

Large-Scale, High-Resolution Comparison of the Core Visual Object Recognition Behavior of Humans, Monkeys, and State-of-the-Art Deep Artificial Neural Networks

 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

Primates, including humans, can typically recognize objects in visual images at a glance despite naturally occurring identity-preserving image transformations (e.g., changes in viewpoint). A primary neuroscience goal is to uncover neuron-level mechanistic models that quantitatively explain this behavior by predicting primate performance for each and every image. Here, we applied this stringent behavioral prediction test to the leading mechanistic models of primate vision (specifically, deep, convolutional, artificial neural networks; ANNs) by directly comparing their behavioral signatures against those of humans and rhesus macaque monkeys. Using high-throughput data collection systems for human and monkey psychophysics, we collected more than one million behavioral trials from 1472 anonymous humans and five male macaque monkeys for 2400 images over 276 binary object discrimination tasks. Consistent with previous work, we observed that state-of-the-art deep, feedforward convolutional ANNs trained for visual categorization (termed DCNN_{IC} models) accurately predicted primate patterns of object-level confusion. However, when we examined behavioral performance for individual images within each object discrimination task, we found that all tested DCNN_{IC} models were significantly nonpredictive of primate performance and that this prediction failure was not accounted for by simple image attributes nor rescued by simple model modifications. These results show that current DCNN_{IC} models cannot account for the image-level behavioral patterns of primates and that new ANN models are needed to more precisely capture the neural mechanisms underlying primate object vision. To this end, large-scale, high-resolution primate behavioral benchmarks such as those obtained here could serve as direct guides for discovering such models.

Key words: deep neural network; human; monkey; object recognition; vision

Significance Statement

Recently, specific feedforward deep convolutional artificial neural networks (ANNs) models have dramatically advanced our quantitative understanding of the neural mechanisms underlying primate core object recognition. In this work, we tested the limits of those ANNs by systematically comparing the behavioral responses of these models with the behavioral responses of humans and monkeys at the resolution of individual images. Using these high-resolution metrics, we found that all tested ANN models significantly diverged from primate behavior. Going forward, these high-resolution, large-scale primate behavioral benchmarks could serve as direct guides for discovering better ANN models of the primate visual system.

Introduction

Primates, both human and nonhuman, can typically recognize objects in visual images at a glance even in the face of naturally

occurring identity-preserving transformations such as changes in viewpoint. This view-invariant visual object recognition ability is thought to be supported primarily by the primate ventral visual stream (Tanaka, 1996; Rolls, 2000; DiCarlo et al., 2012). A pri-

Received Feb. 12, 2018; revised June 6, 2018; accepted July 8, 2018.

Author contributions: R.R., E.B.I., and J.J.D. designed research; R.R., E.B.I., K.K., and K.S. performed research; R.R. and P.B. analyzed data; R.R., E.B.I., and J.J.D. wrote the paper.

This work was supported by the National Institutes of Health—National Eye Institute (Grants R01-EY014970 to J.J.D. and K99-EY022671 to E.B.I.), the Office of Naval Research (Grant MURI-114407 to J.J.D.), Intelligence Advanced Research Projects Activity (Contract D16PC00002 to J.J.D.), the Engineering Research Council of Canada (R.R.), and the McGovern Institute for Brain Research.

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

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

Correspondence should be addressed to James J. DiCarlo, McGovern Institute for Brain Research, Department of Brain and Cognitive Sciences, Massachusetts Institute of Technology, MIT 46-6161, 43 Vassar Street, Cambridge, MA 02139. E-mail: dicarlo@mit.edu.

DOI:10.1523/JNEUROSCI.0388-18.2018

Copyright © 2018 the authors 0270-6474/18/387255-15\$15.00/0

many neuroscience goal is to construct computational models that quantitatively explain the neural mechanisms underlying this ability. That is, our goal is to discover artificial neural networks (ANNs) that accurately predict neuronal firing rate responses at all levels of the ventral stream and its behavioral output. To this end, specific models within a large family of deep, convolutional neural networks (DCNNs), optimized by supervised training on large-scale category-labeled image sets (ImageNet) to match human-level categorization performance (Krizhevsky et al., 2012; LeCun et al., 2015), have been put forth as the leading ANN models of the ventral stream (Kriegeskorte, 2015; Yamins and DiCarlo, 2016). We refer to this subfamily as DCNN_{IC} models (IC to denote ImageNet categorization pre-training) to distinguish them from all possible models in the DCNN family and, more broadly, from the superfamily of all ANNs. To date, it has been shown that DCNN_{IC} models display internal feature representations similar to neuronal representations along the primate ventral visual stream (Yamins et al., 2013, 2014; Cadieu et al., 2014; Khaligh-Razavi and Kriegeskorte, 2014) and they exhibit behavioral patterns similar to the behavioral patterns of pairwise object confusions of primates (Ghodrati et al., 2014; Rajalingham et al., 2015; Jozwik et al., 2016; Kheradpisheh et al., 2016). Therefore, DCNN_{IC} models may provide a quantitative account of the neural mechanisms underlying primate core object recognition behavior.

However, several studies have shown that DCNN_{IC} models can diverge drastically from humans in object recognition behavior, especially with regard to particular images optimized to be adversarial to these networks (Goodfellow et al., 2014; Nguyen et al., 2015). Related work has shown that specific image distortions are disproportionately challenging to current DCNNs compared with humans (RichardWebster et al., 2016; Dodge and Karam, 2017; Geirhos et al., 2017; Hosseini et al., 2017). Such image-specific failures of the current ANN models would likely not be captured by “object-level” behavioral metrics (e.g., the pattern of pairwise object confusions mentioned above) that are computed by pooling over hundreds of images and thus are not sensitive to variation in difficulty across images of the same object. To overcome this limitation of prior work, we here aimed to use scalable behavioral testing methods to precisely characterize primate behavior at the resolution of individual images and to directly compare leading DCNN models to primates over the domain of core object recognition behavior at this high resolution.

We focused on core invariant object recognition, the ability to identify objects in visual images in the central visual field during a single, natural viewing fixation (DiCarlo et al., 2012). We further restricted our behavioral domain to basic-level object discriminations, as defined previously (Rosch et al., 1976). Within this domain, we collected large-scale, high-resolution measurements of human and monkey behavior (over a million behavioral trials) using high-throughput psychophysical techniques, including a novel home cage behavioral system for monkeys. These data enabled us to systematically compare all systems at progressively higher resolution. At lower resolutions, we replicated previous findings that humans, monkeys, and DCNN_{IC} models all share a common pattern of object-level confusion (Rajalingham et al., 2015). However, at the higher resolution of individual images, we found that the behavior of all tested DCNN_{IC} models was significantly different from human and monkey behavior and this model prediction failure could not be easily rescued by simple model modifications. These results show that current DCNN_{IC} models do not fully account for the image-level behavioral patterns of primates, suggesting that new ANN models are needed to

more precisely capture the neural mechanisms underlying primate object vision. To this end, large-scale high-resolution behavioral benchmarks such as those obtained here could serve as a strong top-down constraint for efficiently discovering such models.

Materials and Methods

Visual images. We examined basic-level, core object recognition behavior using a set of 24 broadly sampled objects that we previously found to be reliably labeled by independent human subjects based on the definition of basic level proposed previously (Rosch et al., 1976). For each object, we generated 100 naturalistic synthetic images by first rendering a 3D model of the object with randomly chosen viewing parameters (2D position, 3D rotation, and viewing distance) and then placing that foreground object view onto a randomly chosen, natural image background. To do this, each object was first assigned a canonical position (center of gaze), scale ($\sim 2^\circ$), and pose and then its viewing parameters were randomly sampled uniformly from the following ranges for object translation ($[-3, 3]^\circ$ in both h and v), rotation ($[-180, 180]^\circ$ in all three axes), and scale ($[0.7, 1.7]$). Background images were sampled randomly from a large database of high-dynamic range images of indoor and outdoor scenes obtained from Dosch Design (www.doschdesign.com). This image generation procedure enforces invariant object recognition rather than image matching because it requires the visual recognition system (human, animal or model) to tackle the “invariance problem,” the computational crux of object recognition (Ullman and Humphreys, 1996; Pinto et al., 2008). In particular, we used naturalistic synthetic images with systematic variation in viewpoint parameters and uncorrelated background to remove potential confounds (natural images are often “composed” such that backgrounds covary with the object category) while keeping the task difficult for machine vision systems. We have previously shown that, unlike some photographic image sets, synthetic images of the types we used here are critical to separating some types of computer vision systems from humans (Pinto et al., 2008). Using this procedure, we previously generated 2400 images (100 images per object) rendered at 1024×1024 pixel resolution with 256-level gray scale and subsequently resized to 256×256 pixel resolution for human psychophysics, monkey psychophysics, and model evaluation (Rajalingham et al., 2015). In the current work, we focused our analyses on a randomly subsampled, and then fixed, subset of 240 images (10 images per object; here referred to as the “primary test images”). Figure 1A shows the full list of 24 objects, with two example images of each object.

Because all of the images were generated from synthetic 3D object models, we had explicit knowledge of the viewpoint parameters (position, size, and pose) for each object in each image, as well as perfect segmentation masks. Taking advantage of this feature, we characterized each image based on these high-level attributes consisting of size, eccentricity, relative pose, and contrast of the object in the image. Note that these meta-parameters were independently randomly sampled to generate each image, so there is no correlation among size, eccentricity, and pose over images. The size and eccentricity of the object in each image were computed directly from the corresponding viewpoint parameters under the assumption that the entire image would subtend 6° at the center of visual gaze ($\pm 3^\circ$ in both azimuth and elevation; see below). For each synthetic object, we first defined its “canonical” 3D pose vector based on independent human judgments. To compute the relative pose attribute of each image, we estimated the difference between the object’s 3D pose and its canonical 3D pose. Pose differences were computed as distances in unit quaternion representations: the 3D pose (r_{xy} , r_{xz} , and r_{yz}) was first converted into unit quaternions and distances between quaternions q_1 and q_2 were estimated as $\cos^{-1}[q_1 \cdot q_2]$ (Huynh, 2009). To compute the object contrast, we measured the absolute difference between the mean of the pixel intensities corresponding to the object and the mean of the background pixel intensities in the vicinity of the object (specifically, within 25 pixels of any object pixel, analogous to computing the local foreground-background luminance difference of a foreground object in an image). Figure 5A shows example images with varying values for the four image attributes.

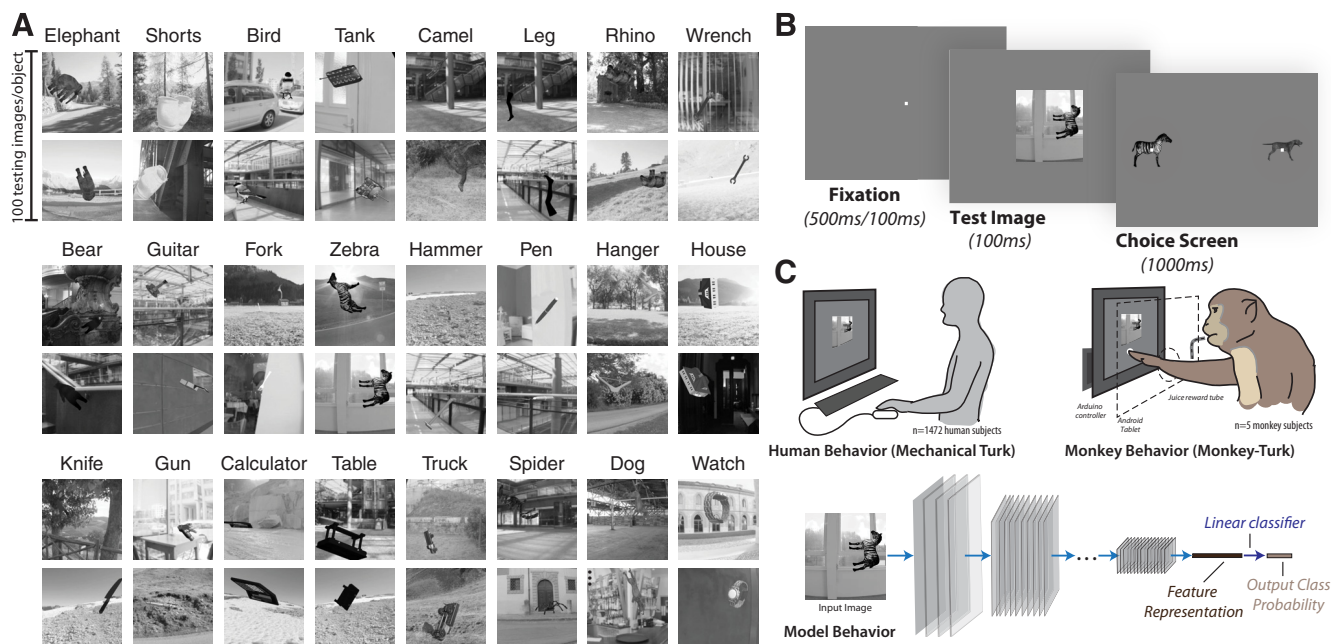


Figure 1. Images and behavioral task. **A**, Two (of 100) example images for each of the 24 basic-level objects. To enforce true invariant object recognition behavior, we generated naturalistic synthetic images, each with one foreground object, by rendering a 3D model of each object with randomly chosen viewing parameters and placing that foreground object view onto a randomly chosen, natural image background. **B**, Time course of example behavioral trial (zebra vs 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, whereas incorrect choices resulted in timeouts of 1–2.5 s. **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 of ~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 multiclass logistic regression classifiers. All behavioral predictions of each visual system model were for images that were not seen in any phase of model training.

Core object recognition behavioral paradigm. Core object discrimination is defined as the ability to discriminate between two or more objects in visual images presented under high view uncertainty in the central visual field (~10°) for durations that approximate the typical primate free-viewing fixation duration (~200 ms) (DiCarlo and Cox, 2007; DiCarlo et al., 2012). As in our previous work (Rajalingham et al., 2015), the behavioral task paradigm consisted of an interleaved set of binary discrimination tasks. Each binary discrimination task is an object discrimination task between a pair of objects (e.g., elephant vs bear). Each such binary task is balanced in that the test image is equally likely (50%) to be of either of the two objects. On each trial, a test image is presented, followed by a choice screen showing canonical views of the two possible objects (the object that was not displayed in the test image is referred to as the “distractor” object, but note that objects are equally likely to be distractors and targets). Here, 24 objects were tested, which resulted in 276 binary object discrimination tasks. To neutralize feature attention, these 276 tasks are randomly interleaved (trial by trial) and the global task is referred to as a basic-level, core object recognition task paradigm.

Testing human behavior. All human behavioral data presented here were collected from 1476 anonymous human subjects on Amazon Mechanical Turk (MTurk) performing the task paradigm described above. Subjects were instructed to report the identity of the foreground object in each presented image from among the two objects presented on the choice screen (see Fig. 1B). Because all 276 tasks were interleaved randomly (trial by trial), subjects could not deploy feature attentional strategies specific to each object or specific to each binary task to process each test image.

Figure 1B illustrates the time course of each behavioral trial, for a particular object discrimination task (e.g., zebra vs dog). Each trial initiated with a central black point for 500 ms, followed by 100 ms presenta-

tion of a test image containing one foreground object presented under high variation in viewing parameters and overlaid on a random background, as described above (see “Visual images” section above). Immediately after extinction of the test image, two choice images, each displaying a single object in a canonical view with no background, were shown to the left and right. One of these two objects was always the same as the object that generated the test image (i.e., the correct object choice) and the location of the correct object (left or right) was randomly chosen on each trial. After clicking on one of the choice images, the subject was queued with another fixation point before the next test image appeared. No feedback was given and human subjects were never explicitly trained on the tasks. Under assumptions of typical computer ergonomics, we estimate that images were presented at a 6–8° of visual angle at the center of gaze and the choice object images were presented at ±6–8° of eccentricity along the horizontal meridian.

We measured human behavior using the online Amazon MTurk platform (see Fig. 1C), which enables efficient collection of large-scale psychophysical data from crowd-sourced “human intelligence tasks” (HITs). The reliability of the online MTurk platform has been validated by comparing results obtained from online and in-laboratory psychophysical experiments (Majaj et al., 2015; Rajalingham et al., 2015). We pooled 927,296 trials from 1472 human subjects to characterize the aggregate human behavior, which we refer to as the “pooled” human (or “archetypal” human). Each human subject performed only a small number of trials (~150) on a subset of the images and binary tasks. All 2400 images were used for behavioral testing but, in some of the HITs, we biased the image selection toward the 240 primary test images (1424 ± 70 trials/image on this subsampled set vs 271 ± 93 trials/image on the remaining images, mean \pm SD) to efficiently characterize behavior at image-level resolution. Images were randomly drawn such that each hu-

man subject was exposed to each image a relatively small number of times (1.5 ± 2.0 trials/image per subject, mean \pm SD), to mitigate potential alternative behavioral strategies (e.g., “memorization” of images) that could arise from a finite image set. Behavioral signatures at the object-level (B.O1, B.O2, see “Behavioral metrics and signatures” section) were measured using all 2400 test images, whereas image-level behavioral signatures (B.I1n, B.I2n, see “Behavioral metrics and signatures” section) were measured using the 240 primary test images. (We observed qualitatively similar results using those metrics on the full 2400 test images, but we here focus on the primary test images because the larger number of trials leads to lower noise levels).

Five other human subjects were separately recruited on MTurk to each perform a large number of trials on the same images and tasks (53,097 \pm 15,278 trials/subject, mean \pm SD). Behavioral data from these five subjects was not included in the characterization of the pooled human described above, but instead aggregated together to characterize a distinct held-out human pool. For the scope of the current work, this held-out human pool, which largely replicated all behavioral signatures of the larger archetypal human (see Figs. 2 and 3), served as an independent validation of our human behavioral measurements.

Testing monkey behavior. Five adult male rhesus macaque monkeys (*Macaca mulatta*, subjects M, Z, N, P, and B) were tested on the same basic-level, core object recognition task paradigm described above, with minor modification as described below. All procedures were performed in compliance with National Institutes of Health guidelines and the standards of the Massachusetts Institute of Technology Committee on Animal Care and the American Physiological Society. To efficiently characterize monkey behavior, we used a novel home cage behavioral system that we developed (termed MonkeyTurk; see Fig. 1C). This system leveraged a tablet touchscreen (9-inch Google Nexus or 10.5-inch Samsung Galaxy Tab S) and used a web application to wirelessly load the task and collect the data (code available from <https://github.com/dicarlolab/mkturk>). Analogous to the online MTurk, which allows for efficient psychophysical assays of a large number (hundreds) of human users in their native environments, MonkeyTurk allowed us to test many monkey subjects simultaneously in their home environment. Each monkey voluntarily initiated trials and each readily performed the task a few hours each day that the task apparatus was made available to it. At an average rate of ~ 2000 trials per day per monkey, we collected a total of 836,117 trials from the 5 monkey subjects over a period of ~ 3 months.

Monkey training was described in detail previously (Rajalingham et al., 2015). Briefly, all monkeys were initially trained on the match test image to object rule using other images and were also trained on discriminating the particular set of 24 objects tested here using a separate set of training images rendered from these objects in the same manner as the main testing images. Two of the monkey subjects (Z and M) were previously trained in the laboratory setting and the remaining three subjects were trained using MonkeyTurk directly in their home cages and did not have significant prior laboratory exposure. Once monkeys reached saturation performance on training images, we began the behavioral testing phase to collect behavior on test images. Monkeys did improve throughout the testing phase, exhibiting an increase in performance between the first and second half of trials of $4 \pm 0.9\%$ (mean \pm SEM over five monkey subjects). However, the image-level behavioral signatures obtained from the first and the second halves of trials were highly correlated to each other (B.I1 noise-adjusted correlation of 0.85 ± 0.06 , mean \pm SEM over five monkey subjects, see “Behavioral metrics and signatures” section below), suggesting that monkeys did not significantly alter strategies (e.g., did not “memorize” images) throughout the behavioral testing phase.

