
940  Pinto N, Cox DD, DiCarlo JJ (2008) Why is real-world visual object recognition hard? PLoS
941  computational biology 4:e27.
942  Rajalingham R, Schmidt K, DiCarlo JJ (2015) Comparison of Object Recognition Behavior in
943  Human and Monkey. The Journal of Neuroscience 35:12127-12136.
944  RichardWebster B, Anthony SE, Scheirer WJ (2016) PsyPhy: A Psychophysics Driven
945  Evaluation Framework for Visual Recognition. arXiv preprint arXiv:161106448.
946  Rosch E, Mervis CB, Gray WD, Johnson DM, Boyes-Braem P (1976) Basic objects in natural
947  categories. Cognitive psychology 8:382-439.
948  Seibert D, Yamins DL, Ardila D, Hong H, DiCarlo JJ, Gardner JL (2016) A performance-
949  optimized model of neural responses across the ventral visual stream. bioRxiv:036475.
950  Simonyan K, Zisserman A (2014) Very deep convolutional networks for large-scale image
951  recognition. arXiv preprint arXiv:14091556.
952  Szegedy C, Zaremba W, Sutskever I, Bruna J, Erhan D, Goodfellow I, Fergus R (2013)
953  Intriguing properties of neural networks. arXiv preprint arXiv:13126199.
954  Ullman S, Humphreys GW (1996) High-level vision: Object recognition and visual cognition: MIT
955  press Cambridge, MA.
956  Wallis TS, Funke CM, Ecker AS, Gatys LA, Wichmann FA, Bethge M (2017) A parametric
957  texture model based on deep convolutional features closely matches texture appearance for
958  humans. Journal of vision 17:5-5.
959  Wen H, Shi J, Zhang Y, Lu K-H, Cao J, Liu Z (2017) Neural encoding and decoding with deep
960  learning for dynamic natural vision. Cerebral Cortex:1-25.
961  Yamins DL, DiCarlo JJ (2016) Using goal-driven deep learning models to understand sensory
962  cortex. Nature neuroscience 19:356-365.
963  Yamins DL, Hong H, Cadieu C, DiCarlo JJ (2013) Hierarchical modular optimization of
964  convolutional networks achieves representations similar to macaque IT and human ventral
965  stream. In: Advances in neural information processing systems, pp 3093-3101.
966  Yamins DL, Hong H, Cadieu CF, Solomon EA, Seibert D, DiCarlo JJ (2014) Performance-
967  optimized hierarchical models predict neural responses in higher visual cortex. Proceedings of
968  the National Academy of Sciences:201403112.
969  Zeiler MD, Fergus R (2014) Visualizing and understanding convolutional networks. In: Computer
970  Vision–ECCV 2014, pp 818-833: Springer.

971
972

973 **TABLES**

974 Table 1

| Behavioral Metric | Hit Rate | False Alarm Rate |
|---|---|---|
| One-versus all object-level performance (B.O1) ($N_{objects}$ x 1) $$O_1(i) = Z\big(HR(i)\big) - Z\big(FAR(i)\big),$$ $$i = 1,2,\ldots,N_{objects}$$ | Proportion of trials when images of object $i$ were correctly labeled as object $i$. | Proportion of trials when any image was incorrectly labeled as object $i$. |
| One-versus-other object-level performance B.O2 ($N_{objects}$ x $N_{objects}$) $$O_2(i,j) = Z\big(HR(i,j)\big) - Z\big(FAR(i,j)\big),$$ $$i = 1,2,\ldots,N_{objects}$$ $$j = 1,2,\ldots,N_{objects}$$ | Proportion of trials when images of object $i$ were correctly labeled as $i$, when presented against distractor object $j$. | Proportion of trials when images of object $j$ were incorrectly labeled as object $i$ |
| One-versus-all image-level performance B.I1 ($N_{images}$ x 1) $$I_1(ii) = Z\big(HR(ii)\big) - Z\big(FAR(ii)\big),$$ $$ii = 1,2,\ldots,N_{images}$$ | Proportion of trials when image $ii$ was correctly classified as object $i$. | Proportion of trials when any image was incorrectly labeled as object $i$. |
| One-versus-other image-level performance B.I2 ($N_{images}$ x $N_{objects}$) $$I_2(ii,j) = Z\big(HR(ii,j)\big) - Z\big(FAR(ii,j)\big),$$ $$ii = 1,2,\ldots,N_{images}$$ $$j = 1,2,\ldots,N_{objects}$$ | Proportion of trials when image $ii$ was correctly classified as object $i$, when presented against distractor object $j$. | Proportion of trials when images of object $j$ were incorrectly labeled as object $i$ |

