1              **Do deep neural networks see the way we do?**

2

3            Georgin Jacob[1,2], R. T. Pramod[1,2], Harish Katti[1] and S. P. Arun[1]*

4     [1]Centre for Neuroscience & [2]Department of Electrical Communication Engineering

5                Indian Institute of Science, Bangalore 560012

6                    *Correspondence to sparun@iisc.ac.in

7

8     Abbreviated Title           :   Perception in deep neural networks

9     Number of Figures          : 6 main and 10 supplementary figures

10

11 **ABSTRACT**

12    Deep neural networks have revolutionized computer vision, and their object

13 representations match coarsely with the brain. As a result, it is widely believed that

14 any fine scale differences between deep networks and brains can be fixed with

15 increased training data or minor changes in architecture. But what if there are

16 qualitative differences between brains and deep networks? Do deep networks even

17 see the way we do? To answer this question, we chose a deep neural network

18 optimized for object recognition and asked whether it exhibits well-known perceptual

19 and neural phenomena despite not being explicitly trained to do so. To our surprise,

20 many phenomena were present in the network, including the Thatcher effect, mirror

21 confusion, Weber's law, relative size, multiple object normalization and sparse coding

22 along multiple dimensions. However, some perceptual phenomena were notably

23 absent, including processing of 3D shape, patterns on surfaces, occlusion, natural

24 parts and a global advantage. Our results elucidate the computational challenges of

25 vision by showing that learning to recognize objects suffices to produce some

26 perceptual phenomena but not others and reveal the perceptual properties that could

27 be incorporated into deep networks to improve their performance.

28

29          **INTRODUCTION**

30                         *How do I know this is true?*
31                         *I look inside myself and see.*
32                                    *Tao Te Ching (Mitchell, 1988)*
33

34          Convolutional or deep neural networks have revolutionized computer vision

35   with their human-like accuracy on vision tasks, and their object representations match

36   coarsely with the brain (see Serre, 2019; Sinz et al., 2019 for detailed reviews). Yet,

37   at a finer scale, they are still outperformed by humans (Katti et al., 2017; Katti and

38   Arun, 2019) and show systematic deviations from human perception (Pramod and

39   Arun, 2016a; Geirhos et al., 2018b; Rajalingham et al., 2018; Dodge and Karam,

40   2019). Even these differences are largely quantitative in that there are no explicit or

41   emergent properties that are present in humans but absent in deep networks. This has

42   given rise to the prevailing belief that any remaining differences between brains and

43   deep networks can be fixed by training on larger datasets, incorporating more

44   constraints (Sinz et al., 2019) or by making relatively minor modifications to network

45   architecture such as by including recurrent feedback (Kar et al., 2019a; Kietzmann et

46   al., 2019).

47          Despite these insights, we do not yet know whether there are qualitative

48   differences between how brains and deep networks see. This is an important question

49   because resolving qualitative differences might require non-trivial changes in network

50   training or architecture. One approach could be to train deep networks on multiple

51   visual tasks and compare them with humans, but the answer would be insightful only

52   if networks fail to learn certain tasks (Fleuret et al., 2011a). Alternatively, we could

53   compare qualitative or emergent properties of our perception with that of deep

54   networks, provided these properties can indeed be checked in any deep network

55   without explicit training for these properties.

56        Fortunately, many classic findings from visual psychology and neuroscience

57    report emergent phenomena and properties that can be directly tested on deep

58    networks. Consider for instance, the classic Thatcher effect (Figure 1A), in which a

59    face with rotated parts looks grotesque in an upright orientation but looks entirely

60    normal when inverted (Thompson, 1980). This effect can be recast as a statement

61    about the underlying face representation: in perceptual space, the distance between

62    the normal and Thatcherized face is presumably larger when they are upright than

63    when they are inverted (Figure 1B). This has been confirmed using dissimilarity ratings

64    in humans (Bartlett and Searcy, 1993). These distances can be compared for any

65    representation, including for a deep network (Figure 1C). Since deep networks are

66    organized layer-wise with increasing complexity across layers, this would also reveal

67    the layers at which the deep network begins to experience or "see" a Thatcher effect

68    (Figure 1D).

69        Knowing whether a deep network exhibits the Thatcher effect can be insightful

70    for a variety of reasons. First, it would confirm that the deep network indeed does see

71    faces the way we do. Second, this question can be asked of any deep network without

72    explicit training to produce a Thatcher effect. For instance, testing this question on

73    face and object detection networks would reveal whether object or face-specific

74    training is sufficient for the emergence of the Thatcher effect. Finally, this question has

75    relevance to neuroscience, because object representations in the early and late layers

76    of deep networks match with early and late visual processing stages in the brain (Cichy

77    et al., 2014; Kar et al., 2019b). The layer at which this effect arises could therefore

78    reveal its underlying computational complexity and offer clues as to its neural

79    substrates.