The monkey task paradigm was nearly identical to the human paradigm (see Fig. 1B) with the exception that trials were initiated by touching a white “fixation” circle horizontally centered on the bottom third of the screen (to avoid occluding centrally presented test images with the hand). This triggered a 100 ms central presentation of a test image, followed immediately by the presentation of the two choice images (see Fig. 1B, location of correct choice randomly assigned on each trial, identical to the human task). Unlike the main human task, monkeys responded by directly touching the screen at the location of one of the two choice

images. Touching the choice image corresponding to the object shown in the test image resulted in the delivery of a drop of juice through a tube positioned at mouth height (but not obstructing view), whereas touching the distractor choice image resulted in a 3 s timeout. Because gaze direction typically follows the hand during reaching movements, we assumed that the monkeys were looking at the screen during touch interactions with the fixation or choice targets. In both the laboratory and in the home cage, we maintained total test image size at $\sim 6^\circ$ of visual angle at the center of gaze and took advantage of the retina-like display qualities of the tablet by presenting images pixel matched to the display (256×256 pixel image displayed using 256×256 pixels on the tablet at a distance of 8 inches) to avoid filtering or aliasing effects.

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

Testing model behavior. We tested a number of different DCNN models on the exact same images and tasks as those presented to humans and monkeys. Importantly, our core object recognition task paradigm is closely analogous to the large-scale ImageNet 1000-way object categorization task for which these networks were optimized and thus expected to perform well. We focused on publicly available DCNN model architectures that have proven highly successful with respect to this computer vision benchmark over the past 5 years: AlexNet (Krizhevsky et al., 2012), NYU (Zeiler and Fergus, 2014), VGG (Simonyan and Zisserman, 2014), GoogleNet (Szegedy et al., 2013), Resnet (He et al., 2016), and Inception-v3 (Szegedy et al., 2013). Because this is only a subset of possible DCNN models, we refer to these as the DCNN_{IC} visual system model subfamily. For each of the publicly available model architectures, we first used ImageNet categorization-trained model instances either using publicly available trained model instances or training them to saturation on the 1000-way classification task in-house. Training took several days on one to two GPUs.

We then performed psychophysical experiments on each ImageNet-trained DCNN model to characterize their behavior on the exact same images and tasks as humans and monkeys. We first adapted these ImageNet-trained models to our 24-way object recognition task by retraining the final class probability layer (initially corresponding to the probability output of the 1000-way ImageNet classification task) while holding all other layers fixed. In practice, this was done by extracting features from the penultimate layer of each DCNN_{IC} (i.e., top-most before class probability layer) on the same images that were presented to humans and monkeys and training back-end multiclass logistic regression classifiers to determine the cross-validated output class probability for each image. This procedure is illustrated in Figure 1C. To estimate the hit rate of a given image in a given binary classification task, we renormalized the 24-way class probabilities of that image, considering only the two relevant classes, to sum to one. Object-level and image-level behavioral metrics were computed based on these hit rate estimates (as described in the “Behavioral metrics and signatures” section below). Importantly, this procedure assumes that the model “retina” layer processes the central 6° of the visual field. It also assumes that linear discrim-

Table 1. Definition of behavioral performance metrics

Behavioral metric	HR	FAR
One-versus-all object-level performance (B.O1) ($N_{\text{objects}} \times 1$) $O_1(i) = Z(\text{HR}(i)) - Z(\text{FAR}(i))$, $i = 1, 2, \dots, N_{\text{objects}}$	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) ($N_{\text{objects}} \times N_{\text{objects}}$) $O_2(i, j) = Z(\text{HR}(i, j)) - Z(\text{FAR}(i, j))$, $i = 1, 2, \dots, N_{\text{objects}}$ $j = 1, 2, \dots, N_{\text{objects}}$	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) ($N_{\text{images}} \times 1$) $I_1(ii) = Z(\text{HR}(ii)) - Z(\text{FAR}(ii))$, $ii = 1, 2, \dots, N_{\text{images}}$	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) ($N_{\text{images}} \times N_{\text{objects}}$) $I_2(ii, j) = Z(\text{HR}(ii, j)) - Z(\text{FAR}(ii, j))$, $ii = 1, 2, \dots, N_{\text{images}}$ $j = 1, 2, \dots, N_{\text{objects}}$	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

The first column provides the name, abbreviation, dimensions, and equations for each of the raw performance metrics. The next two columns provide the definitions for computing the hit rate (HR) and false alarm rate (FAR), respectively.

inants (“readouts”) of the model’s top feature layer are its behavioral output (as intended by the model designers). Manipulating either of these choices (e.g., resizing the input images such that they span only part of the input layer or building linear discriminates for behavior using a different model feature layer) would result in completely new, testable ANN models that we did not test here.

From these analyses, we selected the most human consistent DCNN_{IC} architecture (Inception-v3, see “Behavioral consistency” section below), fixed that architecture, and then performed *post hoc* analyses in which we varied the input image sampling, the initial parameter settings before training, the filter training images, the type of classifiers used to generate the behavior from the model features, and the classifier training images. To examine input image sampling, we retrained the Inception-v3 architecture on images from ImageNet that were first spatially filtered to match the spatial sampling of the primate retina (i.e., an approximately exponential decrease in cone density away from the fovea) by effectively simulating a fisheye transformation on each image. These images were at highest resolution at the “fovea” (i.e., center of the image) with gradual decrease in resolution with increasing eccentricity. To examine the analog of “intersubject variability,” we constructed multiple trained model instances (“subjects”) in which the architecture and training images were held fixed (Inception-v3 and ImageNet, respectively) but the model filter weights initial condition and order of training images were randomly varied for each model instance. Importantly, this procedure is only one possible choice for simulating intersubject variability for DCNN models, which, as a first-order approximation, models different human subjects as random samples from a fixed model class. There are many other possible options for simulating intersubject variability, including: (1) random sampling of different features from a fixed trained model of fixed architecture; (2) random sampling of different trained models from a fixed architecture and optimization; and (3) random sampling of different trained models from a model class varying in architecture (3a), optimization procedure and data (3b), or both (3c), to name a few. This choice is an important open research direction that we do not address here. We do not claim this to be the best model of intersubject variability, but rather a good starting point for that greater goal.

To examine the effect of model training, we first fine-tuned an ImageNet-trained Inception-v3 model on a synthetic image set consisting of ~6.9 million images of 1049 objects (holding out 50,000 images for model validation). These images were generated using the same rendering pipeline as our test images, but the objects were nonoverlapping with the 24 test objects presented here. As expected, fine-tuning on synthetic images led to a small increase in overall performance (see Fig. 5). To push the limits of training on synthetic images, we additionally trained an Inception-v3 architecture exclusively on this synthetic image set. To avoid overfitting on this synthetic image set, we stopped the model training based on maximum performance on a validation set of held-out

objects. These synthetic-trained models were high performing, with only a small decrease in overall performance relative to the ImageNet-trained models (see Fig. 5, last bar). Finally, we tested the effect of different classifiers to generate model behavior by testing both multiclass logistic regression and support vector machine classifiers as well as the effect of varying the number of training images used to train those classifiers (20 vs 50 images per class).

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

The four behavioral metrics we chose are as follows. First, the one-versus-all object-level performance metric (termed B.O1) estimates the discriminability of each object from all other objects, pooling across all distractor object choices. Because we here tested 24 objects, the resulting B.O1 signature has 24 independent values. Second, the one-versus-other object-level performance metric (termed B.O2) estimates the discriminability of each specific pair of objects or the pattern of pairwise object confusions. Because we here tested 276 interleaved binary object discrimination tasks, the resulting B.O2 signature has 276 independent values (the off-diagonal elements on one-half of the 24×24 symmetric matrix). Third, the one-versus-all image-level performance metric (termed B.I1) estimates the discriminability of each image from all other objects, pooling across all possible distractor choices. Because we here focused on the primary image test set of 240 images (10 per object, see above), the resulting B.I1 signature has 240 independent values. Fourth, the one-versus-other image-level performance metric (termed B.I2) estimates the discriminability of each image from each distractor object. Because we here focused on the primary image test set of 240 images (10 per object, see above) with 23 distractors, the resulting B.I2 signature has 5520 independent values.

Naturally, object-level and image-level behavioral signatures are tightly linked. For example, images of a particularly difficult-to-discriminate object would inherit lower performance values on average compared with images from a less difficult-to-discriminate object. To isolate the

behavioral variance that is specifically driven by image variation and not simply predicted by the objects (and thus already captured by B.O1 and B.O2), we defined normalized image-level behavioral metrics (termed B.I1n, B.I2n) by subtracting the mean performance values over all images of the same object and task. This process is schematically illustrated in Figure 3A. We note that the resulting normalized image-level behavioral signatures capture a significant proportion of the total image-level behavioral variance in our data (e.g., 52%, 58% of human B.I1 and B.I2 variance is driven by image variation, independent of object identity). In this study, we use these normalized metrics for image-level behavioral comparisons between models and primates (see Results).

Behavioral consistency. To quantify the similarity between a model visual system and the human visual system with respect to a given behavioral metric, we used a measure called “human consistency” as previously defined (Johnson et al., 2002). Human consistency ($\bar{\rho}$) is computed, for each of the four behavioral metrics, as a noise-adjusted correlation of behavioral signatures (DiCarlo and Johnson, 1999). For each visual system, we randomly split all behavioral trials into two equal halves and applied each behavioral metric to each half, resulting in two independent estimates of the system’s behavioral signature with respect to that metric. We took the Pearson correlation between these two estimates of the behavioral signature as a measure of the reliability of that behavioral signature given the amount of data collected; that is, the split-half internal reliability. To estimate human consistency, we computed the Pearson correlation over all the independent estimates of the behavioral signature from the model (m) and the human (h) and we then divide that raw Pearson correlation by the geometric mean of the split-half internal reliability of the same behavioral signature measured for each system as follows: $\bar{\rho}(m, h) = \frac{\rho(m, h)}{\sqrt{\rho(m, m)\rho(h, h)}}$.

Because all correlations in the numerator and denominator were computed using the same amount of trial data (exactly half of the trial data), we did not need to make use of any prediction formulas (e.g., extrapolation to larger number of trials using Spearman–Brown prediction formula). This procedure was repeated 10 times with different random split-halves of trials. Our rationale for using a reliability-adjusted correlation measure for human consistency was to account for variance in the behavioral signatures that arises from “noise”; that is, variability that is not replicable by the experimental condition (image and task) and thus that no model can be expected to predict (DiCarlo and Johnson, 1999; Johnson et al., 2002). In sum, if the model (m) is a replica of the archetypal human (h), then its expected human consistency is 1.0 regardless of the finite amount of data that are collected. Note that the human consistency value is directly linked, via a squaring operation, to the proportion of the explainable behavioral variance that is explained by models.

Characterization of residuals. In addition to measuring the similarity between the behavioral signatures of primates and models (using human consistency analyses, as described above), we examined the corresponding differences, termed “residual signatures.” Each candidate visual system model’s residual signature was estimated as the residual of a linear least-squares regression of the model’s signature on the corresponding human signature and a free intercept parameter. This procedure effectively captures the differences between human and model signatures after accounting for overall performance differences. Residual signatures were estimated on disjoint split-halves of trials repeated 10 times with random trial permutations. Residuals were computed with respect to the normalized one-versus-all image-level performance metric (B.I1n) because this metric showed a clear difference between DCNN_{IC} models and primates and the behavioral residual can be interpreted based only on the test images (i.e., we can assign a residual per image).

To examine the extent to which the difference between each model and the archetypal human is reliably shared across different models, we measured the Pearson correlation between the residual signatures of pairs of models. Residual similarity was quantified as the proportion of shared variance, defined as the square of the noise-adjusted correlation between residual signatures (the noise adjustment was done as defined in equation above). Correlations of residual signatures were always computed across distinct split-halves of data to avoid introducing spurious correlations from subtracting common noise in the human data. We mea-

sured the residual similarity between all pairs of tested models, holding both architecture and optimization procedure fixed (between instances of the ImageNet categorization-trained Inception-v3 model varying in filter initial conditions), varying the architecture while holding the optimization procedure fixed (between all tested ImageNet categorization-trained DCNN architectures), and holding the architecture fixed while varying the optimization procedure (between ImageNet categorization-trained Inception-v3 and synthetic categorization fine-tuned Inception-v3 models). This analysis addresses not only the reliability of the failure of DCNN_{IC} models to predict human behavior (deviations from humans), but also the relative importance of the characteristics defining similarities within the model subfamily (namely, the architecture and the optimization procedure). We first performed this analysis for residual signatures over the 240 primary test images and subsequently zoomed in on subsets of images that humans found to be particularly difficult. This image selection was made relative to the distribution of image-level performance of held-out human subjects (B.I1 metric from five subjects); difficult images were defined as ones with performance below the 25th percentile of this distribution.

To determine whether the difference between each model and humans can be explained by simple human-interpretable stimulus attributes, we regressed each DCNN_{IC} model’s residual signature on image attributes (object size, eccentricity, pose, and contrast). Briefly, we constructed a design matrix from the image attributes (using individual attributes or all attributes) and used multiple linear least-squares regression to predict the image-level residual signature. The multiple linear regression was tested using twofold cross-validation over trials. The relative importance of each attribute (or groups of attributes) was quantified using the proportion of explainable variance (i.e., variance remaining after accounting for noise variance) explained from the residual signature.

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

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

To make a global statistical inference about whether models sampled from the DCNN_{IC} subfamily meet or fall short of this criterion for success, we attempted to reject the hypothesis that, for a given behavioral metric, the human consistency of DCNN_{IC} models is within the primate zone. To test this hypothesis, we estimated the empirical probability that the distribution of human consistency values, estimated over different model instances within this family, could produce human consistency values within the primate zone. Specifically, we estimated a p -value for each behavioral metric using the following procedure: We first estimated

an empirical distribution of Fisher-transformed human consistency values for this model family (i.e., over all tested DCNN_{IC} models and over all trial resampling of each DCNN_{IC} model). From this empirical distribution, we fit a Gaussian kernel density function, optimizing the bandwidth parameter to minimize the mean squared error to the empirical distribution. This kernel density function was evaluated to compute a p -value by computing the cumulative probability of observing a human consistency value greater than or equal to the criterion of success (i.e., the Fisher transformed $\bar{\rho}_{M\infty}$ value). This p -value indicates the probability that human consistency values sampled from the observed distribution would fall into the primate zone, with smaller p -values indicating stronger evidence against the hypothesis that the human consistency of DCNN models is within the primate zone.

Results

In the present work, we systematically compared the basic level core object recognition behavior of primates and state-of-the-art artificial neural network models using a series of behavioral metrics ranging from low to high resolution within a two-alternative forced choice match-to-sample paradigm. The behavior of each visual system, whether biological or artificial, was tested on the same 2400 images (24 objects, 100 images/object) in the same 276 interleaved binary object recognition tasks. Each system's behavior was characterized at multiple resolutions (see "Behavioral metrics and signatures" section in the Materials and Methods) and directly compared with the corresponding behavioral metric applied on the archetypal human (defined as the average behavior of a large pool of human subjects tested; see Materials and Methods). The overarching logic of this study was that, if two visual systems are equivalent, then they should produce statistically indistinguishable behavioral signatures with respect to these metrics. Specifically, our goal was to compare the behavioral signatures of visual system models with the corresponding behavioral signatures of primates.

Object-level behavioral comparison

We first examined the pattern of one-versus-all object-level behavior (B.O1 metric) computed across all images and possible distractors. Because we tested 24 objects here, the B.O1 signature was 24 dimensional. Figure 2A shows the B.O1 signatures for the pooled human (pooling $n = 1472$ human subjects), pooled monkey (pooling $n = 5$ monkey subjects), and several DCNN_{IC} models as 24-dimensional vectors using a color scale. Each element of the vector corresponds to the system's discriminability of one object against all others that were tested (i.e., all other 23 objects). The color scales span each signature's full performance range and warm colors indicate lower discriminability. For example, red indicates that the tested visual system found the object corresponding to that element of the vector to be very challenging to discriminate from other objects (on average over all 23 discrimination tests and on average over all images). Figure 2B directly compares the B.O1 signatures computed from the behavioral output of two visual system models, a pixel model (Fig. 2B, top) and a DCNN_{IC} model (Inception-v3; Fig. 2B, bottom), against that of the human B.O1 signature. We observe a tighter correspondence to the human behavioral signature for the DCNN_{IC} model visual system than for the baseline pixel model visual system. We quantified that similarity using a noise-adjusted correlation between each pair of B.O1 signatures (termed human consistency; Johnson et al., 2002). The noise adjustment means that a visual system that is identical to the human pool will have an expected human consistency score of 1.0 even if it has irreducible trial-by-trial stochasticity (see Materials and Methods). Figure 2C shows the B.O1 human consistency for each of the tested

model visual systems. We additionally tested the behavior of a held-out pool of five human subjects (black dot) and a pool of five macaque monkey subjects (gray dot) and observed that both yielded B.O1 signatures that were highly human consistent (human consistency, $\bar{\rho} = 0.90, 0.97$ for monkey pool and held-out human pool, respectively). We defined a range of human consistency values, termed the "primate zone" (shaded gray area), delimited by extrapolated human consistency estimates of large pools of macaques (see Materials and Methods; see Fig. 4). We found that the baseline pixel visual system model and the low-level V1 visual system model were not within this zone ($\bar{\rho} = 0.40, 0.67$ for pixels and V1 models, respectively), whereas all tested DCNN_{IC} visual system models were either within or very close to this zone. Indeed, we could not reject the hypothesis that DCNN_{IC} models are primate like ($p = 0.54$, exact test; see Materials and Methods).

Next, we compared the behavior of the visual systems at a slightly higher level of resolution. Specifically, instead of pooling over all discrimination tasks for each object, we computed the mean discriminability of each of the 276 pairwise discrimination tasks (still pooling over images within each of those tasks). This yielded a symmetric matrix that is referred to here as the B.O2 signature. Figure 2D shows the B.O2 signatures of the pooled human, pooled monkey and several DCNN_{IC} visual system models as 24×24 symmetric matrices. Each bin (i, j) corresponds to the system's discriminability of objects i and j , where warmer colors indicate lower performance; color scales are not shown but span each signature's full range. We observed strong qualitative similarities between the pairwise object confusion patterns of all of the high level visual systems (e.g., camel and dog are often confused with each other by all three systems). This similarity is quantified in Figure 2E, which shows the human consistency of all examined visual system models with respect to this metric. Similar to the B.O1 metric, we observed that both a pool of macaque monkeys and a held-out pool of humans are highly human consistent with respect to this metric ($\bar{\rho} = 0.77, 0.94$ for monkeys, humans respectively). Also similar to the B.O1 metric, we found that all DCNN_{IC} visual system models are highly human consistent ($\bar{\rho} > 0.8$), whereas the baseline pixel visual system model and the low-level V1 visual system model were not ($\bar{\rho} = 0.41, 0.57$ for pixels, V1 models, respectively). Indeed, all DCNN_{IC} visual system models are within the defined "primate zone" of human consistency and we could not falsify the hypothesis that DCNN_{IC} models are primate like ($p = 0.99$, exact test).

Together, humans, monkeys, and current DCNN_{IC} models all share similar patterns of object-level behavioral performances (B.O1 and B.O2 signatures) that are not shared with lower-level visual representations (pixels and V1). However, object-level performance patterns do not capture the fact that some images of an object are more challenging than other images of the same object because of interactions of the variation in the object's pose and position with the object's class. To overcome this limitation, we next examined the patterns of behavior at the resolution of individual images on a subsampled set of images where we specifically obtained a large number of behavioral trials to accurately estimate behavioral performance on each image. Note that, from the point of view of the subjects, the behavioral tasks are identical to those already described. We simply aimed to measure and compare their patterns of performance at much higher resolution.