975

976 **Table 1: Definition of behavioral performance metrics.** The first column provides the name,

977 abbreviation, dimensions, and equations for each of the raw performance metrics. The next two

978 columns provide the definitions for computing the hit rate (HR) and false alarm rate (FAR)

979 respectively.

980

**FIGURE LEGENDS**

**Figure 1. Images and behavioral task. (A)** Two (out of 100) example images for each of the 24 basic-level objects. To enforce true invariant object recognition behavior, we generated naturalistic synthetic images, each with one foreground object, by rendering a 3D model of each object with randomly chosen viewing parameters and placing that foreground object view onto a randomly chosen, natural image background. **(B)** Time course of example behavioral trial (zebra versus dog) for human psychophysics. Each trial initiated with a central fixation point for 500 ms, followed by 100 ms presentation of a square test image (spanning 6-8° of visual angle). After extinction of the test image, two choice images were shown to the left and right. Human participants were allowed to freely view the response images for up to 1000 ms and responded by clicking on one of the choice images; no feedback was given. To neutralize top-down feature attention, all 276 binary object discrimination tasks were randomly interleaved on a trial-by-trial basis. The monkey task paradigm was nearly identical to the human paradigm, with the exception that trials were initiated by touching a fixation circle horizontally centered on the bottom third of the screen, and successful trials were rewarded with juice while incorrect choices resulted in timeouts of 1–2.5s. **(C)** Large-scale and high-throughput psychophysics in humans (top left), monkeys (top right), and models (bottom). Human behavior was measured using the online Amazon MTurk platform, which enabled the rapid collection ~1 million behavioral trials from 1472 human subjects. Monkey behavior was measured using a novel custom home-cage behavioral system (MonkeyTurk), which leveraged a web-based behavioral task running on a tablet to test many monkey subjects simultaneously in their home environment. Deep convolutional neural network models were tested on the same images and tasks as those presented to humans and monkeys by extracting features from the penultimate layer of each visual system model and training back-end multi-class logistic regression classifiers. All behavioral predictions of each visual system model were for images that were not seen in any phase of model training.

**Figure 2. Object-level comparison to human behavior. (A)** One-versus-all object-level (B.O1) signatures for the pooled human (n=1472 human subjects), pooled monkey (n=5 monkey subjects), and several $DCNN_{IC}$ models. Each B.O1 signature is shown as a 24-dimensional vector using a color scale; each colored bin corresponds to the system's discriminability of one

1012     object against all others that were tested. The color scales span each signature's full performance
1013     range, and warm colors indicate lower discriminability. **(B)** Direct comparison of the B.O1
1014     signatures of a pixel visual system model (top panel) and a $DCNN_{IC}$ visual system model
1015     (Inception-v3, bottom panel) against that of the human B.O1 signature. **(C)** *Human-consistency*
1016     of B.O1 signatures, for each of the tested model visual systems. The black and gray dots
1017     correspond to a held-out pool of five human subjects and a pool of five macaque monkey
1018     subjects respectively. The shaded area corresponds to the "primate zone," a range of
1019     consistencies delimited by the estimated *human-consistency* of a pool of infinitely many
1020     monkeys (see Figure 4A). **(D)** One-versus-other object-level (B.O2) signatures for pooled
1021     human, pooled monkey, and several $DCNN_{IC}$ models. Each B.O2 signature is shown as a 24x24
1022     symmetric matrices using a color scale, where each bin $(i,j)$ corresponds to the system's
1023     discriminability of objects $i$ and $j$. Color scales similar to (A). **(E)** Human-consistency of B.O2
1024     signatures for each of the tested model visual systems. Format is identical to (C).