80    Here, we identified a number of emergent perceptual and neural properties from

81    visual psychology and neuroscience that can be recast as statements about distances

82    between images in the underlying perceptual/neural representation. We then tested

83    each of these properties on a state-of-the-art deep neural network optimized for object

84    recognition. This revealed a highly interesting and insightful list of properties that were

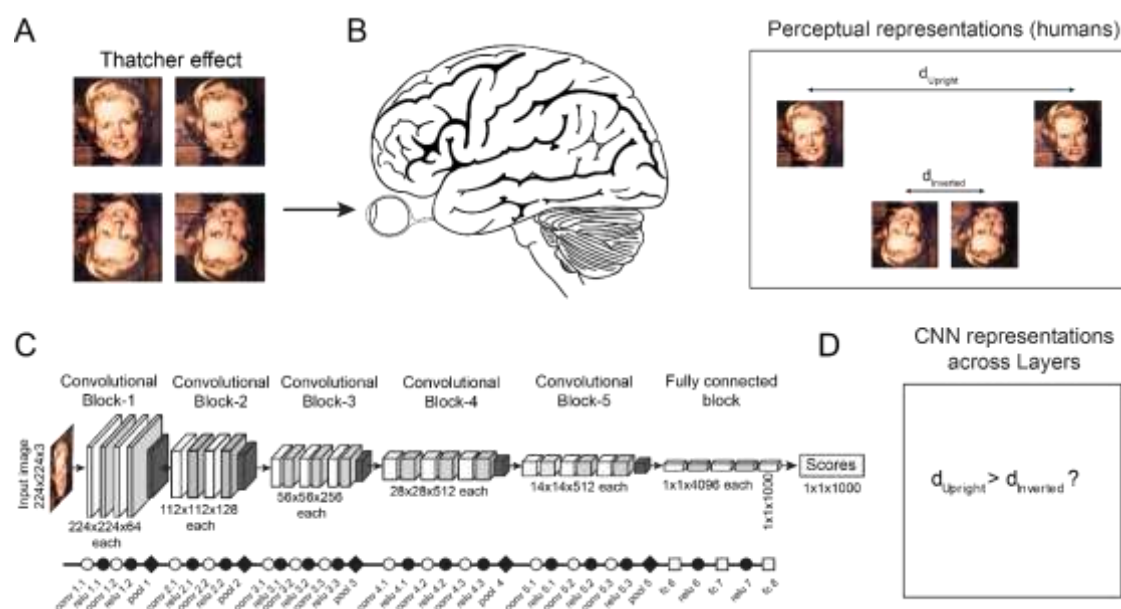85    either present or absent in the network.



86

87    **Figure 1: Evaluating whether deep networks see the way we do**
88    (A) In the classic Thatcher effect, when the parts of a face are individually inverted,
89        the face appears grotesque when upright (*top row*) but not when inverted
90        (*bottom row*). Figure credit: Reproduced with permission from Peter Thompson.
91    (B) When the brain views these images, it presumably extracts specific features
92        from each face so as to give rise to this effect. We can use this idea to recast
93        the Thatcher effect as a statement about the underlying perceptual space. The
94        distance between the normal and Thatcherized face is larger when they are
95        upright compared to when the faces are inverted. This property can easily be
96        checked for any computational model.
97    (C) Architecture of a common deep neural network (VGG-16). Symbols used here
98        and in all subsequent figures indicate the underlying mathematical operations
99        of that layer: *unfilled circle* for convolution, *filled circle* for ReLu, *diamond* for
100       maxpooling and *unfilled square* for fully connected layers. Broadly, unfilled
101       symbols depict linear operations and filled symbols depict non-linear
102       operations.
103    (D) By comparing the distance between upright and inverted Thatcherized images,
104       we can ask whether any given layer of the deep network sees a Thatcher effect.
105

**RESULTS**

106

107    We identified a large number of emergent perceptual and neural properties that

108    can be tested across layers of a deep network. We organized these properties broadly

109    into five groups: (1) those that reveal sensitivity to object or scene statistics,

110    comprising the Thatcher effect, mirror confusion and object-scene incongruence; (2)

111    neural principles observed in high-level visual cortex, related to multiple object tuning

112    and sparseness; (3) relational properties such as Weber's law, relative size and

113    surface invariant pattern processing; (4) encoding of 3d shape or scene structure and

114    (5) processing of object parts and global structure.

115    We evaluated these properties for a state-of-the-art pre-trained deep network,

116    VGG-16, optimized for object classification on the ImageNet dataset (Simonyan and

117    Zisserman, 2014). For each property, we performed an experiment in which we used

118    carefully controlled sets of images as input to the network, obtained the activations of

119    the units in each layer, and asked whether each layer shows that property. We

120    obtained qualitatively similar results on several other pre-trained feedforward networks

121    with diverse architectures: AlexNet, GoogleNet, Resnet-50 and Resnet-152 (Section

122    S1). To ensure that the results are truly due to training and not simply a consequence

123    of the architecture of the network, we also repeated these experiments on a randomly

124    initialized VGG-16 model (Section S2). For simplicity, we report our results only for the

125    VGG-16 network in the sections below, with results for all other networks detailed in

126    supplementary material (Sections S1-2).

127

128    **Experiment 1: Do deep networks see a Thatcher effect?**

129    The Thatcher effect is an elegant demonstration of how upright faces are

130    processed differently from inverted faces, presumably because we encounter mostly

131 upright faces. As detailed in the Introduction, it can be recast as a statement about the

132 underlying distances in perceptual space: that normal vs Thatcherized faces are closer

133 when inverted than when upright (Figure 2A). For each layer of the deep network

134 (VGG-16), we calculated a "Thatcher index" of the form $(d_{upright} - d_{inverted})/(d_{upright} +$

135 $d_{inverted})$, where $d_{upright}$ is the distance between normal and Thatcherized face in the

136 upright orientation, and $d_{inverted}$ is the distance between them in an inverted orientation.