Image-level behavioral comparison

To isolate purely image-level behavioral variance, that is, variance that is not predicted by the object and thus already captured

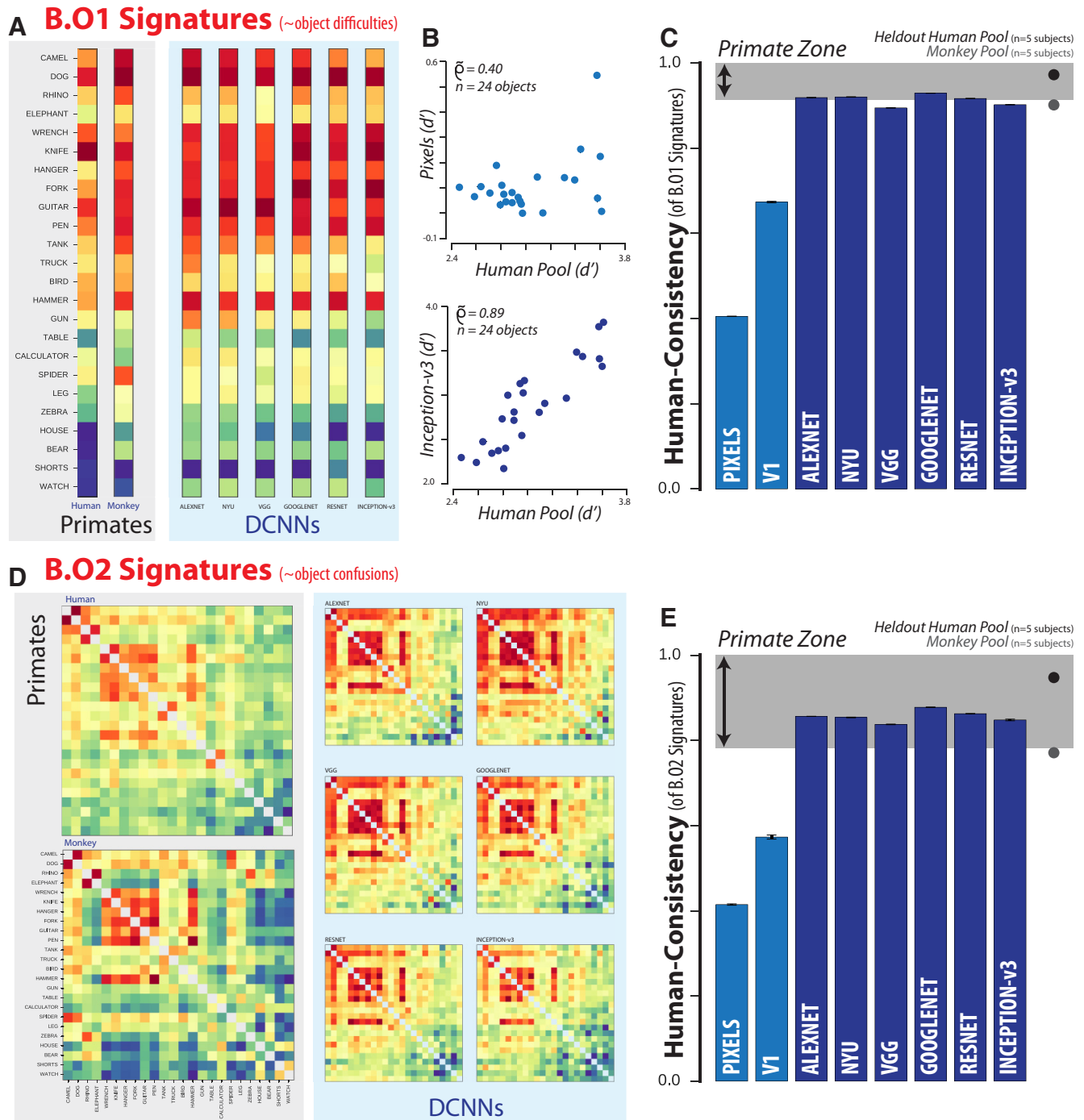


Figure 2. Object-level comparison to human behavior. **A**, One-versus-all object-level (B.O1) signatures for the pooled human ($n = 1472$ human subjects), pooled monkey ($n = 5$ monkey subjects), and several DCNN_{IC} models. Each B.O1 signature is shown as a 24-dimensional vector using a color scale; each colored bin corresponds to the system's discriminability of one object against all others that were tested. The color scales span each signature's full performance range and warm colors indicate lower discriminability. **B**, Direct comparison of the B.O1 signatures of a pixel visual system model (top) and a DCNN_{IC} visual system model (Inception-v3; bottom) against that of the human B.O1 signature. **C**, Human consistency of B.O1 signatures for each of the tested model visual systems. The black and gray dots correspond to a held-out pool of five human subjects and a pool of five macaque monkey subjects, respectively. The shaded area corresponds to the "primate zone," a range of consistencies delimited by the estimated human consistency of a pool of infinitely many monkeys (see Fig. 4A). **D**, One-versus-other object-level (B.O2) signatures for pooled human, pooled monkey, and several DCNN_{IC} models. Each B.O2 signature is shown as a 24×24 symmetric matrices using a color scale, where each bin (i,j) corresponds to the system's discriminability of objects i and j . As in **A**, color scales span each signature's full performance range and warm colors indicate lower discriminability. **E**, Human consistency of B.O2 signatures for each of the tested model visual systems. Format is identical to that in **C**.

by the B.O1 signature, we computed the normalized image-level signature. This normalization procedure is schematically illustrated in Figure 3A, which shows that the one-versus-all image-level signature (240-dimensional, 10 images/object) is used to obtain the normalized one-versus-all image-level signature

(termed B.I1n, see "Behavioral metrics and signatures" section). Figure 3B shows the B.I1n signatures for the pooled human, pooled monkey, and several DCNN_{IC} models as 240 dimensional vectors. Each bin's color corresponds to the discriminability of a single image against all distractor options (after subtraction of

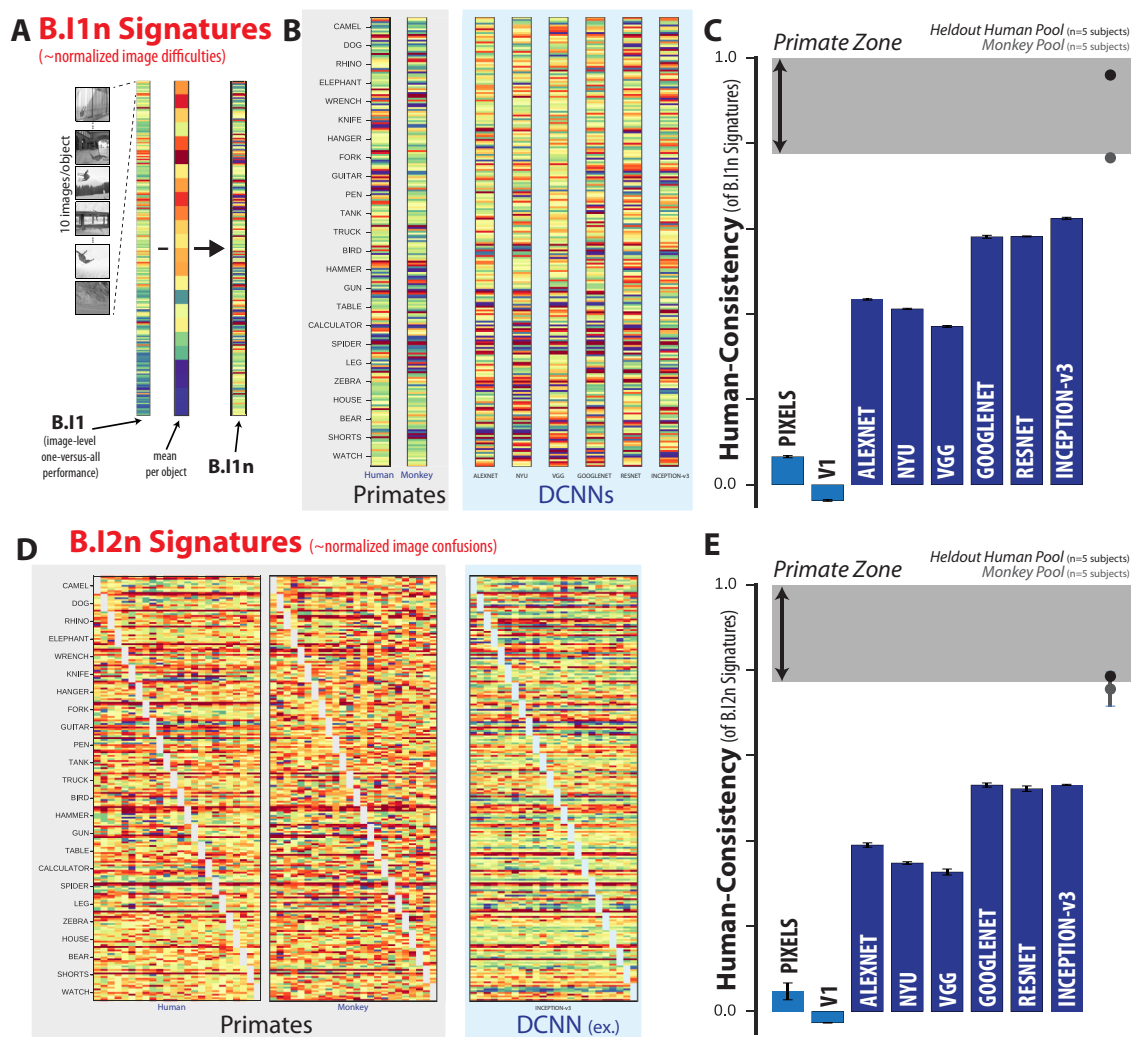


Figure 3. Image-level comparison to human behavior. **A**, Schematic for computing B.11n. First, the one-versus-all image-level signature (B.11) is shown as a 240-dimensional vector (24 objects, 10 images/object) using a color scale, where each colored bin corresponds to the system's discriminability of one image against all distractor objects. From this pattern, the normalized one-versus-all image-level signature (B.11n) is estimated by subtracting the mean performance value over all images of the same object. This normalization procedure isolates behavioral variance that is specifically image driven but not simply predicted by the object. **B**, Normalized one-versus-all object-level (B.11n) signatures for the pooled human, pooled monkey, and several DCNN_{IC} models. Each B.11n signature is shown as a 240-dimensional vector using a color scale formatted as in **A**. **C**, Human consistency of B.11n signatures for each of the tested model visual systems. Format is identical to that in Figure 2C. **D**, Normalized one-versus-other image-level (B.12n) signatures for pooled human, pooled monkey, and several DCNN_{IC} models. Each B.12n signature is shown as a 240 × 24 matrix using a color scale, where each bin (*i,j*) corresponds to the system's discriminability of image *i* against distractor object *j*. Colors scales in **A**, **B** and **D** span each signature's full performance range and warm colors indicate lower discriminability. **E**, Human consistency of B.12n signatures for each of the tested model visual systems. Format is identical to that in Figure 2C.

object-level discriminability; see Fig. 3A), where warmer colors indicate lower values; color scales are not shown but span each signature's full range. Figure 3C shows the human consistency with respect to the B.11n signature for all tested models. Unlike with object-level behavioral metrics, we now observe a divergence between DCNN_{IC} models and primates. Both the monkey pool and the held-out human pool remain highly human consistent ($\bar{p} = 0.77, 0.96$ for monkeys, humans respectively), but all DCNN_{IC} models were significantly less human consistent (Inception-v3: $\bar{p} = 0.62$) and well outside of the defined "primate zone" of B.11n human consistency. Indeed, the hypothesis that the human consistency of DCNN_{IC} models is within the primate zone is strongly rejected ($p = 6.16 \times 10^{-8}$, exact test; see Materials and Methods).

We can zoom in further by examining not only the overall performance for a given image but also the object confusions for each image, that is, the additional behavioral variation that is due, not only to the test image, but also to the interaction of that test image with the alternative (incorrect) object choice that is pro-

vided after the test image (Fig. 1B). This is the highest level of behavioral accuracy resolution that our task design allows. In raw form, it corresponds to one-versus-other image-level confusion matrix, where the size of that matrix is the total number of images by the total number of objects (here, 240×24). Each bin (*i,j*) corresponds to the behavioral discriminability of a single image *i* against distractor object *j*. Again, we isolate variance that is not predicted by object-level performance by subtracting the average performance on this binary task (mean over all images) to convert the raw matrix B.12 above into the normalized matrix, referred to as B.12n. Figure 3D shows the B.12n signatures as 240×24 matrices for the pooled human, pooled monkey, and top DCNN_{IC} visual system models. Color scales are not shown but span each signature's full range; warmer colors correspond to images with lower performance in a given binary task relative to all images of that object in the same task. Figure 3E shows the human consistency with respect to the B.12n metric for all tested visual system models. Extending our observations using

B.I1n, we observe a similar divergence between primates and DCNN_{IC} visual system models on the matrix pattern of image-by-distractor difficulties (B.I2n). Specifically, both the monkey pool and held-out human pool remain highly human consistent ($\bar{p} = 0.75$, 0.77 for monkeys, humans respectively), whereas all tested DCNN_{IC} models are significantly less human consistent (Inception-v3: $\bar{p} = 0.53$) falling well outside of the defined “primate zone” of B.I2n human consistency values. Once again, the hypothesis that the human consistency of DCNN_{IC} models is within the primate zone is strongly rejected ($p = 3.17 \times 10^{-18}$, exact test; see Materials and Methods).

Natural subject-to-subject variation

For each behavioral metric (B.O1, B.O2, B.I1n, B.I2n), we defined a “primate zone” as the range of consistency values delimited by human consistency estimates $\bar{p}_{M\infty}$ and $\bar{p}_{H\infty}$ as lower and upper bounds respectively. $\bar{p}_{M\infty}$ corresponds to the extrapolated estimate of the human consistency of a large (i.e., infinitely many subjects) pool of rhesus macaque monkeys. Therefore, the fact that a particular tested visual system model falls outside of the primate zone can be interpreted as a failure of that visual system model to accurately predict the behavior of the archetypal human at least as well as the archetypal monkey.

However, from the above analyses, it is not yet clear whether a visual system model that fails to predict the archetypal human might nonetheless accurately correspond to one or more individual human subjects found within the natural variation of the human population. Given the difficulty of measuring individual subject behavior at the resolution of single images for large numbers of human and monkey subjects, we could not yet directly test this hypothesis. Instead, we examined it indirectly by asking whether an archetypal model—that is, a pool that includes an increasing number of model “subjects”—would approach the human pool. We simulated model intersubject variability by retraining a fixed DCNN architecture with a fixed training image set with random variation in the initial conditions and order of training images. This procedure results in models that can still perform the task but with slightly different learned weight values. We note that this procedure is only one possible choice of generating intersubject variability within each visual system model type, a choice that is an important open research direction that we do not address here. From this procedure, we constructed multiple trained model instances (“subjects”) for a fixed DCNN architecture and asked whether an increasingly large pool of model “subjects” better captures the behavior of the human pool at least as well as a monkey pool. This *post hoc* analysis was conducted for the most human consistent DCNN architecture (Inception-v3).

Figure 4A shows, for each of the four behavioral metrics, the measured human consistency of subject pools of varying size (number of subjects, n) of rhesus macaque monkeys (black) and ImageNet-trained Inception-v3 models (blue). The human consistency increases with growing number of subjects for both visual systems across all behavioral metrics. To estimate the expected human consistency for a pool of infinitely many monkey or model subjects, we fit an exponential function mapping n to the mean human consistency values and obtained a parameter estimate for the asymptotic value (see Materials and Methods). We note that estimated asymptotic values are not significantly beyond the range of the measured data; the human consistency of a pool of five monkey subjects reaches within 97% of the human consistency of an estimated infinite pool of monkeys for all metrics, giving credence to the extrapolated human consistency values. This analysis suggests that, under this model of intersubject

variability, a pool of Inception-v3 subjects accurately capture archetypal human behavior at the resolution of objects (B.O1, B.O2) by our primate zone criterion (Fig. 4A, first two panels). In contrast, even a large pool of Inception-v3 subjects still fails at its final asymptote to accurately capture human behavior at the image level (B.I1n, B.I2n) (Fig. 4A, last two panels).

Modification of visual system models to try to rescue their human consistency

Next, we wondered whether some relatively simple changes to the DCNN_{IC} visual system models tested here could bring them into better correspondence with the primate visual system behavior (with respect to B.I1n and B.I2n metrics). Specifically, we considered and tested the following modifications to the most human consistent DCNN_{IC} model visual system (Inception-v3): we did the following: (1) changed the input to the model to be more primate-like in its retinal sampling (Inception-v3 + retina-like), (2) changed the transformation (aka “decoder”) from the internal model feature representation into the behavioral output by augmenting the number of decoder training images or changing the decoder type (Inception-v3 + SVM, Inception-v3 + classifier_train), and (3) modified all of the internal filter weights of the model (aka “fine tuning”) by augmenting its ImageNet training with additional images drawn from the same distribution as our test images (Inception-v3 + synthetic-fine-tune), and (4) modified all of the internal filter weights of the model by training exclusively on images drawn from the same distribution as our test images (Inception-v3 + synthetic-train). All model modifications resulted in relatively high-performing models; Figure 5 (gray bars) shows the mean overall performance over object recognition tasks evaluated with a fixed decoder type (logistic classifier, 20 training images per class). However, we found that none of these modifications led to a significant improvement in its human consistency on the behavioral metrics (Fig. 4B). In particular, training exclusively on synthetic images led to a noted decrease in human consistency across all metrics. Figure 5 shows the relationship between model performances (B.O2, averaged over 276 tasks) and human consistency for both object-level and image-level metrics. We observe a strong correlation between performance and image-level human consistency across different model architectures trained on ImageNet (black points, $r = 0.97$, $p < 10^{-4}$), but no such correlation across our modified models within a fixed architecture (gray points, $r = 0.35$, $p = 0.13$). Therefore, the failure of current DCNN_{IC} models to accurately capture the image-level signatures of primates cannot be rescued by simple modifications on a fixed architecture.

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

Together, Figures 2, 3, 4, and 5 suggest that current DCNN_{IC} visual system models fail to accurately capture the image-level behavioral signatures of humans and monkeys. To further examine this failure in the hopes of providing clues for model improvement, we examined the image-level residual signatures of all the visual system models relative to the pooled human. For each model, we computed its residual signature as the difference (positive or negative) of a linear least-squares regression of the model signature on the corresponding human signature. For this analysis, we focused on the B.I1n metric because it showed a clear divergence of DCNN_{IC} models and primates and the behavioral residual can be interpreted based only on the test images (whereas B.I2n depends on the interaction between test images and distractor choice).

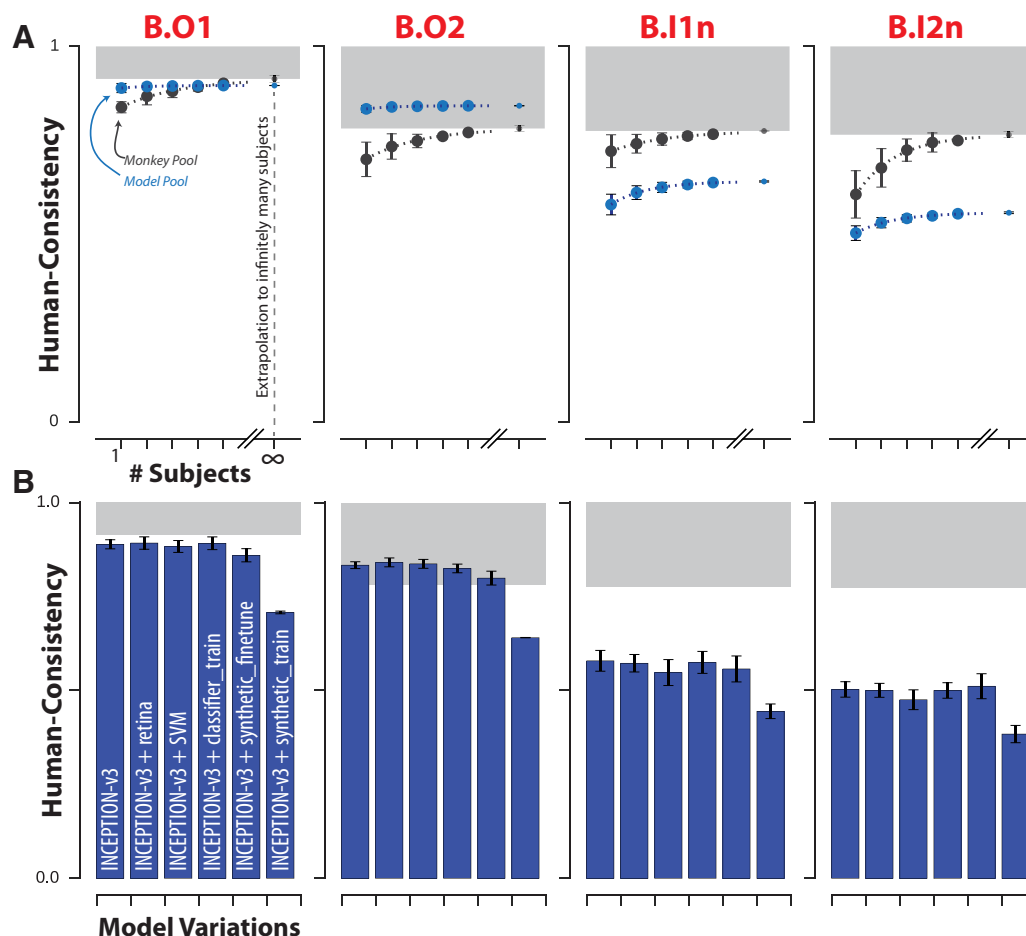


Figure 4. Effect of subject pool size and DCNN model modifications on consistency with human behavior. **A**, For each of the four behavioral metrics, the human consistency distributions of monkey (blue markers) and model (black markers) pools are shown as a function of the number of subjects in the pool (mean ± SD, over subjects). Human consistency increases with growing number of subjects for all visual systems across all behavioral metrics. The dashed lines correspond to fitted exponential functions and the parameter estimate (mean ± SE) of the asymptotic value, corresponding to the estimated human consistency of a pool of infinitely many subjects, is shown at the right most point on each abscissa. **B**, Model modifications that aim to rescue the DCNN_{IC} models. We tested several simple modifications (see Materials and Methods) to the most human consistent DCNN_{IC} visual system model (Inception-v3). Each panel shows the resulting human consistency per modified model (mean ± SD, over different model instances, varying in random filter initializations) for each of the four behavioral metrics.