1025     **Figure 3. Image-level comparison to human behavior. (A)** Schematic for computing B.I1n.
1026     First, the one-versus-all image-level signature (B.I1) is shown as a 240-dimensional vector (24
1027     objects, 10 images/object) using a color scale, where each colored bin corresponds to the
1028     system's discriminability of one image against all distractor objects. From this pattern, the
1029     normalized one-versus-all image-level signature (B.I1n) is estimated by subtracting the mean
1030     performance value over all images of the same object. This normalization procedure isolates
1031     behavioral variance that is specifically image-driven but not simply predicted by the object. **(B)**
1032     Normalized one-versus-all object-level (B.I1n) signatures for the pooled human, pooled monkey,
1033     and several $DCNN_{IC}$ models. Each B.I1n signature is shown as a 240-dimensional vector using a
1034     color scale, formatted as in (A). Color scales similar to Figure 2A. **(C)** *Human-consistency* of
1035     B.I1n signatures for each of the tested model visual systems. Format is identical to Figure 2C.
1036     **(D)** Normalized one-versus-other image-level (B.I2n) signatures for pooled human, pooled
1037     monkey, and several $DCNN_{IC}$ models. Each B.I2n signature is shown as a 240x24 matrix using a
1038     color scale, where each bin $(i,j)$ corresponds to the system's discriminability of image $i$ against
1039     distractor object $j$. Color scales similar to Figure 2A. **(E)** Human-consistency of B.I2n signatures
1040     for each of the tested model visual systems. Format is identical to Figure 2C.

1041     **Figure 4. Effect of subject pool size and DCNN model modifications on consistency with**

1042     **human behavior. (A)** Accounting for natural subject-to-subject variability. For each of the four

1043     behavioral metrics, the *human-consistency* distributions of monkey (blue markers) and model

1044     (black markers) pools are shown as a function of the number of subjects in the pool (mean ± SD,

1045     over subjects). The human consistency increases with growing number of subjects for all visual

1046     systems across all behavioral metrics. The dashed lines correspond to fitted exponential

1047     functions, and the parameter estimate (mean ± SE) of the asymptotic value, corresponding to the

1048     estimated *human-consistency* of a pool of infinitely many subjects, is shown at the right most

1049     point on each abscissa. **(B)** Model modifications that aim to rescue the $DCNN_{IC}$ models. We

1050     tested several simple modifications (see Methods) to the most *human-consistent* $DCNN_{IC}$ visual

1051     system model (Inception-v3). Each panel shows the resulting *human-consistency* per modified

1052     model (mean ± SD over different model instances, varying in random filter initializations) for

1053     each of the four behavioral metrics.

1054

1055     **Figure 5. Analysis of unexplained human behavioral variance. (A)** Residual similarity

1056     between all pairs of human visual system models. The color of bin ($i,j$) indicates the proportion

1057     of explainable variance that is shared between the residual signatures of visual systems $i$ and $j$.

1058     For ease of interpretation, we ordered visual system models based on their architecture and

1059     optimization procedure and partitioned this matrix into four distinct regions. **(B)** Summary of

1060     residual similarity. For each of the four regions in Figure 5A, the similarity to the residuals of

1061     Inception-v3 (region 2 in (A)) is shown (mean ± SD, within each region) for all images (black

1062     dots), and for images that humans found to be particularly difficult (gray dots, selected based on

1063     held-out human data).

1064

1065     **Figure 6. Dependence of primate and DCNN model behavior on image attributes. (A)**

1066     Example images with increasing attribute value, for each of the four pre-defined image attributes

1067     (see Methods). **(B)** Dependence of performance (B.I1n) as a function of four image attributes, for

1068     humans, monkeys and a $DCNN_{IC}$ model (Inception-v3). **(C)** Proportion of explainable variance

1069     of the residual signatures of monkeys (black) and $DCNN_{IC}$ models (blue) that is accounted for by

1070     each of the pre-defined image attributes. Error-bars correspond to SD over trial re-sampling for