137 Note that the Thatcher index for a pixel-like representation (where the activation of

138 each unit is proportional to the brightness of each pixel in the image) will be zero since

139 $d_{upright}$ and $d_{inverted}$ will be identical. For human perception, since $d_{upright} > d_{inverted}$, the

140 Thatcher index will be positive and can be estimated from previous studies (see

141 Methods).

142    Next, we calculated the Thatcher index across layers for two networks with

143 similar architecture but trained on different tasks (see Methods). The first was VGG-

144 16 which is trained for object classification (Simonyan and Zisserman, 2014). The

145 second was VGG-face which is trained for face recognition (Parkhi et al., 2015). The

146 Thatcher index for both networks across layers is shown in Figure 2B. It can be seen

147 that the VGG-16 shows a positive Thatcher index in the conv4 and conv5 layers but

148 eventually does not show a Thatcher effect in the final fully connected layers (Figure

149 2B, red curve). By contrast, the VGG-face network shows a steadily rising Thatcher

150 effect that remains close to human levels even in the fully connected layers. Thus, the

151 Thatcher effect is present only in deep networks trained on upright face recognition

152 but not on generic object recognition.

153

154

155

156    **Experiment 2: Do deep networks show mirror confusion?**

157    Mirror reflections along the vertical axis appear more similar to us than

158    reflections along the horizontal axis (Figure 2C). This effect has been observed both

159    in behaviour as well as in high-level visual cortex in monkeys (Rollenhagen and Olson,

160    2000). To assess whether deep networks show mirror confusion, we calculated a

161    mirror confusion index of the form $(d_{horizontal} - d_{vertical})/(d_{horizontal} + d_{vertical})$, where $d_{horizontal}$

162    and $d_{vertical}$ represent the distance between horizontal mirror image pairs and between

163    vertical mirror image pairs respectively (see Methods). Since vertical mirror images

164    are more similar in perception, this index would be positive. Across the deep network

165    VGG-16, we found an increasing mirror confusion index across layers (Figure 2D).

166    Thus, the deep network does experience (as we do) more mirror confusion for vertical

167    compared to horizontal mirror images.

168

169    **Experiment 3: Do deep networks show scene incongruence?**

170    Our ability to recognize an object is hampered when it is placed in an

171    incongruent context (Davenport and Potter, 2004; Munneke et al., 2013), suggesting

172    that our perception is sensitive to the statistical regularities of objects co-occurring in

173    specific scene context. To explore whether deep networks are also sensitive to scene

174    context, we gave as input the same images as tested on humans, and asked how the

175    deep network classification output changes with scene context (see Methods for more

176    details). An example object (hatchet) placed against a congruent context (forest) and

177    incongruent context (supermarket) are shown in Figure 2E. The VGG-16 network

178    returned a high probability score for the correct target object in the congruent context

179    (Figure 2E, top row) but gave a low probability score for the same object in an

180    incongruent context (Figure 2E bottom row). We obtained similar results across all

181  congruent/incongruent scene pairs: the VGG-16 top-1 accuracy dropped substantially

182  for incongruent compared to congruent contexts (drop in accuracy from congruent to

183  incongruent scenes: 20% for top-5 accuracy; 27% for top-1 accuracy; Figure 2F). On

184  the same scenes, human object naming accuracy has been reported to drop for

185  incongruent scenes, but the drop was smaller compared to the VGG-16 network (drop

186  in human accuracy from congruent to incongruent scenes: 14% in the Davenport &

187  Potter, 2004; 13% in Munneke et. al. 2013; Figure 2G). We note that assessing the

188  statistical significance of the accuracy difference in humans and deep networks is not

189  straightforward since the variations in accuracy reported are across subjects for

190  humans and across scenes for the VGG-16 network. We conclude that deep networks

191  show scene incongruence albeit to a smaller extent than humans.
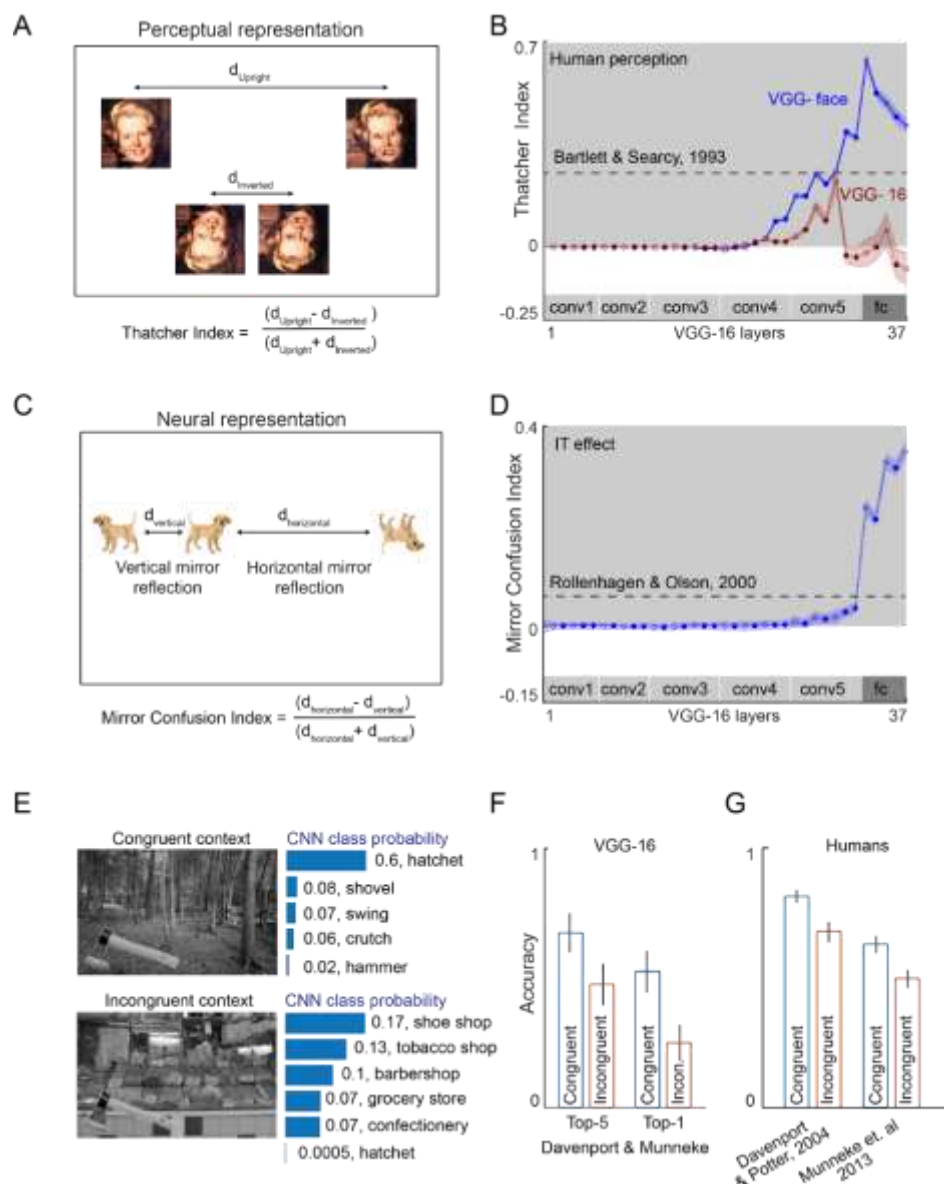
192