We first asked to what extent the residual signatures are shared between different visual system models. Figure 6A shows the similarity between the residual signatures of all pairs of models; the color of bin (i, j) indicates the proportion of explainable variance that is shared between the residual signatures of visual systems i and j . For ease of interpretation, we ordered visual system models based on their architecture and optimization procedure and partitioned this matrix into four distinct regions. Each region compares the residuals of a “source” model group with fixed architecture and optimization procedure (five Inception-v3 models optimized for categorization on ImageNet, varying only in initial conditions and training image order) to a “target” model group. The target groups of models for each of the four regions are as follows: (1) the pooled monkey, (2) other DCNN_{IC} models from the source group, (3) DCNN_{IC} models that differ in architecture but share the optimization procedure of the source group models, and (4) DCNN_{IC} models that differ slightly using an augmented optimization procedure but share the architecture of the source group models. Figure 6B shows the mean (±SD) variance shared in the residuals averaged within these four regions for all images (black dots), as well as for images that humans found to be particularly difficult (gray dots, selected based on held-out human data, see Materials and Methods). First, consistent with

the results shown in Figure 3, we note that the residual signatures of this particular DCNN_{IC} model are not well shared with the pooled monkey ($r^2 = 0.39$ in region 1) and this phenomenon is more pronounced for the images that humans found most difficult ($r^2 = 0.17$ in region 1). However, this relatively low correlation between model and primate residuals is not indicative of spurious model residuals because the model residual signatures were highly reliable between different instances of this fixed DCNN_{IC} model across random training initializations (region 2: $r^2 = 0.79, 0.77$ for all and most difficult images, respectively). Interestingly, residual signatures were still largely shared with other DCNN_{IC} models with vastly different architectures (region 3: $r^2 = 0.70, 0.65$ for all and most difficult images, respectively). However, residual signatures were more strongly altered when the visual training diet of the same architecture was altered (region 4: $r^2 = 0.57, 0.46$ for all and most difficult images respectively, cf. region 3). Together, these results indicate that the images where DCNN_{IC} visual system models diverged from humans (and monkeys) were not spurious, but were rather highly reliable across different model architectures, demonstrating that current DCNN_{IC} models systematically and similarly diverge from primates.

Figure 7A shows example images sorted by B.I1n squared residuals, corresponding to images that models and primates

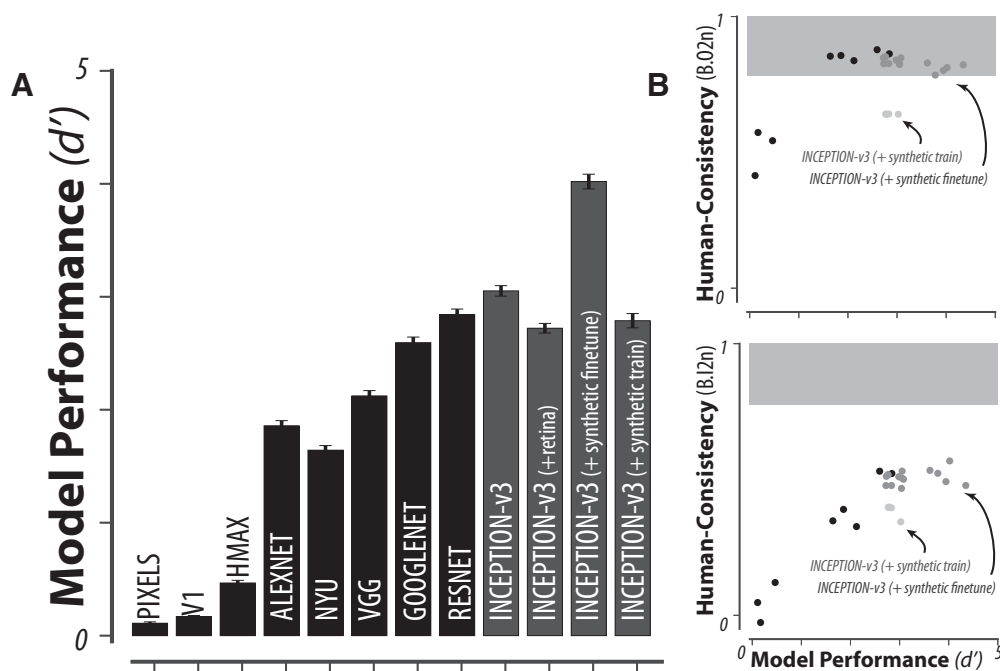


Figure 5. Model performance. **A**, Model performance on synthetic images (average B.02 across 276 tasks) for each of the tested models fixing the number of training images and classifier. Black bars correspond to different model architectures, with fixed optimization, whereas gray bars correspond to different modifications of a fixed model architecture (Inception-v3). **B**, Correlation between model performance and human consistency with respect to object-level (B.02) and image-level (B.12n) behavioral metrics. Each point corresponds to a single instance of a trained DCNN model.

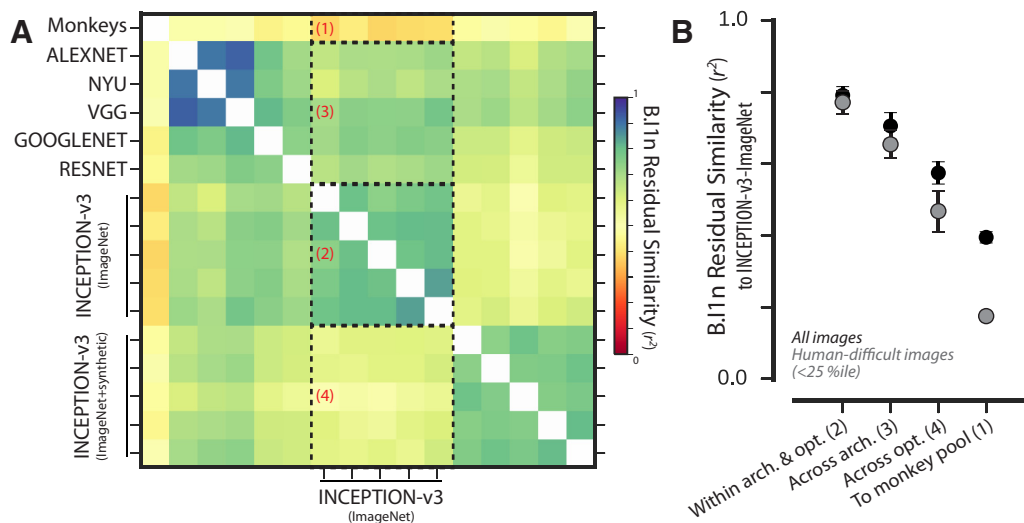


Figure 6. Analysis of unexplained human behavioral variance. **A**, Residual similarity between all pairs of human visual system models. The color of bin (i, j) indicates the proportion of explainable variance that is shared between the residual signatures of visual systems i and j . For ease of interpretation, we ordered visual system models based on their architecture and optimization procedure and partitioned this matrix into four distinct regions. **B**, Summary of residual similarity. For each of the four regions in Figure 5A, the similarity to the residuals of Inception-v3 (region 2 in A) is shown (mean \pm SD, within each region) for all images (black dots) and for images that humans found to be particularly difficult (gray dots, selected based on held-out human data).

largely agreed on (left) and diverged on (right) with respect to the B.11n metric. We observed no qualitative differences between these images. To look for clues for model improvement, we asked what, if any, characteristics of images might account for this divergence of models and primates. We regressed the residual signatures of DCNN_{IC} models on four different image attributes (corresponding to the size, eccentricity, pose, and contrast of the object in each image). We used multiple linear regressions to predict the model residual signatures from all of these image attributes and also considered each attribute individually using simple linear regressions. Figure 7B shows example images (sampled from the full set of 2400 images) with increasing attribute

value for each of these four image attributes. Although the DCNN_{IC} models were not directly optimized to display primate-like performance dependence on such attributes, we observed that the Inception-v3 visual system model nonetheless exhibited qualitatively similar performance dependencies as primates (Fig. 7C). For example, humans (black), monkeys (gray), and the Inception-v3 model (blue) all performed better, on average, for images in which the object is in the center of gaze (low eccentricity) and large in size. Furthermore, all three systems performed better, on average, for images when the pose of the object was closer to the canonical pose (Fig. 7C); this sensitivity to object pose manifested itself as a nonlinear dependence due to the fact

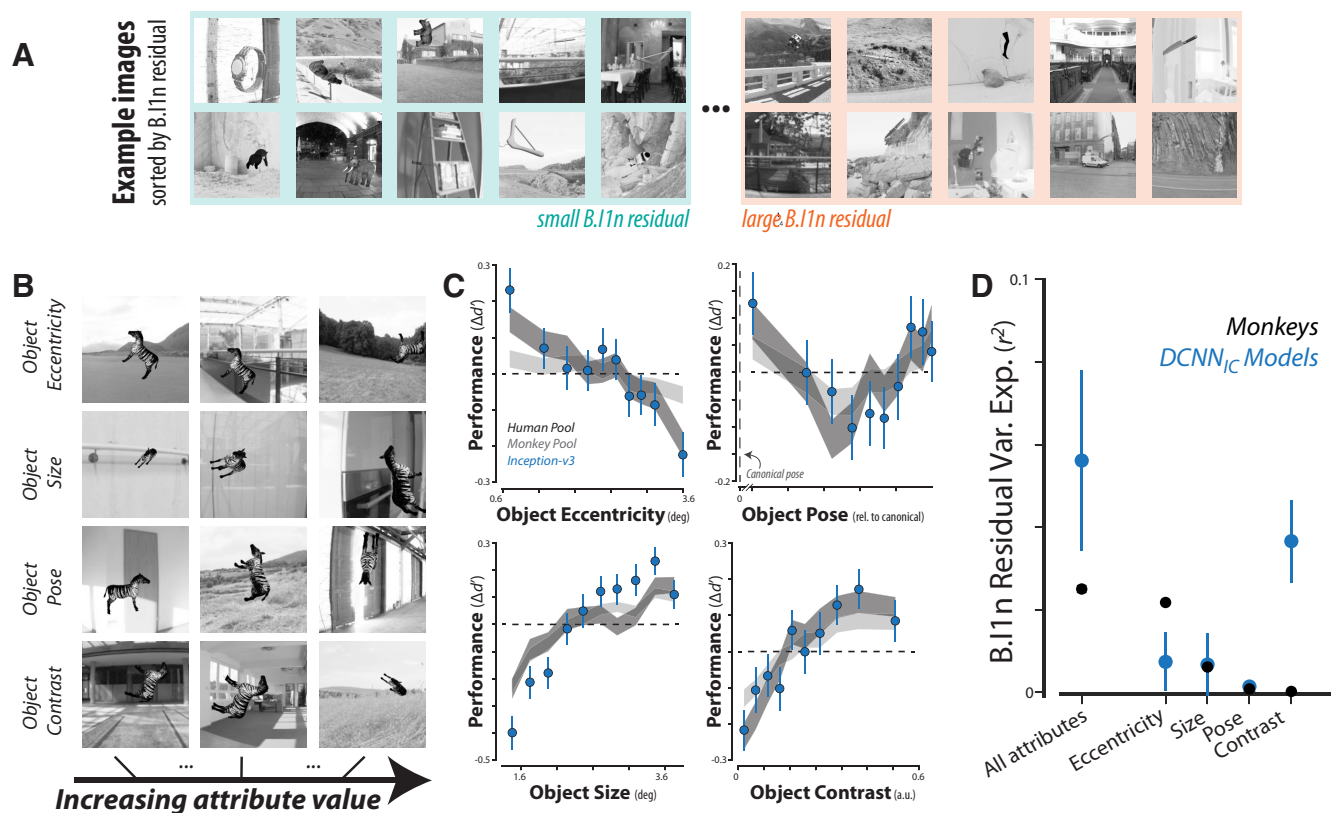


Figure 7. Dependence of primate and DCNN model behavior on image attributes. **A**, Example images showing that models and primates agree on (left) and diverge on (right) with respect to B.11n residuals. **B**, Example images with increasing attribute value for each of the four predefined image attributes (see Materials and Methods). **C**, Dependence of performance (B.11n) as a function of four image attributes for humans, monkeys, and a DCNN_{IC} model (Inception-v3). **D**, Proportion of explainable variance of the residual signatures of monkeys (black) and DCNN_{IC} models (blue) that is accounted for by each of the predefined image attributes. Error bars correspond to SD over trial resampling for monkeys and over different models for DCNN_{IC} models.

that all tested objects exhibited symmetry in at least one axis. The similarity of the patterns in Figure 7C between primates and the DCNN_{IC} visual system models is not perfect but is striking, particularly in light of the fact that these models were not optimized to produce these patterns. However, this similarity is analogous to the similarity in the B.O1 and B.O2 metrics in that it only holds on average over many images. Looking more closely at the image-by-image comparison, we again found that the DCNN_{IC} models failed to capture a large portion of the image-by-image variation (Fig. 3). In particular, Figure 7D shows the proportion of variance explained by specific image attributes for the residual signatures of monkeys (black) and DCNN_{IC} models (blue). We found that, together, all four of these image attributes explained only ~10% of the variance in DCNN_{IC} residual signatures and each individual attribute could explain at most a small amount of residual variance (<5% of the explainable variance). In sum, these analyses show that some behavioral effects that might provide intuitive clues to modify the DCNN_{IC} models are already in place in those models (e.g., a dependence on eccentricity). However, the quantitative image-by-image analyses of the remaining unexplained variance (Fig. 7D) argue that the DCNN_{IC} visual system models' failure to capture primate image-level signatures cannot be further accounted for by these simple image attributes and likely stem from other factors.

Discussion

The current work was motivated by the broad scientific goal of discovering models that quantitatively explain the neuronal mechanisms underlying primate invariant object recognition behavior. To this end, previous work had shown that specific ANNs

drawn from a large family of DCNNs and optimized to achieve high levels of object categorization performance on large-scale image sets capture a large fraction of the variance in primate visual recognition behaviors (Ghodrati et al., 2014; Rajalingham et al., 2015; Jozwik et al., 2016; Kheradpisheh et al., 2016; Kubilius et al., 2016; Peterson et al., 2016; Battleday et al., 2017; Wallis et al., 2017) and the internal hidden neurons of those same models also predict a large fraction of the image-driven response variance of brain activity at multiple stages of the primate ventral visual stream (Yamins et al., 2013, 2014; Cadieu et al., 2014; Khaligh-Razavi and Kriegeskorte, 2014; Güçlü and van Gerven, 2015; Cichy et al., 2016, 2017; Hong et al., 2016; Seibert et al., 2016; Cadena et al., 2017; Eickenberg et al., 2017; Khaligh-Razavi et al., 2017; Seeliger et al., 2017; Wen et al., 2017). For clarity, we here referred to this subfamily of models as DCNN_{IC} (to denote ImageNet categorization training) to distinguish them from all possible models in the DCNN family and, more broadly, from the superfamily of all ANNs. In this work, we directly compared leading DCNN_{IC} models to primates (humans and monkeys) with respect to their behavioral signatures at both object-level and image-level resolution in the domain of core object recognition. To do so, we measured and characterized primate behavior at larger scale and higher resolution than previously possible. We first replicate prior work (Rajalingham et al., 2015) showing that, at the object level, DCNN_{IC} models produce statistically indistinguishable behavior from primates and we extend that work by showing that these models also match the average primate sensitivities to object contrast, eccentricity, size, and pose, a noteworthy similarity in light of the fact that these models were not optimized

to produce these performance patterns. This similarity, which we speculate may largely reflect the training history of these models (i.e., optimization on photographer-framed images, which may be biased with respect to these attributes), is a good direction for future work. However, our primary novel result is that, examining behavior at the higher resolution of individual images, all leading DCNN_{IC} models failed to replicate the image-level behavioral signatures of primates. An important related claim is that rhesus monkeys are more consistent with the archetypal human than any of the tested DCNN_{IC} models (at the image level).

We compared human, monkey, and model behavior on a set of naturalistic synthetic images sampled with systematic variation in viewpoint parameters and uncorrelated background. Although these images may seem non-natural, synthetic images of this type are, unlike some photographic image sets, very powerful at separating humans from low-level computer vision systems (Pinto et al., 2008). Furthermore, given our definition of success (a candidate visual system model that accurately captures primate behavior for all images), the inference that models differ from primates can be supported by testing models on any image set, including the synthetic set that we used. Importantly, this inference is with respect to a trained model's behavior at "run time," agnostic to its learning procedure. Although we do not make any claims about learning, we note that neither humans nor monkeys were separately optimized to perform on our synthetic images (beyond the minimal training described in the Materials and Methods section). Interestingly, we observed that synthetic-image-optimized models (new ANN models constructed by fine-tuning or training an existing network architecture exclusively on large sets of synthetic image) were no more similar to primates than ANN models optimized only on ImageNet, suggesting that the tested ANN architectures have one or more fundamental flaws that cannot be readily overcome by manipulating the training environment. Together, these results support the general inference, agnostic to natural choices of image set, that DCNN_{IC} models diverge from primates in their core object recognition behavior.

This inference is consistent with previous work showing that DCNN_{IC} models can diverge from human behavior on specifically chosen adversarial images (Szegedy et al., 2013). A strength of our work is that we did not optimize images to induce failure, but instead randomly sampled a broadly defined image generative parameter space highlighting a general rather than adversarial-induced divergence of DCNN_{IC} from primate core object recognition behavior. Furthermore, we showed that this failure could not be rescued by a number of simple modifications to the model class, providing qualitative insight into what does and does not explain the gap between primates and DCNN models. Together, these results suggest that new ANN models are needed to more precisely capture the neural mechanisms underlying primate object vision.

With regard to new ANN models, we can attempt to make prospective inferences about future possible DCNN_{IC} models from the data presented here. Based on the observed distribution of image-level human consistency values for the DCNN_{IC} models tested here, we infer that as yet untested model instances sampled identically (i.e., from the DCNN_{IC} model subfamily) are very likely to have similarly inadequate image-level human consistency. Although we cannot rule out the possibility that at least one model instance within the DCNN_{IC} subfamily would fully match the image-level behavioral signatures, the probability of sampling such a model is vanishingly small ($p < 10^{-17}$ for B.12n human consistency, estimated using exact test using Gaussian kernel density estimation; see Materials and Methods and Results). An important

caveat of this inference is that we may have a biased estimate of the human consistency distribution of this model subfamily because we did not exhaustively sample the subfamily. In particular, if the model sampling process is nonstationary over time (e.g., increases in computational power over time allows larger models to be successfully trained), then the human consistency of new (i.e., yet to be sampled) models may lie outside the currently estimated distribution. Therefore, it is possible that the evolution of "next-generation" models within the DCNN_{IC} subfamily could meet our criteria for successful matching primate-like behavior.

Alternatively, it is possible—and we think likely—that future DCNN_{IC} models will also fail to capture primate-like image-level behavior, suggesting that either the architectural limitations (e.g., convolutional, feedforward) and/or the optimization procedure (including the diet of visual images) that define this model subfamily are fundamentally limiting. Therefore, ANN model subfamilies using different architectures (e.g., recurrent neural networks) and/or optimized for different behavioral goals (e.g., loss functions other than object classification performance, and/or images other than category-labeled ImageNet images) may be necessary to accurately capture primate behavior. To this end, we propose that testing even individual changes to the DCNN_{IC} models, each creating a new ANN model subfamily, may be the best way forward because DCNN_{IC} models currently offer the best explanations (in a predictive sense) of both the behavioral and neural phenomena of core object recognition.

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

References

- Battleday RM, Peterson JC, Griffiths TL (2017) Modeling human categorization of natural images using deep feature representations. arXiv preprint arXiv:1711.04855. Advance online publication. Retrieved Nov 13, 2017. Available at <https://arxiv.org/abs/1711.04855>.