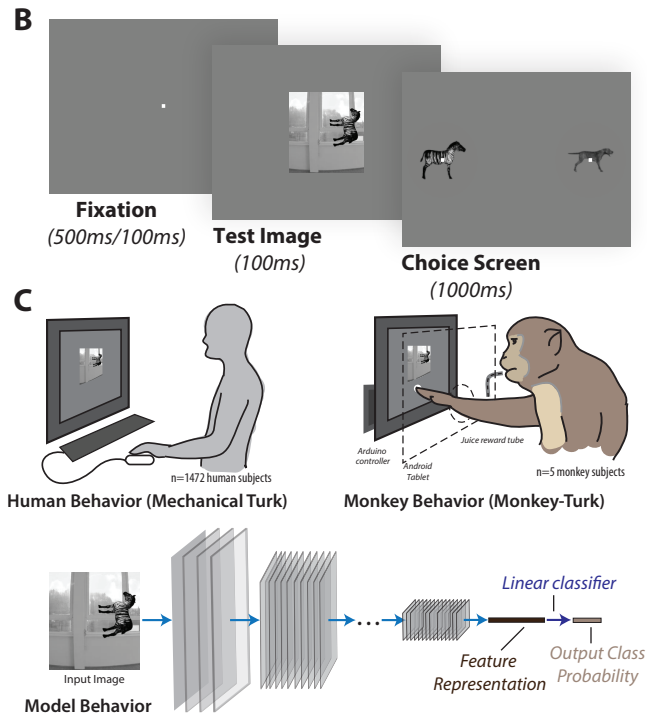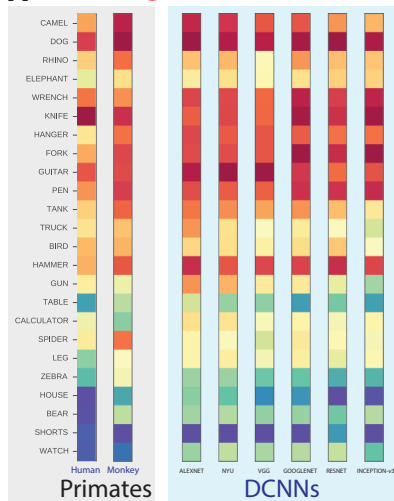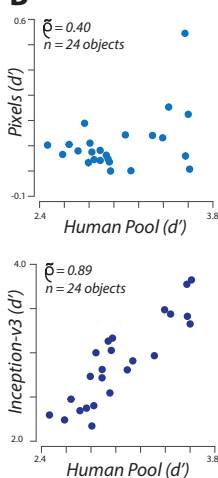1071     monkeys, and over different models for $DCNN_{IC}$ models.
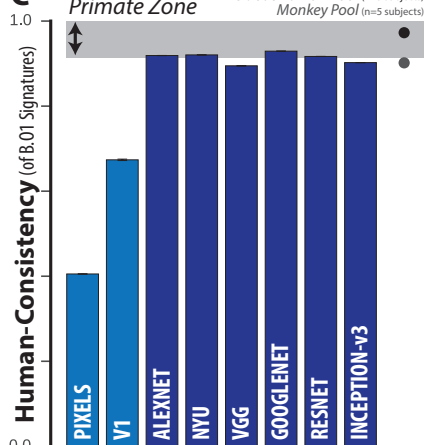
# Figure 1



**A**

100 testing images/object

| Elephant | Shorts | Bird | Tank | Camel | Leg | Rhino | Wrench |

| Bear | Guitar | Fork | Zebra | Hammer | Pen | Hanger | House |

| Knife | Gun | Calculator | Table | Truck | Spider | Dog | Watch |

**B**

**Fixation**
*(500ms/100ms)*

**Test Image**
*(100ms)*

**Choice Screen**
*(1000ms)*

**C**

**Human Behavior (Mechanical Turk)**
n=1472 human subjects

**Monkey Behavior (Monkey-Turk)**
Arduino controller
Android Tablet
Juice reward tube
n=5 monkey subjects

**Model Behavior**
Input Image

Feature Representation

*Linear classifier*

Output Class Probability

Figure 2



**A** **B.O1 Signatures** (~object difficulties)

Primates: Human, Monkey

Objects: CAMEL, DOG, RHINO, ELEPHANT, WRENCH, KNIFE, HANGER, FORK, GUITAR, PEN, TANK, TRUCK, BIRD, HAMMER, GUN, TABLE, CALCULATOR, SPIDER, LEG, ZEBRA, HOUSE, BEAR, SHORTS, WATCH

DCNNs: ALEXNET, NYU, VGG, GOOGLENET, RESNET, INCEPTION-v3

**B**

$\tilde{\rho} = 0.40$
$n = 24$ objects

Pixels (d')

Human Pool (d')