193

**Figure 2: Object and scene regularities in deep networks**

(A) Schematic showing the perceptual arrangement of normal and Thatcherized faces in the upright and inverted orientations. A Thatcherized face is more similar to its normal counterpart in the inverted orientation.

(B) Thatcher index (averaged across 20 pairs of normal-Thatcherized face pairs) plotted across layers for two pre-trained networks, VGG-16 (red) and VGG-face (blue). Shaded error bars indicated s.e.m across 20 face pairs. The gray zone indicates human-like performance, with dotted lines indicating the strength of the Thatcher effect estimated from humans (see Methods).

(C) Schematic showing the representation of vertical and horizontal mirror images in perception. Vertical mirror image pairs are closer than horizontal mirror image pairs. The same effect holds even if the dog is rotated by 90°, showing that this effect is not simply due to the object's axis of elongation.

(D) Mirror Confusion Index (averaged across 50 naturalistic objects and their 90° rotated versions) across layers for the VGG-16 network. Dotted lines show the strength of the mirror confusion index estimated from monkey inferior temporal neurons (see Methods).

211    (E) An example object (*hatchet*) embedded in a congruent context (*forest*) and an
212         incongruent context (*supermarket*), highlighting the vulnerability of deep
213         network classification to scene context. The CNN class probability returned by
214         the VGG-16 network is shown beside each image. It can be seen that the deep
215         network returns the correct class label on the congruent but not the incongruent
216         scene.
217    (F) Accuracy of object classification by the VGG-16 network across congruent
218         (blue) and incongruent(red) context scenes, for top-5 accuracy (*left*) and top-1
219         accuracy (*right*).  Error bars represent s.e.m across 40 congruent/incongruent
220         scene pairs.
221    (G) Accuracy of object naming by humans in congruent (blue) and incongruent (red)
222         contexts across two separate studies (Davenport and Potter, 2004; Munneke
223         et al., 2013). Humans are less accurate on naming objects placed on
224         incongruent scene contexts, but show a smaller drop compared to the VGG-16
225         network. Error bars represent s.e.m reproduced from the published studies.
226

227 **Experiment 4: Do deep networks show multiple object normalization?**

228       Next we asked whether individual units in deep networks conform to two

229 general principles observed in single neurons of high-level visual cortex. The first one

230 is multiple object normalization, whereby the neural response to multiple objects is the

231 average of the individual object responses at those locations (Zoccolan et al., 2005).

232 This principle is illustrated in Figure 3A. Note that this analysis is meaningful only for

233 units that respond to all three locations: a unit in an early layer with a small receptive

234 field would respond to objects at only one location regardless of how many other

235 objects were present in the image. To identify units that are responsive to objects at

236 each location, we calculated the variance of its activation across all objects presented

237 at that location. We then performed this analysis on units that showed a non-zero

238 response variance at all three locations, which meant units in Layer 23 (conv4.3)

239 onwards.

240       To assess whether deep networks show multiple object normalization, we

241 plotted for each unit in a given layer its response to multiple objects against the sum

242 of its responses to the individual objects. If there is multiple object normalization, the

243 slope of the resulting plot should be 1/2 for two objects and 1/3 for three objects. The

244    resulting plot is shown for Layer 37 of the VGG-16 network (Figure 3B). The overall

245    slope was 0.60 for two objects and 0.42 for three objects for all units. To evaluate this

246    effect across layers, we plotted the two-object and three-object slopes obtained in this

247    manner across layers (Figure 3C). For the later layers we observed a nearly monotonic

248    decrease in the slopes, approaching the levels observed in monkey high-level visual

249    areas (Figure 3C). We conclude that deep networks exhibit multiple object

250    normalization.