- Cadena SA, Denfield GH, Walker EY, Gatys LA, Tolas AS, Bethge M, Ecker AS (2017) Deep convolutional models improve predictions of macaque V1 responses to natural images. *bioRxiv*:201764. [CrossRef](#)
- Cadiou 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 Comput Biol* 10:e1003963. [CrossRef](#) [Medline](#)
- 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. *Sci Rep* 6:27755. [CrossRef](#) [Medline](#)
- Cichy RM, Khosla A, Pantazis D, Oliva A (2017) Dynamics of scene representations in the human brain revealed by magnetoencephalography and deep neural networks. *Neuroimage* 153:346–358. [CrossRef](#) [Medline](#)
- DiCarlo JJ, Cox DD (2007) Untangling invariant object recognition. *Trends Cogn Sci* 11:333–341. [CrossRef](#)
- DiCarlo JJ, Johnson KO (1999) Velocity invariance of receptive field structure in somatosensory cortical area 3b of the alert monkey. *J Neurosci* 19:401–419. [CrossRef](#) [Medline](#)
- DiCarlo JJ, Zoccolan D, Rust NC (2012) How does the brain solve visual object recognition? *Neuron* 73:415–434. [CrossRef](#) [Medline](#)
- Dodge S, Karam L (2017) A study and comparison of human and deep learning recognition performance under visual distortions. *arXiv preprint arXiv:170502498*. Advance online publication. Retrieved May 6, 2017. Available at <https://arxiv.org/abs/1705.02498>
- Eickenberg M, Gramfort A, Varoquaux G, Thirion B (2017) Seeing it all: convolutional network layers map the function of the human visual system. *Neuroimage* 152:184–194. [CrossRef](#) [Medline](#)
- 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*. Advance online publication. Retrieved June 21, 2017. Available at <https://arxiv.org/abs/1706.06969>.
- Ghodrati M, Farzmaidi A, Rajaei K, Ebrahimpour R, Khaligh-Razavi S-M (2014) Feedforward object-vision models only tolerate small image variations compared to human. *Front Comput Neurosci* 8:74. [Medline](#)
- Goodfellow IJ, Shlens J, Szegedy C (2014) Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*. Advance online publication. Retrieved Dec 20, 2014. Available at <https://arxiv.org/abs/1412.6572>.
- Güçlü U, van Gerven MA (2015) Deep neural networks reveal a gradient in the complexity of neural representations across the ventral stream. *J Neurosci* 35:10005–10014. [CrossRef](#) [Medline](#)
- 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. *Nat Neurosci* 19:613–622. [CrossRef](#) [Medline](#)
- 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. [CrossRef](#)
- Johnson KO, Hsiao SS, Yoshioka T (2002) Neural coding and the basic law of psychophysics. *Neuroscientist* 8:111–121. [CrossRef](#) [Medline](#)
- 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. [CrossRef](#) [Medline](#)
- Khaligh-Razavi S-M, Kriegeskorte N (2014) Deep supervised, but not unsupervised, models may explain IT cortical representation. *PLoS Comput Biol* 10:e1003915. [CrossRef](#) [Medline](#)
- Khaligh-Razavi S-M, Henriksson L, Kay K, Kriegeskorte N (2017) Fixed versus mixed RSA: explaining visual representations by fixed and mixed feature sets from shallow and deep computational models. *Journal of Mathematical Psychology* 76:184–197. [CrossRef](#) [Medline](#)
- Kheradpisheh SR, Ghodrati M, Ganjtabesh M, Masquelier T (2016) Deep networks can resemble human feed-forward vision in invariant object recognition. *Sci Rep* 6:32672. [CrossRef](#) [Medline](#)
- Kriegeskorte N (2015) Deep neural networks: a new framework for modeling biological vision and brain information processing. *Annu Rev Vis Sci* 1:417–446. [CrossRef](#) [Medline](#)
- 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 Comput Biol* 12:e1004896. [CrossRef](#) [Medline](#)
- LeCun Y, Bengio Y, Hinton G (2015) Deep learning. *Nature* 521:436–444. [CrossRef](#) [Medline](#)
- Majaj NJ, Hong H, Solomon EA, DiCarlo JJ (2015) Simple learned weighted sums of inferior temporal neuronal firing rates accurately predict human core object recognition performance. *J Neurosci* 35:13402–13418. [CrossRef](#) [Medline](#)
- Nguyen A, Yosinski J, Clune J (2015) Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp 427–436.
- Peterson JC, Abbott JT, Griffiths TL (2016) Adapting deep network features to capture psychological representations. *arXiv preprint arXiv:160802164*. Advance online publication. Retrieved August 6, 2016. Available at <https://arxiv.org/abs/1608.02164>.
- Pinto N, Cox DD, DiCarlo JJ (2008) Why is real-world visual object recognition hard? *PLoS Comput Biol* 4:e27. [CrossRef](#) [Medline](#)
- Rajalingham R, Schmidt K, DiCarlo JJ (2015) Comparison of object recognition behavior in human and monkey. *J Neurosci* 35:12127–12136. [CrossRef](#) [Medline](#)
- RichardWebster B, Anthony SE, Scheirer WJ (2016) PsyPhy: a psychophysics driven evaluation framework for visual recognition. *arXiv preprint arXiv:161106448*. Advance online publication. Retrieved Nov 19, 2016. Available at <https://arxiv.org/abs/1611.06448>.
- Rolls ET (2000) Functions of the primate temporal lobe cortical visual areas in invariant visual object and face recognition. *Neuron* 27:205–218. [CrossRef](#) [Medline](#)
- Rosch E, Mervis CB, Gray WD, Johnson DM, Boyes-Braem P (1976) Basic objects in natural categories. *Cogn Psychol* 8:382–439. [CrossRef](#)
- Seeliger K, Fritsche M, Güçlü U, Schoenmakers S, Schoffelen J-M, Bosch S, van Gerven M (2017) Convolutional neural network-based encoding and decoding of visual object recognition in space and time. *Neuroimage pii: S1053-8119(17)30586-4*. [CrossRef](#) [Medline](#)
- Seibert D, Yamins DL, Ardila D, Hong H, DiCarlo JJ, Gardner JL (2016) A performance-optimized model of neural responses across the ventral visual stream. *bioRxiv*:036475. [CrossRef](#)
- Simonyan K, Zisserman A (2014) Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*. Advance online publication. Retrieved September 4, 2014. Available at <https://arxiv.org/abs/1409.1556>.
- Szegedy C, Zaremba W, Sutskever I, Bruna J, Erhan D, Goodfellow I, Fergus R (2013) Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*. Advance online publication. Retrieved December 21, 2013. Available at <https://arxiv.org/abs/1312.6199>.
- Tanaka K (1996) Inferotemporal cortex and object vision. *Annu Rev Neurosci* 19:109–139. [CrossRef](#) [Medline](#)
- Ullman S, Humphreys GW (1996) High-level vision: object recognition and visual cognition. Cambridge, MA: MIT.
- Wallis TS, Funke CM, Ecker AS, Gatys LA, Wichmann FA, Bethge M (2017) A parametric texture model based on deep convolutional features closely matches texture appearance for humans. *J Vis* 17:5. [CrossRef](#) [Medline](#)
- Wen H, Shi J, Zhang Y, Lu KH, Cao J, Liu Z (2017) Neural encoding and decoding with deep learning for dynamic natural vision. *Cereb Cortex* 20:1–25. [CrossRef](#) [Medline](#)
- Yamins DL, DiCarlo JJ (2016) Using goal-driven deep learning models to understand sensory cortex. *Nat Neurosci* 19:356–365. [CrossRef](#) [Medline](#)
- Yamins DL, Hong H, Cadiou C, DiCarlo JJ (2013) Hierarchical modular optimization of convolutional networks achieves representations similar to macaque IT and human ventral stream. In: *Advances in neural information processing systems*, pp 3093–3101.
- Yamins DL, Hong H, Cadiou CF, Solomon EA, Seibert D, DiCarlo JJ (2014) Performance-optimized hierarchical models predict neural responses in higher visual cortex. *Proc Natl Acad Sci U S A* 111:8619–8624. [CrossRef](#) [Medline](#)
- Zeiler MD, Fergus R (2014) Visualizing and understanding convolutional networks. In: *Computer Vision ECCV 2014*, pp 818–833. New York, NY: Springer.

Research Articles: Behavioral/Cognitive

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

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

DOI: 10.1523/JNEUROSCI.0388-18.2018

Received: 12 February 2018

Revised: 6 June 2018

Accepted: 8 July 2018

Published: 13 July 2018

Author contributions: R.R., E.B.I., and J.J.D. designed research; R.R., E.B.I., K.K., and K.S. performed research; R.R. and P.B. analyzed data; R.R., E.B.I., and J.J.D. wrote the paper.

Conflict of Interest: The authors declare no competing financial interests.

This research was supported by US National Eye Institute grants R01-EY014970 (J.J.D.) and K99-EY022671 (E.B.I.), Office of Naval Research MURI-114407 (J.J.D.), IARPA Contract D16PC00002 (J.J.D.), Engineering Research Council of Canada (R.R.), and The McGovern Institute for Brain Research

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

Cite as: J. Neurosci ; 10.1523/JNEUROSCI.0388-18.2018

Alerts: Sign up at www.jneurosci.org/cgi/alerts to receive customized email alerts when the fully formatted version of this article is published.

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.

ACKNOWLEDGEMENTS

This research was supported by US National Eye Institute grants R01-EY014970 (J.J.D.) and K99-EY022671 (E.B.I.), Office of Naval Research MURI-114407 (J.J.D.), IARPA Contract D16PC00002 (J.J.D.), Engineering Research Council of Canada (R.R.), and The McGovern Institute for Brain Research

COMPETING FINANCIAL INTERESTS

The authors declare no competing financial interests.

ABSTRACT

Primates—including humans—can typically recognize objects in visual images at a glance even in spite of naturally occurring identity-preserving image transformations (e.g. changes in viewpoint). A primary neuroscience goal is to uncover neuron-level mechanistic models that quantitatively explain this behavior by predicting primate performance for each and every image. Here, we applied this stringent behavioral prediction test to the leading mechanistic models of primate vision (specifically, deep, convolutional, artificial neural networks; ANNs) by directly comparing their behavioral signatures against those of humans and rhesus macaque monkeys. Using high-throughput data collection systems for human and monkey psychophysics, we collected over one million behavioral trials from 1472 anonymous humans and five male macaque monkeys for 2400 images over 276 binary object discrimination tasks. Consistent with previous work, we observed that state-of-the-art deep, feed-forward convolutional ANNs trained for visual categorization (termed DCNN_{IC} models) accurately predicted primate patterns of object-level confusion. However, when we examined behavioral performance for individual images within each object discrimination task, we found that all tested DCNN_{IC} models were significantly non-predictive of primate performance, and that this prediction failure was not accounted for by simple image attributes, nor rescued by simple model modifications. These results show that current DCNN_{IC} models cannot account for the image-level behavioral patterns of primates, and that new ANN models are needed to more precisely capture the neural mechanisms underlying primate object vision. To this end, large-scale, high-resolution primate behavioral benchmarks—such as those obtained here—could serve as direct guides for discovering such models.

SIGNIFICANCE STATEMENT

Recently, specific feed-forward deep convolutional artificial neural networks (ANNs) models have dramatically advanced our quantitative understanding of the neural mechanisms underlying primate core object recognition. In this work, we tested the limits of those ANNs by systematically comparing the behavioral responses of these models with the behavioral responses of humans and monkeys, at the resolution of individual images. Using these high-resolution metrics, we found that all tested ANN models significantly diverged from primate behavior. Going forward, these high-resolution, large-scale primate behavioral benchmarks could serve as direct guides for discovering better ANN models of the primate visual system.

INTRODUCTION

Primates—both human and non-human—can typically recognize objects in visual images at a glance, even in the face of naturally occurring identity-preserving transformations such as changes in viewpoint. This view-invariant visual object recognition ability is thought to be supported primarily by the primate ventral visual stream (Tanaka, 1996; Rolls, 2000; DiCarlo et al., 2012). A primary neuroscience goal is to construct computational models that quantitatively explain the neural mechanisms underlying this ability. That is, our goal is to discover artificial neural networks (ANNs) that accurately predict neuronal firing rate responses at all levels of the ventral stream and its behavioral output. To this end, specific models within a large family of deep, convolutional neural networks (DCNNs), optimized by supervised training on large-scale category-labeled image-sets (ImageNet) to match human-level categorization performance (Krizhevsky et al., 2012; LeCun et al., 2015), have been put forth as the leading ANN models of the ventral stream (Kriegeskorte, 2015; Yamins and DiCarlo, 2016). We refer to this sub-family as DCNN_{IC} models (IC to denote ImageNet-categorization pre-training), so as to distinguish them from all possible models in the DCNN family, and more broadly, from the super-family of all ANNs. To date, it has been shown that DCNN_{IC} models display internal feature representations similar to neuronal representations along the primate ventral visual stream (Yamins et al., 2013; Cadieu et al., 2014; Khaligh-Razavi and Kriegeskorte, 2014; Yamins et al., 2014), and they exhibit behavioral patterns similar to the behavioral patterns of pairwise object confusions of primates (Ghodrati et al., 2014; Rajalingham et al., 2015; Jozwik et al., 2016; Kheradpisheh et al., 2016). Thus, DCNN_{IC} models may provide a quantitative account of the neural mechanisms underlying primate core object recognition behavior.

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

confusions mentioned above) that are computed by pooling over hundreds of images and thus are not sensitive to variation in difficulty across images of the same object. To overcome this limitation of prior work, we here aimed to use scalable behavioral testing methods to precisely characterize primate behavior at the resolution of individual images and to directly compare leading DCNN models to primates over the domain of core object recognition behavior at this high resolution.

We focused on *core invariant object recognition*—the ability to identify objects in visual images in the central visual field during a single, natural viewing fixation (DiCarlo et al., 2012). We further restricted our behavioral domain to *basic-level* object discriminations, as defined previously (Rosch et al., 1976). Within this domain, we collected large-scale, high-resolution measurements of human and monkey behavior (over a million behavioral trials) using high-throughput psychophysical techniques—including a novel home-cage behavioral system for monkeys. These data enabled us to systematically compare all systems at progressively higher resolution. At lower resolutions, we replicated previous findings that humans, monkeys, and DCNN_{IC} models all share a common pattern of object-level confusion (Rajalingham et al., 2015). However, at the higher resolution of individual images, we found that the behavior of all tested DCNN_{IC} models was significantly different from human and monkey behavior, and this model prediction failure could not be easily rescued by simple model modifications. These results show that current DCNN_{IC} models do not fully account for the image-level behavioral patterns of primates, suggesting that new ANN models are needed to more precisely capture the neural mechanisms underlying primate object vision. To this end, large-scale high-resolution behavioral benchmarks, such as those obtained here, could serve as a strong top-down constraint for efficiently discovering such models.

MATERIALS & METHODS

Visual images

We examined basic-level, core object recognition behavior using a set of 24 broadly-sampled objects that we previously found to be reliably labeled by independent human subjects, based on the definition of basic-level proposed by (Rosch et al., 1976). For each object, we

145 generated 100 naturalistic synthetic images by first rendering a 3D model of the object with
 146 randomly chosen viewing parameters (2D position, 3D rotation and viewing distance), and then
 147 placing that foreground object view onto a randomly chosen, natural image background. To do
 148 this, each object was first assigned a canonical position (center of gaze), scale (~ 2 degrees) and
 149 pose, and then its viewing parameters were randomly sampled uniformly from the following
 150 ranges for object translation ($[-3, 3]$ degrees in both h and v), rotation ($[-180, 180]$ degrees in all
 151 three axes) and scale ($[x0.7, x1.7]$). Background images were sampled randomly from a large
 152 database of high-dynamic range images of indoor and outdoor scenes obtained from Dosch
 153 Design (www.doschdesign.com). This image generation procedure enforces invariant object
 154 recognition, rather than image matching, as it requires the visual recognition system (human,
 155 animal or model) to tackle the “invariance problem,” the computational crux of object
 156 recognition (Ullman and Humphreys, 1996; Pinto et al., 2008). In particular, we used naturalistic
 157 synthetic images with systematic variation in viewpoint parameters and uncorrelated background
 158 to remove potential confounds (natural images are often “composed” such that backgrounds co-
 159 vary with the object category) while keeping the task difficult for machine vision systems. We
 160 have previously shown that, unlike some photographic image sets, synthetic images of the types
 161 we used here are critical to separating some types of computer vision systems from humans
 162 (Pinto et al., 2008). Using this procedure, we previously generated 2400 images (100 images per
 163 object) rendered at 1024x1024 pixel resolution with 256-level gray scale and subsequently
 164 resized to 256x256 pixel resolution for human psychophysics, monkey psychophysics and model
 165 evaluation (Rajalingham et al., 2015). In the current work, we focused our analyses on a
 166 randomly subsampled, and then fixed, sub-set of 240 images (10 images per object; here referred
 167 to as the “primary test images”). Figure 1A shows the full list of 24 objects, with two example
 168 images of each object.

169
 170 Because all of the images were generated from synthetic 3D object models, we had
 171 explicit knowledge of the viewpoint parameters (position, size, and pose) for each object in each
 172 image, as well as perfect segmentation masks. Taking advantage of this feature, we characterized
 173 each image based on these high-level attributes, consisting of size, eccentricity, relative pose and
 174 contrast of the object in the image. Note that these meta-parameters were independently
 175 randomly sampled to generate each image, and thus there is no correlation between size,

eccentricity and pose over images. The size and eccentricity of the object in each image were computed directly from the corresponding viewpoint parameters, under the assumption that the entire image would subtend 6° at the center of visual gaze ($\pm 3^\circ$ in both azimuth and elevation; see below). For each synthetic object, we first defined its “canonical” 3D pose vector, based on independent human judgments. To compute the relative pose attribute of each image, we estimated the difference between the object’s 3D pose and its canonical 3D pose. Pose differences were computed as distances in unit quaternion representations: the 3D pose (r_{xy} , r_{xz} , r_{yz}) was first converted into unit quaternions, and distances between quaternions q_1 , q_2 were estimated as $\cos^{-1}|q_1 \cdot q_2|$ (Huynh, 2009). To compute the object contrast, we measured the absolute difference between the mean of the pixel intensities corresponding to the object and the mean of the background pixel intensities in the vicinity of the object (specifically, within 25 pixels of any object pixel, analogous to computing the local foreground-background luminance difference of a foreground object in an image). Figure 5A shows example images with varying values for the four image attributes.

Core object recognition behavioral paradigm

Core object discrimination is defined as the ability to discriminate between two or more objects in visual images presented under high view uncertainty in the central visual field ($\sim 10^\circ$), for durations that approximate the typical primate, free-viewing fixation duration (~ 200 ms) (DiCarlo and Cox, 2007; DiCarlo et al., 2012). As in our previous work (Rajalingham et al., 2015), the behavioral task paradigm consisted of a interleaved set of binary discrimination tasks. Each binary discrimination task is an object discrimination task between a pair of objects (e.g. elephant vs. bear). Each such binary task is balanced in that the test image is equally likely (50%) to be of either of the two objects. On each trial, a test image is presented, followed by a choice screen showing canonical views of the two possible objects (the object that was not displayed in the test image is referred to as the “distractor” object, but note that objects are equally likely to be distractors and targets). Here, 24 objects were tested, which resulted in 276 binary object discrimination tasks. To neutralize feature attention, these 276 tasks are randomly interleaved (trial by trial), and the global task is referred to as a basic-level, core object recognition task paradigm.

206

207 *Testing human behavior*

208 All human behavioral data presented here were collected from 1476 anonymous human
 209 subjects on Amazon Mechanical Turk (MTurk) performing the task paradigm described above.
 210 Subjects were instructed to report the identity of the foreground object in each presented image
 211 from among the two objects presented on the choice screen (Fig 1B). Because all 276 tasks were
 212 interleaved randomly (trial-by-trial), subjects could not deploy feature attentional strategies
 213 specific to each object or specific to each binary task to process each test image.

214

215 Figure 1B illustrates the time course of each behavioral trial, for a particular object
 216 discrimination task (zebra versus dog). Each trial initiated with a central black point for 500 ms,
 217 followed by 100 ms presentation of a test image containing one foreground object presented
 218 under high variation in viewing parameters and overlaid on a random background, as described
 219 above (see *Visual images* above). Immediately after extinction of the test image, two choice
 220 images, each displaying a single object in a canonical view with no background, were shown to
 221 the left and right. One of these two objects was always the same as the object that generated the
 222 test image (i.e., the correct object choice), and the location of the correct object (left or right) was
 223 randomly chosen on each trial. After clicking on one of the choice images, the subject was
 224 queued with another fixation point before the next test image appeared. No feedback was given;
 225 human subjects were never explicitly trained on the tasks. Under assumptions of typical
 226 computer ergonomics, we estimate that images were presented at 6–8° of visual angle at the
 227 center of gaze, and the choice object images were presented at ± 6 –8° of eccentricity along the
 228 horizontal meridian.