$\tilde{\rho} = 0.89$
$n = 24$ objects

Inception-v3 (d')

Human Pool (d')

**C**

*Primate Zone*

Heldout Human Pool (n=5 subjects)
Monkey Pool (n=5 subjects)

**Human-Consistency** (of B.O1 Signatures)

PIXELS, V1, ALEXNET, NYU, VGG, GOOGLENET, RESNET, INCEPTION-v3

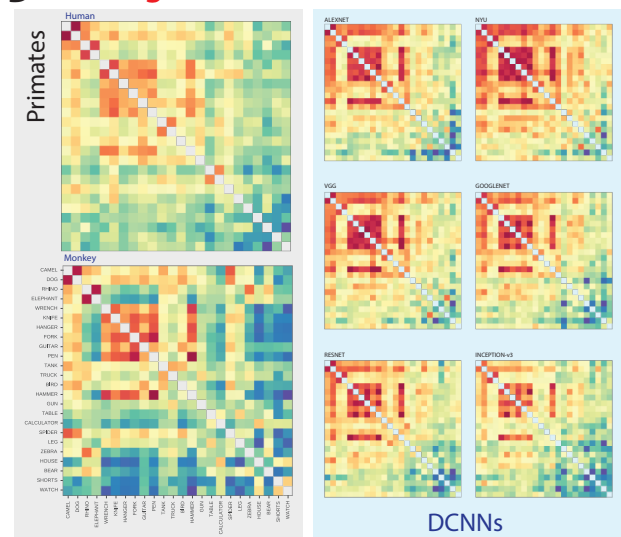**D** **B.O2 Signatures** (~object confusions)

Primates: Human, Monkey

DCNNs: ALEXNET, NYU, VGG, GOOGLENET, RESNET, INCEPTION-v3

**E**

*Primate Zone*

Heldout Human Pool (n=5 subjects)
Monkey Pool (n=5 subjects)

**Human-Consistency** (of B.O2 Signatures)

PIXELS, V1, ALEXNET, NYU, VGG, GOOGLENET, RESNET, INCEPTION-v3

Figure 3



**A** **B.I1n Signatures**
(~normalized image difficulties)

10 images/object

B.I1
(image-level
one-versus-all
performance)

mean
per object

B.I1n

**B**

CAMEL
DOG
RHINO
ELEPHANT
WRENCH
KNIFE
HANGER
FORK
GUITAR
PEN
TANK
TRUCK
BIRD
HAMMER
GUN
TABLE
CALCULATOR
SPIDER
LEG
ZEBRA
HOUSE
BEAR
SHORTS
WATCH

Human    Monkey

Primates

ALEXNET  NYU  VGG  GOOGLENET  RESNET  INCEPTION-v3

DCNNs

**C**

*Primate Zone*

Heldout Human Pool (n=5 subjects)
Monkey Pool (n=5 subjects)

1.0

Human-Consistency (of B.I1n Signatures)

0.0

PIXELS
V1
ALEXNET
NYU
VGG
GOOGLENET
RESNET
INCEPTION-v3

**D** **B.I2n Signatures** (~normalized image confusions)

CAMEL
DOG
RHINO
ELEPHANT
WRENCH
KNIFE
HANGER
FORK
GUITAR
PEN
TANK
TRUCK
BIRD
HAMMER
GUN
TABLE
CALCULATOR
SPIDER
LEG
ZEBRA
HOUSE
BEAR
SHORTS
WATCH

Human    Monkey

Primates

INCEPTION-v3

DCNN (ex.)

**E**

*Primate Zone*

Heldout Human Pool (n=5 subjects)
Monkey Pool (n=5 subjects)

1.0

Human-Consistency (of B.I2n Signatures)

0.0

PIXELS
V1
ALEXNET
NYU
VGG
GOOGLENET
RESNET
INCEPTION-v3

Figure 4



**A**

Human-Consistency

B.O1   B.O2   B.I1n   B.I2n

Monkey Pool
Model Pool

Extrapolation to infinitely many subjects

# Subjects

**B**

Human-Consistency

INCEPTION-v3
INCEPTION-v3 + retina
INCEPTION-v3 + SVM
INCEPTION-v3 + classifier_train
INCEPTION-v3 + synthetic_train

**Model Variations**

# Figure 5

## Figure 6