251

252    **Experiment 5: Do deep networks show selectivity across multiple dimensions?**

253    In a recent study we showed that neurons in the monkey inferior temporal cortex have

254    intrinsic constraints on their selectivity that manifests in two ways (Zhivago and Arun,

255    2016). First, neurons that respond to fewer shapes have sharper tuning to parametric

256    changes in these shapes. To assess whether units in the deep network VGG-16 show

257    this pattern, we calculated the sparseness of each unit across a reference set of

258    disparate shapes (Figure 3D), and its sparseness for parametric changes between

259    pairs of these stimuli (an example morph line is shown in Figure 3D). This revealed a

260    consistently high correlation across units of each layer in the VGG-16 network (Figure

261    3E). Second, we found that neurons that are sharply tuned across textures are also

262    sharply tuned to shapes. To assess this effect across layers, we calculated the

263    correlation across units between sparseness on textures with the sparseness on

264    shapes.  Although there was no such consistently positive correlation in the early

265    layers, we did find a positive correlation in the later (conv5 & fc) layers (Figure 3F).

266    We conclude that deep networks show selectivity along multiple dimensions just like

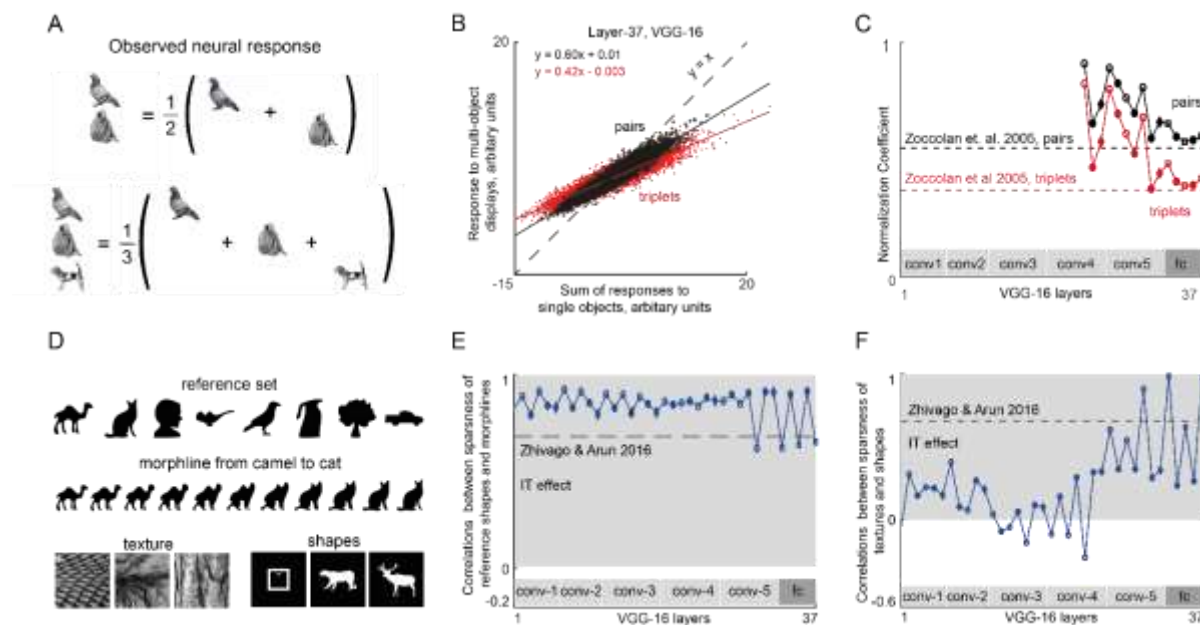267    neurons in high-level visual cortex.

**Figure 3: Single unit properties of deep networks**

(A) Schematic illustrating a general principle observed in neurons in high-level visual areas of the brain. The response of a neuron to multiple objects is typically the average of its responses to the individual objects at those locations.

(B) Response to multiple object displays plotted against the sum of the individual object responses for two-object displays (black) and three-object displays (red), across 10,000 units randomly selected from Layer 37 of the VGG-16 network.

(C) Normalization slope plotted across layers for two object displays (black) and three-object displays (red).

(D) Stimuli used to compare selectivity across multiple dimensions in a previous study of IT neurons (Zhivago and Arun, 2016).

(E) Correlation between sparseness on the reference set and maximum sparseness along morphlines across units of each layer in the VGG-16 network.

(F) Correlation between sparseness along textures and sparseness along shapes plotted across layers of the VGG-16 network.

**Experiment 6: Do deep networks show Weber's law?**

Next we asked whether deep networks are sensitive to relational properties in visual displays. The first and most widely known of these is Weber's law, which states that sensitivity to changes in any sensory quantity is proportional to the baseline level being used for comparison. Weber's law for line length is illustrated in Figure 4A. This in turn predicts that the distance between any two lines differing in length should be proportional to the relative but not absolute change in length. In a previous study, we

292    showed that this is true for humans in visual search for both length and intensity

293    changes (Pramod and Arun, 2014).

294    We therefore asked for the deep network VGG-16, whether pairwise distances

295    between lines of varying length are correlated with absolute or relative changes in

296    length. Specifically, if the correlation between pairwise distances and relative changes

297    in length is larger than the correlation with absolute changes in length, we deemed

298    that layer to exhibit Weber's law. This difference in correlation is positive for humans

299    in visual search, and we plotted this difference across layers of the VGG-16 network

300    (Figure 4B). The correlation difference was initially negative in the early layers of the

301    network, meaning that the early layers were more sensitive to absolute changes in

302    length. To our surprise, however, distances in the later layers were sensitive to relative

303    changes in length in accordance to Weber's law (Figure 4B).