229

230 We measured human behavior using the online Amazon MTurk platform (see Figure 1C),
 231 which enables efficient collection of large-scale psychophysical data from crowd-sourced
 232 “human intelligence tasks” (HITs). The reliability of the online MTurk platform has been
 233 validated by comparing results obtained from online and in-lab psychophysical experiments
 234 (Majaj et al., 2015; Rajalingham et al., 2015). We pooled 927,296 trials from 1472 human
 235 subjects to characterize the aggregate human behavior, which we refer to as the “pooled” human
 236 (or “archetypal” human). Each human subject performed only a small number of trials (~150) on

a subset of the images and binary tasks. All 2400 images were used for behavioral testing, but in some of the HITs, we biased the image selection towards the 240 primary test images (1424 ± 70 trials/image on this subsampled set, versus 271 ± 93 trials/image on the remaining images, mean \pm SD) to efficiently characterize behavior at image level resolution. Images were randomly drawn such that each human subject was exposed to each image a relatively small number of times (1.5 ± 2.0 trials/image per subject, mean \pm SD), in order to mitigate potential alternative behavioral strategies (e.g. “memorization” of images) that could arise from a finite image set. Behavioral signatures at the object-level (B.O1, B.O2, see *Behavioral metrics and signatures*) were measured using all 2400 test images, while image-level behavioral signatures (B.I1n, B.I2n, see *Behavioral metrics and signatures*) were measured using the 240 primary test images. (We observed qualitatively similar results using those metrics on the full 2400 test images, but we here focus on the primary test images as the larger number of trials leads to lower noise levels).

Five other human subjects were separately recruited on MTurk to each perform a large number of trials on the same images and tasks ($53,097 \pm 15,278$ trials/subject, mean \pm SD). Behavioral data from these five subjects was not included in the characterization of the pooled human described above, but instead aggregated together to characterize a distinct held-out human pool. For the scope of the current work, this held-out human pool—which largely replicated all behavioral signatures of the larger archetypal human (see Figures 2 and 3)—served as an independent validation of our human behavioral measurements.

Testing monkey behavior

Five adult male rhesus macaque monkeys (*Macaca mulatta*, subjects *M*, *Z*, *N*, *P*, *B*) were tested on the same basic-level, core object recognition task paradigm described above, with minor modification as described below. All procedures were performed in compliance with National Institutes of Health guidelines and the standards of the Massachusetts Institute of Technology Committee on Animal Care and the American Physiological Society. To efficiently characterize monkey behavior, we used a novel home-cage behavioral system developed in our lab (termed MonkeyTurk, see Fig. 1C). This system leveraged a tablet touchscreen (9” Google Nexus or 10.5” Samsung Galaxy Tab S) and used a web application to wirelessly load the task and collect the data (code available from private github repository upon request). Analogous to

the online Amazon Mechanical Turk, which allows for efficient psychophysical assays of a large number (hundreds) of human users in their native environments, MonkeyTurk allowed us to test many monkey subjects simultaneously in their home environment. Each monkey voluntarily initiated trials, and each readily performed the task a few hours each day that the task apparatus was made available to it. At an average rate of ~2,000 trials per day per monkey, we collected a total of 836,117 trials from the five monkey subjects over a period of ~3 months.

Monkey training is described in detail elsewhere (Rajalingham et al., 2015). Briefly, all monkeys were initially trained on the match-test-image-to-object rule using other images and were also trained on discriminating the particular set of 24 objects tested here using a separate set of training images rendered from these objects, in the same manner as the main testing images. Two of the monkeys subjects (Z and M) were previously trained in the lab setting, and the remaining three subjects were trained using MonkeyTurk directly in their home cages and did not have significant prior lab exposure. Once monkeys reached saturation performance on training images, we began the behavioral testing phase to collect behavior on test images. Monkeys did improve throughout the testing phase, exhibiting an increase in performance between the first and second half of trials of $4\% \pm 0.9\%$ (mean \pm SEM over five monkey subjects). However, the image-level behavioral signatures obtained from the first and the second halves of trials were highly correlated to each other (B.II noise-adjusted correlation of 0.85 ± 0.06 , mean \pm SEM over five monkey subjects, see *Behavioral metrics and signatures* below), suggesting that monkeys did not significantly alter strategies (e.g. did not “memorize” images) throughout the behavioral testing phase.

The monkey task paradigm was nearly identical to the human paradigm (see Figure 1B), with the exception that trials were initiated by touching a white “fixation” circle horizontally centered on the bottom third of the screen (to avoid occluding centrally-presented test images with the hand). This triggered a 100ms central presentation of a test image, followed immediately by the presentation of the two choice images (Fig. 1B, location of correct choice randomly assigned on each trial, identical to the human task). Unlike the main human task, monkeys responded by directly touching the screen at the location of one of the two choice images. Touching the choice image corresponding to the object shown in the test image resulted

299 in the delivery of a drop of juice through a tube positioned at mouth height (but not obstructing
 300 view), while touching the distractor choice image resulted in a three second timeout. Because
 301 gaze direction typically follows the hand during reaching movements, we assumed that the
 302 monkeys were looking at the screen during touch interactions with the fixation or choice targets.
 303 In both the lab and in the home cage, we maintained total test image size at ~6 degrees of visual
 304 angle at the center of gaze, and we took advantage of the retina-like display qualities of the tablet
 305 by presenting images pixel matched to the display (256 x 256 pixel image displayed using 256 x
 306 256 pixels on the tablet at a distance of 8 inches) to avoid filtering or aliasing effects.

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

323 324 *Testing model behavior*

325 We tested a number of different deep convolutional neural network (DCNN) models on
 326 the exact same images and tasks as those presented to humans and monkeys. Importantly, our
 327 core object recognition task paradigm is closely analogous to the large-scale ImageNet 1000-way
 328 object categorization task for which these networks were optimized and thus expected to perform
 329 well. We focused on publicly available DCNN model architectures that have proven highly

330 successful with respect to this computer vision benchmark over the past five years: AlexNet
 331 (Krizhevsky et al., 2012), NYU (Zeiler and Fergus, 2014), VGG (Simonyan and Zisserman,
 332 2014), GoogleNet (Szegedy et al., 2013), Resnet (He et al., 2016), and Inception-v3 (Szegedy et
 333 al., 2013). As this is only a subset of possible DCNN models, we refer to these as the DCNN_{IC}
 334 (to denote ImageNet-Categorization) visual system model sub-family. For each of the publicly
 335 available model architectures, we first used ImageNet-categorization-trained model instances,
 336 either using publicly available trained model instances or training them to saturation on the 1000-
 337 way classification task in-house. Training took several days on 1-2 GPUs.

338
 339 We then performed psychophysical experiments on each ImageNet-trained DCNN model
 340 to characterize their behavior on the exact same images and tasks as humans and monkeys. We
 341 first adapted these ImageNet-trained models to our 24-way object recognition task by re-training
 342 the final class probability layer (initially corresponding to the probability output of the 1000-way
 343 ImageNet classification task) while holding all other layers fixed. In practice, this was done by
 344 extracting features from the penultimate layer of each DCNN_{IC} (i.e. top-most prior to class
 345 probability layer), on the same images that were presented to humans and monkeys, and training
 346 back-end multi-class logistic regression classifiers to determine the cross-validated output class
 347 probability for each image. This procedure is illustrated in Figure 1C. To estimate the hit rate of
 348 a given image in a given binary classification task, we renormalized the 24-way class
 349 probabilities of that image, considering only the two relevant classes, to sum to one. Object-level
 350 and image-level behavioral metrics were computed based on these hit rate estimates (as
 351 described in *Behavioral metrics and signatures* below). Importantly, this procedure assumes that
 352 the model “retina” layer processes the central 6 degrees of the visual field. It also assumes that
 353 linear discriminants (“readouts”) of the model’s top feature layer are its behavioral output (as
 354 intended by the model designers). Manipulating either of these choices (e.g. resizing the input
 355 images such that they span only part of the input layer, or building linear discriminates for
 356 behavior using a different model feature layer) would result in completely new, testable ANN
 357 models that we do not test here.

358
 359 From these analyses, we selected the most *human-consistent* DCNN_{IC} architecture
 360 (Inception-v3, see *Behavioral consistency* below), fixed that architecture, and then performed

361 post-hoc analyses in which we varied: the input image sampling, the initial parameter settings
 362 prior to training, the filter training images, the type of classifiers used to generate the behavior
 363 from the model features, and the classifier training images. To examine input image sampling,
 364 we re-trained the Inception-v3 architecture on images from ImageNet that were first spatially
 365 filtered to match the spatial sampling of the primate retina (i.e. an approximately exponential
 366 decrease in cone density away from the fovea) by effectively simulating a fish-eye
 367 transformation on each image. These images were at highest resolution at the “fovea” (i.e. center
 368 of the image) with gradual decrease in resolution with increasing eccentricity. To examine the
 369 analog of “inter-subject variability”, we constructed multiple trained model instances
 370 (“subjects”), where the architecture and training images were held fixed (Inception-v3 and
 371 ImageNet, respectively) but the model filter weights initial condition and order of training
 372 images were randomly varied for each model instance. Importantly, this procedure is only one
 373 possible choice for simulating inter-subject variability for DCNN models, which, as a first order
 374 approximation, models different human subjects as random samples from a fixed model class.
 375 There are many other possible options for simulating inter-subject variability, including 1)
 376 random sampling of different features from a fixed trained model of fixed architecture, 2)
 377 random sampling of different trained models from a fixed architecture and optimization, 3)
 378 random sampling of different trained models from a model class, varying in either architecture
 379 (3a) or optimization procedure and data (3b) or both (3c), to name a few. This choice is an
 380 important open research direction that we do not address here. We do not claim this to be the
 381 best model of inter-subject variability, but rather a good starting point for that greater goal.

382 To examine the effect of model training, we first fine-tuned an ImageNet-trained
 383 Inception-v3 model on a synthetic image set consisting of ~6.9 million images of 1049 objects
 384 (holding out 50,000 images for model validation). These images were generated using the same
 385 rendering pipeline as our test images, but the objects were non-overlapping with the 24 test
 386 objects presented here. As expected, fine-tuning on synthetic images led to a small increase in
 387 overall performance (see Figure 5). To push the limits of training on synthetic images, we
 388 additionally trained an Inception-v3 architecture exclusively on this synthetic image set. To
 389 avoid over-fitting on this synthetic image set, we stopped the model training based on maximum
 390 performance on a validation set of held-out objects. These synthetic-trained models were high
 391 performing, with only a small decrease in overall performance relative to the ImageNet trained

models (see Figure 5, last bar). Finally, we tested the effect of different classifiers to generate model behavior by testing both multi-class logistic regression and support vector machine classifiers, as well as the effect of varying the number of training images used to train those classifiers (20 versus 50 images per class).

Behavioral metrics and signatures

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

The four behavioral metrics we chose are as follows: First, the one-versus-all object-level performance metric (termed B.O1) estimates the discriminability of each object from all other objects, pooling across all distractor object choices. Since we here tested 24 objects, the resulting B.O1 signature has 24 independent values. Second, the one-versus-other object-level performance metric (termed B.O2) estimates the discriminability of each specific pair of objects, or the pattern of pairwise object confusions. Since we here tested 276 interleaved binary object discrimination tasks, the resulting B.O2 signature has 276 independent values (the off-diagonal elements on one half of the 24x24 symmetric matrix). Third, the one-versus-all image-level performance metric (termed B.I1) estimates the discriminability of each image from all other objects, pooling across all possible distractor choices. Since we here focused on the primary image test set of 240 images (10 per object, see above), the resulting B.I1 signature has 240 independent values. Fourth, the one-versus-other image-level performance metric (termed B.I2)

estimates the discriminability of each image from each distractor object. Since we here focused on the primary image test set of 240 images (10 per object, see above) with 23 distractors, the resulting B.I2 signature has 5520 independent values.

Naturally, object-level and image-level behavioral signatures are tightly linked. For example, images of a particularly difficult-to-discriminate object would inherit lower performance values on average as compared to images from a less difficult-to-discriminate object. To isolate the behavioral variance that is specifically driven by image variation and not simply predicted by the objects (and thus already captured by B.O1 and B.O2), we defined normalized image-level behavioral metrics (termed B.I1n, B.I2n) by subtracting the mean performance values over all images of the same object and task. This process is schematically illustrated in Figure 3A. We note that the resulting normalized image-level behavioral signatures capture a significant proportion of the total image-level behavioral variance in our data (e.g. 52%, 58% of human B.I1 and B.I2 variance is driven by image variation, independent of object identity). In this study, we use these normalized metrics for image-level behavioral comparisons between models and primates (see Results).

Behavioral Consistency

To quantify the similarity between a model visual system and the human visual system with respect to a given behavioral metric, we used a measure called the “*human-consistency*” as previously defined (Johnson et al., 2002). *Human-consistency* ($\tilde{\rho}$) is computed, for each of the four behavioral metrics, as a noise-adjusted correlation of behavioral signatures (DiCarlo and Johnson, 1999). For each visual system, we randomly split all behavioral trials into two equal halves and applied each behavioral metric to each half, resulting in two independent estimates of the system’s behavioral signature with respect to that metric. We took the Pearson correlation between these two estimates of the behavioral signature as a measure of the reliability of that behavioral signature given the amount of data collected, i.e. the split-half internal reliability. To estimate the *human-consistency*, we computed the Pearson correlation over all the independent estimates of the behavioral signature from the model (\mathbf{m}) and the human (\mathbf{h}), and we then divide that raw Pearson correlation by the geometric mean of the split-half internal reliability of the same behavioral signature measured for each system: $\tilde{\rho}(\mathbf{m}, \mathbf{h}) = \frac{\rho(\mathbf{m}, \mathbf{h})}{\sqrt{\rho(\mathbf{m}, \mathbf{m})\rho(\mathbf{h}, \mathbf{h})}}$.

454

455 Since all correlations in the numerator and denominator were computed using the same
 456 amount of trial data (exactly half of the trial data), we did not need to make use of any prediction
 457 formulas (e.g. extrapolation to larger number of trials using Spearman-Brown prediction
 458 formula). This procedure was repeated 10 times with different random split-halves of trials. Our
 459 rationale for using a reliability-adjusted correlation measure for *human-consistency* was to
 460 account for variance in the behavioral signatures that arises from “noise,” i.e., variability that is
 461 not replicable by the experimental condition (image and task) and thus that no model can be
 462 expected to predict (DiCarlo and Johnson, 1999; Johnson et al., 2002). In sum, if the model (m)
 463 is a replica of the archetypal human (h), then its expected human-consistency is 1.0, regardless of
 464 the finite amount of data that are collected. Note that the human-consistency value is directly
 465 linked, via a squaring operation, to the proportion of the explainable behavioural variance that is
 466 explained by models.

467

468 *Characterization of Residuals*

469 In addition to measuring the similarity between the behavioral signatures of primates and
 470 models (using *human-consistency* analyses, as described above), we examined the corresponding
 471 differences, termed “residual signatures.” Each candidate visual system model’s residual
 472 signature was estimated as the residual of a linear least squares regression of the model’s
 473 signature on the corresponding human signature and a free intercept parameter. This procedure
 474 effectively captures the differences between human and model signatures after accounting for
 475 overall performance differences. Residual signatures were estimated on disjoint split-halves of
 476 trials, repeating 10 times with random trial permutations. Residuals were computed with respect
 477 to the normalized one-versus-all image-level performance metric (B.IIn) as this metric showed a
 478 clear difference between DCNN_{IC} models and primates, and the behavioral residual can be
 479 interpreted based only the test images (i.e. we can assign a residual per image).

480

481 To examine the extent to which the difference between each model and the archetypal
 482 human is reliably shared across different models, we measured the Pearson correlation between
 483 the residual signatures of pairs of models. Residual similarity was quantified as the proportion of
 484 shared variance, defined as the square of the noise-adjusted correlation between residual

signatures (the noise-adjustment was done as defined in equation above). Correlations of residual signatures were always computed across distinct split-halves of data, to avoid introducing spurious correlations from subtracting common noise in the human data. We measured the residual similarity between all pairs of tested models, holding both architecture and optimization procedure fixed (between instances of the ImageNet-categorization trained Inception-v3 model, varying in filter initial conditions), varying the architecture while holding the optimization procedure fixed (between all tested ImageNet-categorization trained DCNN architectures), and holding the architecture fixed while varying the optimization procedure (between ImageNet-categorization trained Inception-v3 and synthetic-categorization fine-tuned Inception-v3 models). This analysis addresses not only the reliability of the failure of DCNN_{IC} models to predict human behavior (deviations from humans), but also the relative importance of the characteristics defining similarities within the model sub-family (namely, the architecture and the optimization procedure). We first performed this analysis for residual signatures over the 240 primary test images, and subsequently zoomed in on subsets of images that humans found to be particularly difficult. This image selection was made relative to the distribution of image-level performance of held-out human subjects (B.II metric from five subjects); difficult images were defined as ones with performance below the 25th percentile of this distribution.

To examine whether the difference between each model and humans can be explained by simple human-interpretable stimulus attributes, we regressed each DCNN_{IC} model's residual signature on image attributes (object size, eccentricity, pose, and contrast). Briefly, we constructed a design matrix from the image attributes (using individual attributes, or all attributes), and used multiple linear least squares regression to predict the image-level residual signature. The multiple linear regression was tested using two-fold cross-validation over trials. The relative importance of each attribute (or groups of attributes) was quantified using the proportion of explainable variance (i.e. variance remaining after accounting for noise variance) explained from the residual signature.

Primate behavior zone

In this work, we are primarily concerned with the behavior of an “archetypal human”, rather than the behavior of any given individual human subject. We operationally defined this

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

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

To make a global statistical inference about whether models sampled from the DCNN_{IC} sub-family meet or fall short of this criterion for success, we attempted to reject the hypothesis that, for a given behavioral metric, the *human-consistency* of DCNN_{IC} models is within the primate zone. To test this hypothesis, we estimated the empirical probability that the distribution of *human-consistency* values, estimated over different model instances within this family, could produce *human-consistency* values within the primate zone. Specifically, we estimated a p-value for each behavioral metric using the following procedure: We first estimated an empirical distribution of Fisher-transformed *human-consistency* values for this model family (i.e. over all tested DCNN_{IC} models and over all trial-resampling of each DCNN_{IC} model). From this

empirical distribution, we fit a Gaussian kernel density function, optimizing the bandwidth parameter to minimize the mean squared error to the empirical distribution. This kernel density function was evaluated to compute a p-value, by computing the cumulative probability of observing a *human-consistency* value greater than or equal to the criterion of success (i.e. the Fisher transformed $\tilde{\rho}_{M\infty}$ value). This p-value indicates the probability that *human-consistency* values sampled from the observed distribution would fall into the primate zone, with smaller p-values indicating stronger evidence against the hypothesis that the *human-consistency* of DCNN models is within the primate zone.

RESULTS

In the present work, we systematically compared the basic level core object recognition behavior of primates and state-of-the-art artificial neural network models using a series of behavioral metrics ranging from low to high resolution within a two-alternative forced choice match-to-sample paradigm. The behavior of each visual system, whether biological or artificial, was tested on the same 2400 images (24 objects, 100 images/object) in the same 276 interleaved binary object recognition tasks. Each system's behavior was characterized at multiple resolutions (see *Behavioral metrics and signatures* in Methods) and directly compared to the corresponding behavioral metric applied on the archetypal human (defined as the average behavior of a large pool of human subjects tested; see Methods). The overarching logic of this study was that, if two visual systems are equivalent, they should produce statistically indistinguishable behavioral signatures with respect to these metrics. Specifically, our goal was to compare the behavioral signatures of visual system models with the corresponding behavioral signatures of primates.