304    We conclude that deep networks exhibit Weber's law for length.

305

306    **Experiment 7: Do deep networks encode relative size?**

307    We have previously shown that neurons in high-level visual areas are sensitive

308    to the relative size of items in a display (Vighneshvel and Arun, 2015). Specifically, we

309    found that, when two items in a display undergo congruent changes in size, the neural

310    response is more similar than expected given the two individual changes in size. This

311    pattern was present only in a small fraction (7%) of the neurons. This effect is

312    illustrated in Figure 4C. To explore whether this effect can be observed in a given layer

313    of the deep network VGG-16, we performed a similar analysis. We selected units in a

314    given layer with the strongest interactions (see Methods) and calculated a relative size

315    index of the form (d1-d2)/(d1+d2) where d1 is the distance between stimuli with

316    incongruent changes in size, and d2 is the distance between stimuli with congruent

317    size changes (i.e. where both items or parts are scaled up or down in size by the same

318    factor). The relative size index was calculated for each tetrad as to replicate the same

319    analysis done in the previous study(Vighneshvel and Arun, 2015). The relative size

320    index for the VGG-16 network remained close to zero in the initial layers and increased

321    modestly to a positive level in the later layers (Figure 4D). However the size of this

322    effect was far smaller than that observed in IT neurons, but in the same direction. We

323    conclude that the deep network VGG-16 encodes relative size.

324

325    **Experiment 8: Do deep networks decouple pattern shape from surface shape?**

326         A recent study showed that IT neurons respond more similarly when a pattern

327    and a surface undergo congruent changes in curvature or tilt (Ratan Murty and Arun,

328    2017). This effect is illustrated for a pattern surface pair in Figure 4E, where it can be

329    seen that the distance between incongruent pattern-surface pairs (where the pattern

330    and surface change in opposite directions) is larger than the distance between

331    congruent pairs where the pattern and surface undergo congruent changes. To assess

332    whether the deep network VGG-16 shows this property, we identified units with

333    increased interactions (see Methods) and calculated a surface invariance index of the

334    form (d1-d2)/(d1+d2), where d1 is the distance between incongruent pairs, and d2 is

335    the distance between congruent pairs. A positive value of this index for a given layer

336    implies that the layer shows surface invariance. However the surface invariance index

337    was consistently below zero across layers for the VGG-16 network (Figure 4F). We

338    conclude that the VGG-16 network does not show surface invariance.
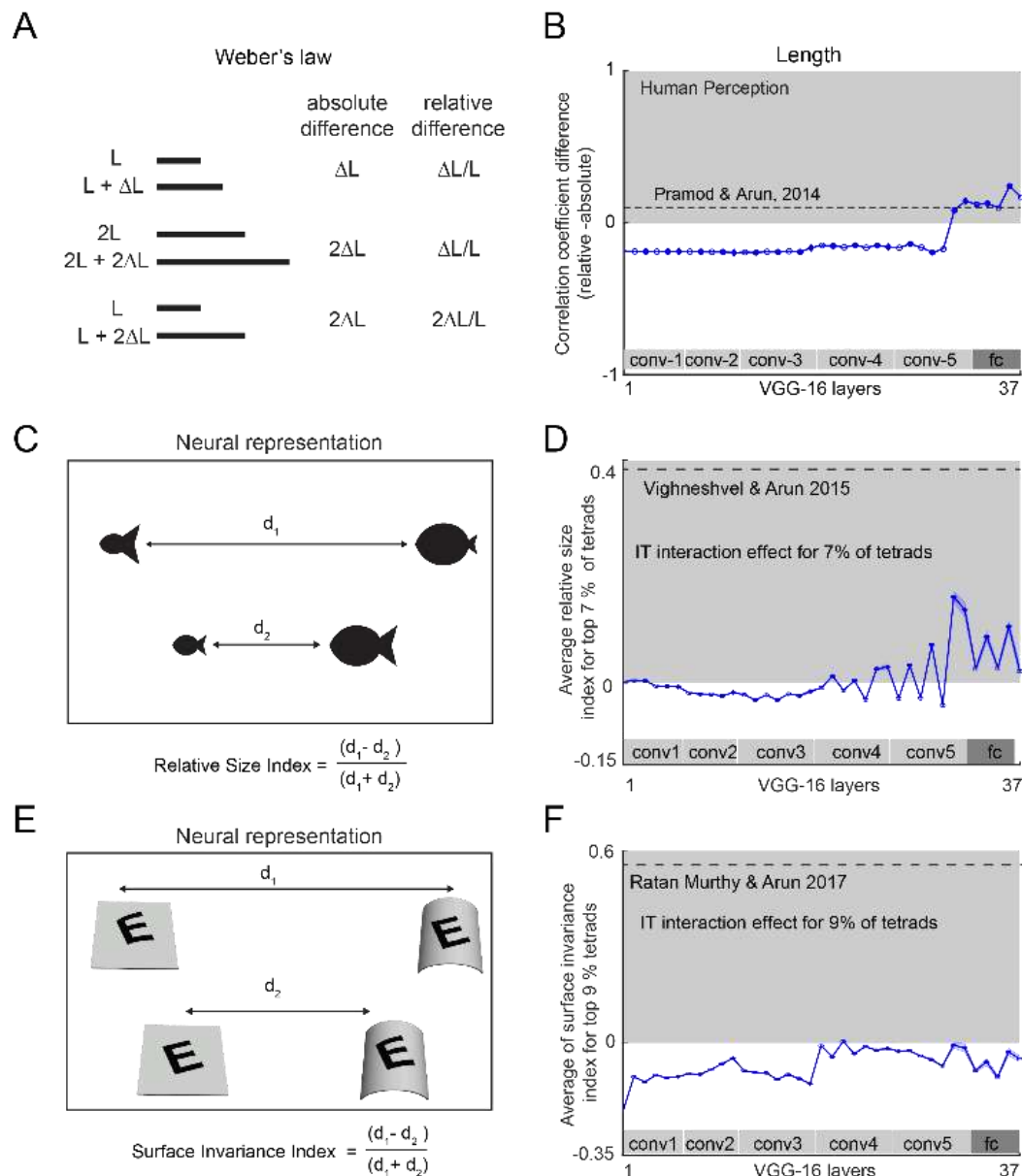
339

**Figure 4: Relational properties in deep networks**

(A) Example illustrating the Weber's law for line length. Although the original statement of Weber's law is that the just-noticeable difference in length will depend on the baseline length, it also applies to perceptual distances between any two stimuli (Pramod and Arun, 2014). In this formulation, the perceptual distance between two lines differing in length will be proportional to the relative change in length rather than the absolute change.

(B) To calculate a single quantity that measures adherence to Weber's law, we calculated the difference in correlation between pairwise distances between lines with the relative and absolute differences in line length. A positive difference indicates adherence to Weber's law (indicated by the gray shading). This difference in correlation is plotted across layers for line length.

(C) Schematic of the relative size encoding observed in monkey IT neurons (Vighneshvel and Arun, 2015). For a fraction of neurons, the distance between two-part objects when both parts covary in size is smaller than the distance

356       when they show inconsistent changes in size. Thus, these neurons are
357       sensitive to the relative size of items in a display.

358 (D) Relative size index across units with interaction effects (averaged across top
359       7% tetrads  tetrads, error bars representing s.e.m) across layers of the VGG-
360       16 network. Dotted lines show the strength of the relative index estimated from
361       monkey inferior temporal neurons (Vighneshvel and Arun, 2015).

362 (E) Schematic of the surface invariance index observed in monkey IT neurons
363       (Ratan Murty and Arun, 2017). For a fraction of neurons, the distance between
364       two stimuli with congruent changes in pattern and surface curvature is smaller
365       than between two stimuli with incongruent pattern/surface changes. Thus,
366       these neurons decouple pattern shape from surface shape.

367 (F) Surface invariance index across units with interaction effects (averaged across
368       top 9% pattern/surface tetrads, error bars representing s.e.m) across layers of
369       the VGG-16 network. Dotted lines show the strength of the surface invaraince
370       index estimated from monkey inferior temporal neurons (Ratan Murty and Arun,
371       2017).

372

## Experiment 9: Do deep networks show 3d processing?

We are sensitive to three-dimensional shape and not simply two-dimensional contours in the image. This was demonstrated in an elegant experiment in which search for a target differing in 3d structure is easy whereas search for a target with the same difference in 2d features is hard (Enns and Rensink, 1990, 1991). This effect can be recast as a statement about distances in perceptual space as illustrated in Figure 5A. All three pairs of shapes depicted in Figure 5A differ in the same Y-shaped feature, but the two cuboids are more dissimilar because they differ also in 3d shape. To assess whether units in a given layer of the deep network show this effect, we calculated a 3d processing index of the form $(d1-d2)/(d1+d2)$ where $d1$ is the distance between the cuboids and $d2$ is the distance between the two equivalent 2d conditions. A positive 3D processing index indicates an effect similar to human perception. However we found that the 3D processing index was consistently near zero or negative across all layers of the VGG-16 network (Figure 5B). We conclude that deep networks are not sensitive to 3d shape. We speculate that explicit training of deep networks on 3d shape processing may be required for the network to exhibit this effect.

389 **Experiment 10: Do deep networks understand occlusions?**

390       A classic finding in human perception is that we automatically process occlusion

391 relations between objects (Rensink and Enns, 1998). Specifically, search for a target

392 containing occluded objects among distractors that contain the same objects

393 unoccluded is hard, whereas searching for the equivalent 2d feature difference is

394 much easier (Figure 5C, top row). Likewise, searching for a target that is different in

395 the order of occlusion is hard, whereas searching for the equivalent 2d feature

396 difference is easy (Figure 5C, bottom row). These observations demonstrate that our

397 visual system has a similar representation for occluded and unoccluded objects.

398       We therefore asked whether similar effects could be observed in the VGG-16

399 network, by calculating an occlusion index of the form $(d2-d1)/(d2+d1)$ where d1 is the

400 distance between two displays that are equivalent except for occlusion, and d2 is the

401 distance between equivalent displays with the same 2d feature difference. A positive

402 occlusion index implies an effect similar to human perception. The occlusion index

403 remained consistently below across all layers of the VGG-16 network, for both the

404 occlusion effect and occlusion ordering (Figure 5D). We conclude that deep networks