Object-level behavioral comparison

We first examined the pattern of one-versus-all object-level behavior (termed “B.O1 metric”) computed across all images and possible distractors. Since we tested 24 objects here, the B.O1 signature was 24 dimensional. Figure 2A shows the B.O1 signatures for the pooled human (pooling $n=1472$ human subjects), pooled monkey (pooling $n=5$ monkey subjects), and several DCNN_{IC} models as 24-dimensional vectors using a color scale. Each element of the vector corresponds to the system's discriminability of one object against all others that were tested (i.e.

all other 23 objects). The color scales span each signature's full performance range, and warm colors indicate lower discriminability. For example, red indicates that the tested visual system found the object corresponding to that element of the vector to be very challenging to discriminate from other objects (on average over all 23 discrimination tests, and on average over all images). Figure 2B directly compares the B.O1 signatures computed from the behavioral output of two visual system models—a pixel model (top panel) and a DCNN_{IC} model (Inception-v3, bottom panel)—against that of the human B.O1 signature. We observe a tighter correspondence to the human behavioral signature for the DCNN_{IC} model visual system than for the baseline pixel model visual system. We quantified that similarity using a noise-adjusted correlation between each pair of B.O1 signatures (termed *human-consistency*, following (Johnson et al., 2002)); the noise adjustment means that a visual system that is identical to the human pool will have an expected *human-consistency* score of 1.0, even if it has irreducible trial-by-trial stochasticity; see Methods). Figure 2C shows the B.O1 *human-consistency* for each of the tested model visual systems. We additionally tested the behavior of a held-out pool of five human subjects (black dot) and a pool of five macaque monkey subjects (gray dot), and we observed that both yielded B.O1 signatures that were highly human-consistent (*human-consistency* $\tilde{p} = 0.90, 0.97$ for monkey pool and held-out human pool, respectively). We defined a range of *human-consistency* values, termed the “primate zone” (shaded gray area), delimited by extrapolated *human-consistency* estimates of large pools of macaques (see Methods, Figure 4). We found that the baseline pixel visual system model and the low-level V1 visual system model were not within this zone ($\tilde{p} = 0.40, 0.67$ for pixels and V1 models, respectively), while all tested DCNN_{IC} visual system models were either within or very close to this zone. Indeed, we could not reject the hypothesis that DCNN_{IC} models are primate-like ($p = 0.54$, exact test, see Methods).

Next, we compared the behavior of the visual systems at a slightly higher level of resolution. Specifically, instead of pooling over all discrimination tasks for each object, we computed the mean discriminability of each of the 276 pairwise discrimination tasks (still pooling over images within each of those tasks). This yielded a symmetric matrix that is referred to here as the B.O2 signature. Figure 2D shows the B.O2 signatures of the pooled human, pooled monkey, and several DCNN_{IC} visual system models as 24x24 symmetric matrices. Each bin (i, j) corresponds to the system's discriminability of objects i and j , where warmer colors indicate

609 lower performance; color scales are not shown but span each signature's full range. We observed
 610 strong qualitative similarities between the pairwise object confusion patterns of all of the high
 611 level visual systems (e.g. camel and dog are often confused with each other by all three systems).
 612 This similarity is quantified in Figure 2E, which shows the *human-consistency* of all examined
 613 visual system models with respect to this metric. Similar to the B.O1 metric, we observed that
 614 both a pool of macaque monkeys and a held-out pool of humans are highly *human-consistent*
 615 with respect to this metric ($\tilde{p} = 0.77, 0.94$ for monkeys, humans respectively). Also similar to the
 616 B.O1 metric, we found that all DCNN_{IC} visual system models are highly *human-consistent* ($\tilde{p} >$
 617 0.8) while the baseline pixel visual system model and the low-level V1 visual system model were
 618 not ($\tilde{p} = 0.41, 0.57$ for pixels, V1 models respectively). Indeed, all DCNN_{IC} visual system
 619 models are within the defined "primate zone" of *human-consistency*, and we could not falsify the
 620 hypothesis that DCNN_{IC} models are primate-like ($p = 0.99$, exact test).

621

622 Taken together, humans, monkeys, and current DCNN_{IC} models all share similar patterns
 623 of object-level behavioral performances (B.O1 and B.O2 signatures) that are not shared with
 624 lower-level visual representations (pixels and V1). However, object-level performance patterns
 625 do not capture the fact that some images of an object are more challenging than other images of
 626 the same object because of interactions of the variation in the object's pose and position with the
 627 object's class. To overcome this limitation, we next examined the patterns of behavior at the
 628 resolution of individual images on a subsampled set of images where we specifically obtained a
 629 large number of behavioral trials to accurately estimate behavioral performance on each image.
 630 Note that, from the point of view of the subjects, the behavioral tasks are identical to those
 631 already described. We simply aimed to measure and compare their patterns of performance at
 632 much higher resolution.

633

634 *Image-level behavioral comparison*

635 To isolate purely image-level behavioral variance, i.e. variance that is not predicted by
 636 the object and thus already captured by the B.O1 signature, we computed the normalized image-
 637 level signature. This normalization procedure is schematically illustrated in Figure 3A which
 638 shows that the one-versus-all image-level signature (240-dimensional, 10 images/object) is used
 639 to obtain the normalized one-versus-all image-level signature (termed B.I1n, see *Behavioral*

640 *metrics and signatures*). Figure 3B shows the B.IIn signatures for the pooled human, pooled
 641 monkey, and several DCNN_{IC} models as 240 dimensional vectors. Each bin's color corresponds
 642 to the discriminability of a single image against all distractor options (after subtraction of object-
 643 level discriminability, see Figure 3A), where warmer colors indicate lower values; color scales
 644 are not shown but span each signature's full range. Figure 3C shows the *human-consistency* with
 645 respect to the B.IIn signature for all tested models. Unlike with object-level behavioral metrics,
 646 we now observe a divergence between DCNN_{IC} models and primates. Both the monkey pool and
 647 the held-out human pool remain highly *human-consistent* ($\tilde{p} = 0.77, 0.96$ for monkeys, humans
 648 respectively), but all DCNN_{IC} models were significantly less *human-consistent* (Inception-
 649 v3: $\tilde{p} = 0.62$) and well outside of the defined "primate zone" of B.IIn *human-consistency*.
 650 Indeed, the hypothesis that the *human-consistency* of DCNN_{IC} models is within the primate zone
 651 is strongly rejected ($p = 6.16\text{e-}8$, exact test, see Methods).

652
 653 We can zoom in further by examining not only the overall performance for a given image
 654 but also the object confusions for each image, i.e. the additional behavioral variation that is due
 655 not only to the test image but to the interaction of that test image with the alternative (incorrect)
 656 object choice that is provided after the test image (see Fig. 1B). This is the highest level of
 657 behavioral accuracy resolution that our task design allows. In raw form, it corresponds to one-
 658 versus-other image-level confusion matrix, where the size of that matrix is the total number of
 659 images by the total number of objects (here, 240x24). Each bin (i,j) corresponds to the behavioral
 660 discriminability of a single image i against distractor object j . Again, we isolate variance that is
 661 not predicted by object-level performance by subtracting the average performance on this binary
 662 task (mean over all images) to convert the raw matrix B.I2 above into the normalized matrix,
 663 referred to as B.I2n. Figure 3D shows the B.I2n signatures as 240x24 matrices for the pooled
 664 human, pooled monkey and top DCNN_{IC} visual system models. Color scales are not shown but
 665 span each signature's full range; warmer colors correspond to images with lower performance in
 666 a given binary task, relative to all images of that object in the same task. Figure 3E shows the
 667 *human-consistency* with respect to the B.I2n metric for all tested visual system models.
 668 Extending our observations using B.IIn, we observe a similar divergence between primates and
 669 DCNN_{IC} visual system models on the matrix pattern of image-by-distractor difficulties (B.I2n).
 670 Specifically, both the monkey pool and held-out human pool remain highly *human-consistent*

671 ($\tilde{p} = 0.75, 0.77$ for monkeys, humans respectively), while all tested DCNN_{IC} models are
 672 significantly less *human-consistent* (Inception-v3: $\tilde{p} = 0.53$) falling well outside of the defined
 673 “primate zone” of B.I2n *human-consistency* values. Once again, the hypothesis that the *human-*
 674 *consistency* of DCNN_{IC} models is within the primate zone is strongly rejected ($p = 3.17\text{e-}18$,
 675 exact test, see Methods).

676

677 *Natural subject-to-subject variation*

678 For each behavioral metric (B.O1, BO2, B.I1n, BI2n), we defined a “primate zone” as the
 679 range of consistency values delimited by *human-consistency* estimates $\tilde{p}_{M\infty}$ and $\tilde{p}_{H\infty}$ as lower
 680 and upper bounds respectively. $\tilde{p}_{M\infty}$ corresponds to the extrapolated estimate of the *human-*
 681 *consistency* of a large (i.e. infinitely many subjects) pool of rhesus macaque monkeys. Thus, the
 682 fact that a particular tested visual system model falls outside of the primate zone can be
 683 interpreted as a failure of that visual system model to accurately predict the behavior of the
 684 archetypal human at least as well as the archetypal monkey.

685

686 However, from the above analyses, it is not yet clear whether a visual system model that
 687 fails to predict the archetypal human might nonetheless accurately correspond to one or more
 688 individual human subjects found within the natural variation of the human population. Given the
 689 difficulty of measuring individual subject behavior at the resolution of single images for large
 690 numbers of human and monkey subjects, we could not yet directly test this hypothesis. Instead,
 691 we examined it indirectly by asking whether an archetypal model—that is a pool that includes an
 692 increasing number of model “subjects”—would approach the human pool. We simulated model
 693 inter-subject variability by retraining a fixed DCNN architecture with a fixed training image set
 694 with random variation in the initial conditions and order of training images. This procedure
 695 results in models that can still perform the task but with slightly different learned weight values.
 696 We note that this procedure is only one possible choice of generating inter-subject variability
 697 within each visual system model type, a choice that is an important open research direction that
 698 we do not address here. From this procedure, we constructed multiple trained model instances
 699 (“subjects”) for a fixed DCNN architecture, and asked whether an increasingly large pool of
 700 model “subjects” better captures the behavior of the human pool, at least as well as a monkey

pool. This post-hoc analysis was conducted for the most *human-consistent* DCNN architecture (Inception-v3).

Figure 4A shows, for each of the four behavioral metrics, the measured *human-consistency* of subject pools of varying size (number of subjects n) of rhesus macaque monkeys (black) and ImageNet-trained Inception-v3 models (blue). The *human-consistency* increases with growing number of subjects for both visual systems across all behavioral metrics. To estimate the expected *human-consistency* for a pool of infinitely many monkey or model subjects, we fit an exponential function mapping n to the mean *human-consistency* values and obtained a parameter estimate for the asymptotic value (see Methods). We note that estimated asymptotic values are not significantly beyond the range of the measured data—the *human-consistency* of a pool of five monkey subjects reaches within 97% of the *human-consistency* of an estimated infinite pool of monkeys for all metrics—giving credence to the extrapolated *human-consistency* values. This analysis suggests that under this model of inter-subject variability, a pool of Inception-v3 subjects accurately capture archetypal human behavior at the resolution of objects (B.O1, B.O2) by our primate zone criterion (see Figure 4A, first two panels). In contrast, even a large pool of Inception-v3 subjects still fails at its final asymptote to accurately capture human behavior at the image-level (B.I1n, B.I2n) (Figure 4A, last two panels).

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

Next, we wondered if some relatively simple changes to the DCNN_{IC} visual system models tested here could bring them into better correspondence with the primate visual system behavior (with respect to B.I1n and B.I2n metrics). Specifically, we considered and tested the following modifications to the most *human-consistent* DCNN_{IC} model visual system (Inception-v3): we (1) changed the input to the model to be more primate-like in its retinal sampling (Inception-v3 + retina-like), (2) changed the transformation (aka “decoder”) from the internal model feature representation into the behavioral output by augmenting the number of decoder training images or changing the decoder type (Inception-v3 + SVM, Inception-v3 + classifier_train), and (3) modified all of the internal filter weights of the model (aka “fine tuning”) by augmenting its ImageNet training with additional images drawn from the same distribution as our test images (Inception-v3 + synthetic-fine-tune) and (4) modified all of the

internal filter weights of the model by training exclusively on images drawn from the same distribution as our test images (Inception-v3 + synthetic-train). All model modifications resulted in relatively high-performing models; Figure 5 (gray bars) shows the mean overall performance over object recognition tasks, evaluated with a fixed decoder type (logistic classifier, 20 training images per class). However, we found that none of these modifications led to a significant improvement in its *human-consistency* on the behavioral metrics (Figure 4B). In particular, training exclusively on synthetic images led to a noted decrease in human consistency, across all metrics. Figure 5 shows the relationship between model performances (B.O2, averaged over 276 tasks) and human-consistency, for both object-level and image-level metrics. We observe a strong correlation between performance and image-level human-consistency across different model architectures trained on ImageNet (black points, $r=0.97$, $p<10^{-4}$), but no such correlation across our modified models within a fixed architecture (grey points, $r=0.35$, $p=0.13$). Thus, the failure of current DCNN_{IC} models to accurately capture the image-level signatures of primates cannot be rescued by simple modifications on a fixed architecture.

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

Taken together, Figures 2, 3, 4 and 5 suggest that current DCNN_{IC} visual system models fail to accurately capture the image-level behavioral signatures of humans and monkeys. To further examine this failure in the hopes of providing clues for model improvement, we examined the image-level residual signatures of all the visual system models, relative to the pooled human. For each model, we computed its residual signature as the difference (positive or negative) of a linear least squares regression of the model signature on the corresponding human signature. For this analysis, we focused on the B.I1n metric as it showed a clear divergence of DCNN_{IC} models and primates, and the behavioral residual can be interpreted based only on the test images (whereas B.I2n depends on the interaction between test images and distractor choice).

We first asked to what extent the residual signatures are shared between different visual system models. Figure 6A shows the similarity between the residual signatures of all pairs of models; the color of bin (i,j) indicates the proportion of explainable variance that is shared between the residual signatures of visual systems i and j . For ease of interpretation, we ordered

visual system models based on their architecture and optimization procedure and partitioned this matrix into four distinct regions. Each region compares the residuals of a “source” model group with fixed architecture and optimization procedure (five Inception-v3 models optimized for categorization on ImageNet, varying only in initial conditions and training image order) to a “target” model group. The target groups of models for each of the four regions are: 1) the pooled monkey, 2) other DCNN_{IC} models from the source group, 3) DCNN_{IC} models that differ in architecture but share the optimization procedure of the source group models and 4) DCNN_{IC} models that differ slightly using an augmented optimization procedure but share the architecture of the source group models. Figure 6B shows the mean (\pm SD) variance shared in the residuals averaged within these four regions for all images (black dots), as well as for images that humans found to be particularly difficult (gray dots, selected based on held-out human data, see Methods). First, consistent with the results shown in Figure 3, we note that the residual signatures of this particular DCNN_{IC} model are not well shared with the pooled monkey ($r^2=0.39$ in region 1), and this phenomenon is more pronounced for the images that humans found most difficult ($r^2=0.17$ in region 1). However, this relatively low correlation between model and primate residuals is not indicative of spurious model residuals, as the model residual signatures were highly reliable between different instances of this fixed DCNN_{IC} model, across random training initializations (region 2: $r^2=0.79$, 0.77 for all and most difficult images, respectively). Interestingly, residual signatures were still largely shared with other DCNN_{IC} models with vastly different architectures (region 3: $r^2=0.70$, 0.65 for all and most difficult images, respectively). However, residual signatures were more strongly altered when the visual training diet of the same architecture was altered (region 4: $r^2=0.57$, 0.46 for all and most difficult images respectively, cf. region 3). Taken together, these results indicate that the images where DCNN_{IC} visual system models diverged from humans (and monkeys) were not spurious but were rather highly reliable across different model architectures, demonstrating that current DCNN_{IC} models systematically and similarly diverge from primates.

Figure 7A shows example images, sorted by B.IIn squared residuals, corresponding to images that models and primates largely agreed on (left) and diverged on (right), with respect to the B.IIn metric. We observed no qualitative differences between these images. To look for clues for model improvement, we asked what, if any, characteristics of images might account for

794 this divergence of models and primates. We regressed the residual signatures of DCNN_{IC} models
 795 on four different image attributes (corresponding to the size, eccentricity, pose, and contrast of
 796 the object in each image). We used multiple linear regressions to predict the model residual
 797 signatures from all of these image attributes, and also considered each attribute individually
 798 using simple linear regressions. Figure 7B shows example images (sampled from the full set of
 799 2400 images) with increasing attribute value for each of these four image attributes. While the
 800 DCNN_{IC} models were not directly optimized to display primate-like performance dependence on
 801 such attributes, we observed that the Inception-v3 visual system model nonetheless exhibited
 802 qualitatively similar performance dependencies as primates (see Figure 7C). For example,
 803 humans (black), monkeys (gray) and the Inception-v3 model (blue) all performed better, on
 804 average, for images in which the object is in the center of gaze (low eccentricity) and large in
 805 size. Furthermore, all three systems performed better, on average, for images when the pose of
 806 the object was closer to the canonical pose (see Figure 7C); this sensitivity to object pose
 807 manifested itself as a non-linear dependence due to the fact that all tested objects exhibited
 808 symmetry in at least one axis. The similarity of the patterns in Figure 7C between primates and
 809 the DCNN_{IC} visual system models is not perfect but is striking, particularly in light of the fact
 810 that these models were not optimized to produce these patterns. However, this similarity is
 811 analogous to the similarity in the B.O1 and B.O2 metrics in that it only holds on average over
 812 many images. Looking more closely at the image-by-image comparison, we again found that the
 813 DCNN_{IC} models failed to capture a large portion of the image-by-image variation (Figure 3). In
 814 particular, Figure 7D shows the proportion of variance explained by specific image attributes for
 815 the residual signatures of monkeys (black) and DCNN_{IC} models (blue). We found that, taken
 816 together, all four of these image attributes explained only ~10% of the variance in DCNN_{IC}
 817 residual signatures, and each individual attribute could explain at most a small amount of
 818 residual variance (<5% of the explainable variance). In sum, these analyses show that some
 819 behavioral effects that might provide intuitive clues to modify the DCNN_{IC} models are already in
 820 place in those models (e.g. a dependence on eccentricity). But the quantitative image-by-image
 821 analyses of the remaining unexplained variance (Figure 7D) argue that the DCNN_{IC} visual
 822 system models' failure to capture primate image-level signatures cannot be further accounted for
 823 by these simple image attributes and likely stem from other factors.

824

825 **DISCUSSION**

826

827 The current work was motivated by the broad scientific goal of discovering models that
828 quantitatively explain the neuronal mechanisms underlying primate invariant object recognition
829 behavior. To this end, previous work had shown that specific artificial neural network models
830 (ANNs), drawn from a large family of deep convolutional neural networks (DCNNs) and
831 optimized to achieve high levels of object categorization performance on large-scale image-sets,
832 capture a large fraction of the variance in primate visual recognition behaviors (Ghodrati et al.,
833 2014; Rajalingham et al., 2015; Jozwik et al., 2016; Kheradpisheh et al., 2016; Kubilius et al.,
834 2016; Peterson et al., 2016; Battleday et al., 2017; Wallis et al., 2017), and the internal hidden
835 neurons of those same models also predict a large fraction of the image-driven response variance
836 of brain activity at multiple stages of the primate ventral visual stream (Yamins et al., 2013;
837 Cadieu et al., 2014; Khaligh-Razavi and Kriegeskorte, 2014; Yamins et al., 2014; Güçlü and van
838 Gerven, 2015; Cichy et al., 2016; Hong et al., 2016; Seibert et al., 2016; Cadena et al., 2017;
839 Cichy et al., 2017; Eickenberg et al., 2017; Khaligh-Razavi et al., 2017; Seeliger et al., 2017;
840 Wen et al., 2017). For clarity, we here referred to this sub-family of models as DCNN_{IC} (to
841 denote ImageNet-Categorization training), so as to distinguish them from all possible models in
842 the DCNN family, and more broadly, from the super-family of all ANNs. In this work, we
843 directly compared leading DCNN_{IC} models to primates (humans and monkeys) with respect to
844 their behavioral signatures at both object and image level resolution in the domain of core object
845 recognition. In order to do so, we measured and characterized primate behavior at larger scale
846 and higher resolution than previously possible. We first replicate prior work (Rajalingham et al.,
847 2015) showing that, at the object level, DCNN_{IC} models produce statistically indistinguishable
848 behavior from primates, and we extend that work by showing that these models also match the
849 *average* primate sensitivities to object contrast, eccentricity, size, and pose, a noteworthy
850 similarity in light of the fact that these models were not optimized to produce these performance
851 patterns. This similarity, which we speculate may largely reflect the training history of these
852 models (i.e. optimization on photographer-framed images, which may be biased with respect to
853 these attributes), is a good direction for future work. However, our primary novel result is that,
854 examining behavior at the higher resolution of individual images, all leading DCNN_{IC} models
855 failed to replicate the image-level behavioral signatures of primates. An important related claim

856 is that rhesus monkeys are more consistent with the archetypal human than any of the tested
 857 DCNN_{IC} models (at the image-level).

858

859 We compared human, monkey and model behavior on a set of naturalistic synthetic
 860 images, sampled with systematic variation in viewpoint parameters and uncorrelated
 861 background. While these images may seem non-natural, synthetic images of this types are –
 862 unlike some photographic image sets – very powerful at separating humans from low-level
 863 computer vision systems (Pinto et al., 2008). Furthermore, given our definition of success (a
 864 candidate visual system model that accurately captures primate behavior *for all images*), the
 865 inference that models differ from primates can be supported by testing models on any image set,
 866 including the synthetic set we used. Importantly, this inference is with respect to a trained
 867 model’s behavior at “run time”, agnostic to its learning procedure. While we do not make any
 868 claims about learning, we note that neither humans nor monkeys were separately optimized to
 869 perform on our synthetic images (beyond the minimal training described in the Methods section).
 870 Interestingly, we observed that synthetic-image-optimized models (new ANN models
 871 constructed by fine-tuning or training an existing network architecture exclusively on large sets
 872 of synthetic image) were no more similar to primates than ANN models optimized only on
 873 ImageNet, suggesting that the tested ANN architectures have one or more fundamental flaws that
 874 cannot be readily overcome by manipulating the training environment. Taken together, these
 875 results support the general inference, agnostic to natural choices of image-set, that DCNN_{IC}
 876 models diverge from primates in their core object recognition behavior.

877

878 This inference is consistent with previous work showing that DCNN_{IC} models can
 879 diverge from human behavior on specifically chosen adversarial images (Szegedy et al., 2013). A
 880 strength of our work is that we did not optimize images to induce failure but instead randomly
 881 sampled a broadly-defined image generative parameter space highlighting a general, rather than
 882 adversarial-induced, divergence of DCNN_{IC} from primate core object recognition behavior.
 883 Furthermore, we showed that this failure could not be rescued by a number of simple
 884 modifications to the model class, providing qualitative insight into what does and does not
 885 explain the gap between primates and DCNN models. Taken together, these results suggest that

new ANN models are needed to more precisely capture the neural mechanisms underlying primate object vision.

With regards to new ANN models, we can attempt to make prospective inferences about future possible DCNN_{IC} models from the data presented here. Based on the observed distribution of image-level *human-consistency* values for the DCNN_{IC} models tested here, we infer that yet untested model instances sampled identically (i.e. from the DCNN_{IC} model sub-family) are very likely to have similarly inadequate image-level *human-consistency*. While we cannot rule out the possibility that at least one model instance within the DCNN_{IC} sub-family would fully match the image-level behavioral signatures, the probability of sampling such a model is vanishingly small ($p < 10^{-17}$ for B.I2n *human-consistency*, estimated using exact test using Gaussian kernel density estimation, see Methods, Results). An important caveat of this inference is that we may have a biased estimate of the *human-consistency* distribution of this model sub-family, as we did not exhaustively sample the sub-family. In particular, if the model sampling process is non-stationary over time (e.g. increases in computational power over time allows larger models to be successfully trained), the *human-consistency* of new (i.e. yet to be sampled) models may lie outside the currently estimated distribution. Thus, it is possible that the evolution of “next-generation” models within the DCNN_{IC} sub-family could meet our criteria for successful matching primate-like behavior.

Alternatively, it is possible—and we think likely—that future DCNN_{IC} models will also fail to capture primate-like image-level behavior, suggesting that either the architectural limitations (e.g. convolutional, feed-forward) and/or the optimization procedure (including the diet of visual images) that define this model sub-family are fundamentally limiting. Thus, ANN model sub-families utilizing different architectures (e.g. recurrent neural networks) and/or optimized for different behavioral goals (e.g. loss functions other than object classification performance, and/or images other than category-labeled ImageNet images) may be necessary to accurately capture primate behavior. To this end, we propose that testing even individual changes to the DCNN_{IC} models—each creating a new ANN model sub-family—may be the best way forward, because DCNN_{IC} models currently offer the best explanations (in a predictive sense) of both the behavioral and neural phenomena of core object recognition.

917

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

939

940

941 REFERENCES

- 942
943 Battleday RM, Peterson JC, Griffiths TL (2017) Modeling human categorization of natural
944 images using deep feature representations. arXiv preprint arXiv:171104855.
945 Cadena SA, Denfield GH, Walker EY, Gatys LA, Tolia AS, Bethge M, Ecker AS (2017) Deep
946 convolutional models improve predictions of macaque V1 responses to natural images.
947 bioRxiv:201764.
948 Cadieu CF, Hong H, Yamins DL, Pinto N, Ardila D, Solomon EA, Majaj NJ, DiCarlo JJ (2014)
949 Deep neural networks rival the representation of primate IT cortex for core visual object
950 recognition. PLoS computational biology 10:e1003963.
951 Cichy RM, Khosla A, Pantazis D, Oliva A (2017) Dynamics of scene representations in the
952 human brain revealed by magnetoencephalography and deep neural networks. Neurolmage
953 153:346-358.
954 Cichy RM, Khosla A, Pantazis D, Torralba A, Oliva A (2016) Comparison of deep neural
955 networks to spatio-temporal cortical dynamics of human visual object recognition reveals
956 hierarchical correspondence. Scientific reports 6:27755.
957 DiCarlo JJ, Johnson KO (1999) Velocity invariance of receptive field structure in somatosensory
958 cortical area 3b of the alert monkey. Journal of Neuroscience 19:401-419.
959 DiCarlo JJ, Cox DD (2007) Untangling invariant object recognition. Trends in cognitive sciences
960 11:333-341.
961 DiCarlo JJ, Zoccolan D, Rust NC (2012) How does the brain solve visual object recognition?
962 Neuron 73:415-434.
963 Dodge S, Karam L (2017) A Study and Comparison of Human and Deep Learning Recognition
964 Performance Under Visual Distortions. arXiv preprint arXiv:170502498.
965 Eickenberg M, Gramfort A, Varoquaux G, Thirion B (2017) Seeing it all: Convolutional network
966 layers map the function of the human visual system. Neurolmage 152:184-194.
967 Geirhos R, Janssen DH, Schütt HH, Rauber J, Bethge M, Wichmann FA (2017) Comparing
968 deep neural networks against humans: object recognition when the signal gets weaker. arXiv
969 preprint arXiv:170606969.
970 Ghodrati M, Farzmaadi A, Rajaei K, Ebrahimpour R, Khaligh-Razavi S-M (2014) Feedforward
971 object-vision models only tolerate small image variations compared to human. Frontiers in
972 computational neuroscience 8:74.
973 Goodfellow IJ, Shlens J, Szegedy C (2014) Explaining and harnessing adversarial examples.
974 arXiv preprint arXiv:14126572.
975 Güçlü U, van Gerven MA (2015) Deep neural networks reveal a gradient in the complexity of
976 neural representations across the ventral stream. Journal of Neuroscience 35:10005-10014.
977 He K, Zhang X, Ren S, Sun J (2016) Deep residual learning for image recognition. In:
978 Proceedings of the IEEE conference on computer vision and pattern recognition, pp 770-778.
979 Hong H, Yamins DL, Majaj NJ, DiCarlo JJ (2016) Explicit information for category-orthogonal
980 object properties increases along the ventral stream. Nature neuroscience 19:613-622.
981 Hosseini H, Xiao B, Jaiswal M, Poovendran R (2017) On the Limitation of Convolutional Neural
982 Networks in Recognizing Negative Images. human performance 4:6.
983 Huynh DQ (2009) Metrics for 3D rotations: Comparison and analysis. Journal of Mathematical
984 Imaging and Vision 35:155-164.
985 Johnson KO, Hsiao SS, Yoshioka T (2002) Neural coding and the basic law of psychophysics.
986 The Neuroscientist 8:111-121.
987 Jozwik KM, Kriegeskorte N, Mur M (2016) Visual features as stepping stones toward semantics:
988 Explaining object similarity in IT and perception with non-negative least squares.
989 Neuropsychologia 83:201-226.

- 990 Khaligh-Razavi S-M, Kriegeskorte N (2014) Deep supervised, but not unsupervised, models
 991 may explain IT cortical representation. *PLoS computational biology* 10:e1003915.
- 992 Khaligh-Razavi S-M, Henriksson L, Kay K, Kriegeskorte N (2017) Fixed versus mixed RSA:
 993 Explaining visual representations by fixed and mixed feature sets from shallow and deep
 994 computational models. *Journal of mathematical psychology* 76:184-197.
- 995 Kheradpisheh SR, Ghodrati M, Ganjtabesh M, Masquelier T (2016) Deep networks can
 996 resemble human feed-forward vision in invariant object recognition. *Scientific reports* 6:32672.
- 997 Kriegeskorte N (2015) Deep neural networks: a new framework for modeling biological vision
 998 and brain information processing. *Annual review of vision science* 1:417-446.
- 999 Krizhevsky A, Sutskever I, Hinton GE (2012) Imagenet classification with deep convolutional
 1000 neural networks. In: *Advances in neural information processing systems*, pp 1097-1105.
- 1001 Kubilius J, Bracci S, de Beeck HPO (2016) Deep neural networks as a computational model for
 1002 human shape sensitivity. *PLoS computational biology* 12:e1004896.
- 1003 LeCun Y, Bengio Y, Hinton G (2015) Deep learning. *Nature* 521:436-444.
- 1004 Majaj NJ, Hong H, Solomon EA, DiCarlo JJ (2015) Simple Learned Weighted Sums of Inferior
 1005 Temporal Neuronal Firing Rates Accurately Predict Human Core Object Recognition
 1006 Performance. *The Journal of Neuroscience* 35:13402-13418.
- 1007 Nguyen A, Yosinski J, Clune J (2015) Deep neural networks are easily fooled: High confidence
 1008 predictions for unrecognizable images. In: *Proceedings of the IEEE Conference on Computer
 1009 Vision and Pattern Recognition*, pp 427-436.
- 1010 Peterson JC, Abbott JT, Griffiths TL (2016) Adapting deep network features to capture
 1011 psychological representations. *arXiv preprint arXiv:160802164*.
- 1012 Pinto N, Cox DD, DiCarlo JJ (2008) Why is real-world visual object recognition hard? *PLoS
 1013 computational biology* 4:e27.
- 1014 Rajalingham R, Schmidt K, DiCarlo JJ (2015) Comparison of Object Recognition Behavior in
 1015 Human and Monkey. *The Journal of Neuroscience* 35:12127-12136.
- 1016 RichardWebster B, Anthony SE, Scheirer WJ (2016) *PsyPhy: A Psychophysics Driven
 1017 Evaluation Framework for Visual Recognition*. *arXiv preprint arXiv:161106448*.
- 1018 Rolls ET (2000) Functions of the primate temporal lobe cortical visual areas in invariant visual
 1019 object and face recognition. *Neuron* 27:205-218.
- 1020 Rosch E, Mervis CB, Gray WD, Johnson DM, Boyes-Braem P (1976) Basic objects in natural
 1021 categories. *Cognitive psychology* 8:382-439.
- 1022 Seeliger K, Fritsche M, Güçlü U, Schoenmakers S, Schoffelen J-M, Bosch S, van Gerven M
 1023 (2017) Convolutional neural network-based encoding and decoding of visual object recognition
 1024 in space and time. *NeuroImage*.
- 1025 Seibert D, Yamins DL, Ardila D, Hong H, DiCarlo JJ, Gardner JL (2016) A performance-
 1026 optimized model of neural responses across the ventral visual stream. *bioRxiv:036475*.
- 1027 Simonyan K, Zisserman A (2014) Very deep convolutional networks for large-scale image
 1028 recognition. *arXiv preprint arXiv:14091556*.
- 1029 Szegedy C, Zaremba W, Sutskever I, Bruna J, Erhan D, Goodfellow I, Fergus R (2013)
 1030 Intriguing properties of neural networks. *arXiv preprint arXiv:13126199*.
- 1031 Tanaka K (1996) Inferotemporal cortex and object vision. *Annual review of neuroscience*
 1032 19:109-139.
- 1033 Ullman S, Humphreys GW (1996) *High-level vision: Object recognition and visual cognition*: MIT
 1034 press Cambridge, MA.
- 1035 Wallis TS, Funke CM, Ecker AS, Gatys LA, Wichmann FA, Bethge M (2017) A parametric
 1036 texture model based on deep convolutional features closely matches texture appearance for
 1037 humans. *Journal of vision* 17:5-5.
- 1038 Wen H, Shi J, Zhang Y, Lu K-H, Cao J, Liu Z (2017) Neural encoding and decoding with deep
 1039 learning for dynamic natural vision. *Cerebral Cortex*:1-25.

1040 Yamins DL, DiCarlo JJ (2016) Using goal-driven deep learning models to understand sensory
1041 cortex. *Nature neuroscience* 19:356-365.
1042 Yamins DL, Hong H, Cadieu C, DiCarlo JJ (2013) Hierarchical modular optimization of
1043 convolutional networks achieves representations similar to macaque IT and human ventral
1044 stream. In: *Advances in neural information processing systems*, pp 3093-3101.
1045 Yamins DL, Hong H, Cadieu CF, Solomon EA, Seibert D, DiCarlo JJ (2014) Performance-
1046 optimized hierarchical models predict neural responses in higher visual cortex. *Proceedings of*
1047 *the National Academy of Sciences*:201403112.
1048 Zeiler MD, Fergus R (2014) Visualizing and understanding convolutional networks. In: *Computer*
1049 *Vision—ECCV 2014*, pp 818-833: Springer.

1050
1051

1052 **TABLES**1053 **Table 1**

Behavioral Metric	Hit Rate	False Alarm Rate
One-versus all object-level performance (B.O1) ($N_{\text{objects}} \times 1$) $O_1(i) = Z(HR(i)) - Z(FAR(i)),$ $i = 1, 2, \dots, N_{\text{objects}}$	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 ($N_{\text{objects}} \times N_{\text{objects}}$) $O_2(i, j) = Z(HR(i, j)) - Z(FAR(i, j)),$ $i = 1, 2, \dots, N_{\text{objects}}$ $j = 1, 2, \dots, N_{\text{objects}}$	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 ($N_{\text{images}} \times 1$) $I_1(ii) = Z(HR(ii)) - Z(FAR(ii)),$ $ii = 1, 2, \dots, N_{\text{images}}$	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 ($N_{\text{images}} \times N_{\text{objects}}$) $I_2(ii, j) = Z(HR(ii, j)) - Z(FAR(ii, j)),$ $ii = 1, 2, \dots, N_{\text{images}}$ $j = 1, 2, \dots, N_{\text{objects}}$	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

1054

1055 **Table 1: Definition of behavioral performance metrics.** The first column provides the name,
1056 abbreviation, dimensions, and equations for each of the raw performance metrics. The next two
1057 columns provide the definitions for computing the hit rate (HR) and false alarm rate (FAR)
1058 respectively.

1059

1060 **FIGURE LEGENDS**

1061

1062 **Figure 1. Images and behavioral task. (A)** Two (out of 100) example images for each of the 24
 1063 basic-level objects. To enforce true invariant object recognition behavior, we generated
 1064 naturalistic synthetic images, each with one foreground object, by rendering a 3D model of each
 1065 object with randomly chosen viewing parameters and placing that foreground object view onto a
 1066 randomly chosen, natural image background. **(B)** Time course of example behavioral trial (zebra
 1067 versus dog) for human psychophysics. Each trial initiated with a central fixation point for 500
 1068 ms, followed by 100 ms presentation of a square test image (spanning 6-8° of visual angle).
 1069 After extinction of the test image, two choice images were shown to the left and right. Human
 1070 participants were allowed to freely view the response images for up to 1000 ms and responded
 1071 by clicking on one of the choice images; no feedback was given. To neutralize top-down feature
 1072 attention, all 276 binary object discrimination tasks were randomly interleaved on a trial-by-trial
 1073 basis. The monkey task paradigm was nearly identical to the human paradigm, with the
 1074 exception that trials were initiated by touching a fixation circle horizontally centered on the
 1075 bottom third of the screen, and successful trials were rewarded with juice while incorrect choices
 1076 resulted in timeouts of 1–2.5s. **(C)** Large-scale and high-throughput psychophysics in humans
 1077 (top left), monkeys (top right), and models (bottom). Human behavior was measured using the
 1078 online Amazon MTurk platform, which enabled the rapid collection ~1 million behavioral trials
 1079 from 1472 human subjects. Monkey behavior was measured using a novel custom home-cage
 1080 behavioral system (MonkeyTurk), which leveraged a web-based behavioral task running on a
 1081 tablet to test many monkey subjects simultaneously in their home environment. Deep
 1082 convolutional neural network models were tested on the same images and tasks as those
 1083 presented to humans and monkeys by extracting features from the penultimate layer of each
 1084 visual system model and training back-end multi-class logistic regression classifiers. All
 1085 behavioral predictions of each visual system model were for images that were not seen in any
 1086 phase of model training.

1087 **Figure 2. Object-level comparison to human behavior. (A)** One-versus-all object-level (B.O1)
 1088 signatures for the pooled human (n=1472 human subjects), pooled monkey (n=5 monkey
 1089 subjects), and several DCNN_{IC} models. Each B.O1 signature is shown as a 24-dimensional
 1090 vector using a color scale; each colored bin corresponds to the system's discriminability of one

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

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

1120 **Figure 4. Effect of subject pool size and DCNN model modifications on consistency with**
 1121 **human behavior. (A)** Accounting for natural subject-to-subject variability. For each of the four
 1122 behavioral metrics, the *human-consistency* distributions of monkey (blue markers) and model
 1123 (black markers) pools are shown as a function of the number of subjects in the pool (mean \pm SD,
 1124 over subjects). The human consistency increases with growing number of subjects for all visual
 1125 systems across all behavioral metrics. The dashed lines correspond to fitted exponential
 1126 functions, and the parameter estimate (mean \pm SE) of the asymptotic value, corresponding to the
 1127 estimated *human-consistency* of a pool of infinitely many subjects, is shown at the right most
 1128 point on each abscissa. **(B)** Model modifications that aim to rescue the DCNN_{IC} models. We
 1129 tested several simple modifications (see Methods) to the most *human-consistent* DCNN_{IC} visual
 1130 system model (Inception-v3). Each panel shows the resulting *human-consistency* per modified
 1131 model (mean \pm SD over different model instances, varying in random filter initializations) for
 1132 each of the four behavioral metrics.

1133
 1134 **Figure 5. Model performance. (A)** Model performance on synthetic images (average B.O2
 1135 across 276 tasks) for each of the tested models, fixing the number of training images and
 1136 classifier. Black bars correspond to different model architectures, with fixed optimization,
 1137 whereas grey bars correspond to different modifications of a fixed model architecture (Inception-
 1138 v3). **(B)** Correlation between model performance and human-consistency, with respect to object-
 1139 level (B.O2) and image-level (B.I2n) behavioral metrics. Each point corresponds to a single
 1140 instance of a trained DCNN model.

1141
 1142 **Figure 6. Analysis of unexplained human behavioral variance. (A)** Residual similarity
 1143 between all pairs of human visual system models. The color of bin (i,j) indicates the proportion
 1144 of explainable variance that is shared between the residual signatures of visual systems i and j .
 1145 For ease of interpretation, we ordered visual system models based on their architecture and
 1146 optimization procedure and partitioned this matrix into four distinct regions. **(B)** Summary of
 1147 residual similarity. For each of the four regions in Figure 5A, the similarity to the residuals of
 1148 Inception-v3 (region 2 in (A)) is shown (mean \pm SD, within each region) for all images (black
 1149 dots), and for images that humans found to be particularly difficult (gray dots, selected based on
 1150 held-out human data).

1151

1152 **Figure 7. Dependence of primate and DCNN model behavior on image attributes. (A)**
1153 Example images that models and primates agree on (left) and diverge on (right) with respect to
1154 B.IIn residuals. **(B)** Example images with increasing attribute value, for each of the four pre-
1155 defined image attributes (see Methods). **(C)** Dependence of performance (B.IIn) as a function of
1156 four image attributes, for humans, monkeys and a DCNN_{IC} model (Inception-v3). **(D)** Proportion
1157 of explainable variance of the residual signatures of monkeys (black) and DCNN_{IC} models (blue)
1158 that is accounted for by each of the pre-defined image attributes. Error-bars correspond to SD
1159 over trial re-sampling for monkeys, and over different models for DCNN_{IC} models.