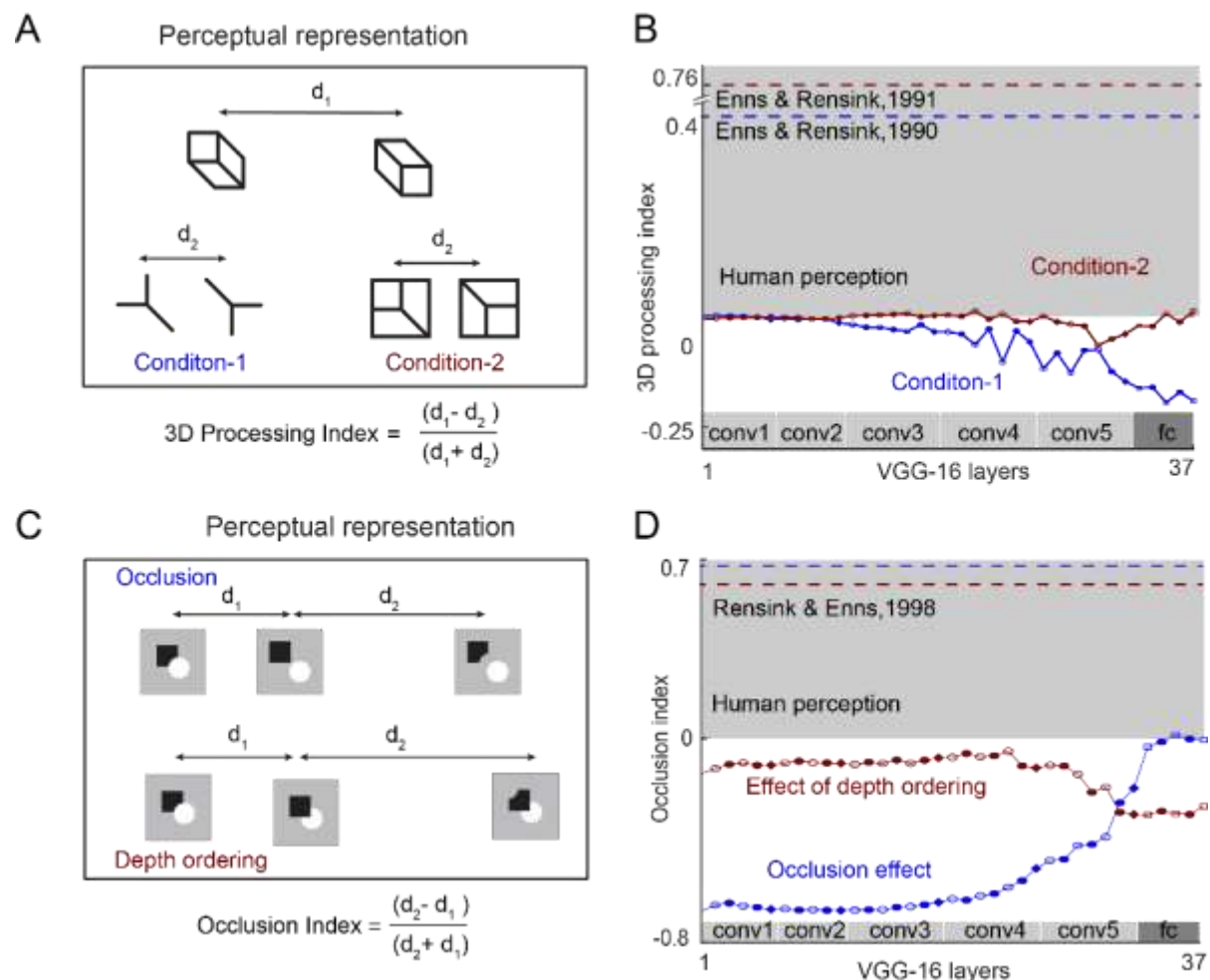405 do not understand occlusions.

406

**Figure 5: 3D processing in deep networks**

(A) Schematic demonstrating sensitivity to 3d shape in our perception (Enns and Rensink, 1991). Three equivalent image pairs are shown in perceptual space. The first image pair (with distance marked d1) consists of two cuboids at different orientations, and corresponds to an easy visual search i.e. the two objects are less similar. The second pair (marked with distance d2) contains the same feature difference as in the first pair, but represents a hard search i.e. is perceived as much more similar. Likewise, the third pair, with the same feature difference as the other two, is a hard search i.e. perceived as similar.

(B) 3D processing index for the VGG-16 network across layers, for condition 1 ($d_1$ vs $d_2$, *blue*) and condition 2 ($d_1$ vs $d_3$, *red*). Dotted lines of the corresponding color represent the estimated human effect size.

(C) Schematic showing processing of occlusions in perception. *Top:* A square alongside a disk is perceptually similar to a display with the same objects but occluded, but dissimilar to a 2d control image with an equivalent feature difference. *Bottom:* A square occluding a disk or disk occluding square are perceptually similar, but dissimilar to an equivalent 2d control with the same set of feature differences.

(D) Occlusion index for both basic occlusion (*blue*) and depth ordering (*red*) for each layer of the VGG-16 network.

**Experiment 11: Do deep networks understand object parts?**

430

431     We not only recognize objects but are able to easily describe their parts. We

432     conducted two related experiments to investigate part processing in deep networks.

433     In Experiment 11A, we compared deep network feature representations for whole

434     objects and for the same object broken down into either natural or unnatural parts. In

435     perception, searching for an object broken into its natural parts with the original object

436     as distractors is much harder than searching for the same object broken at an

437     unnatural location (Xu and Singh, 2002). This result is depicted schematically in terms

438     of the underlying distances (Figure 6A). To assess whether the VGG-16 network also

439     shows this part decomposition, we calculated a part processing index of the form ($d_u$

440     $- d_n$)/($d_u + d_n$) where $d_u$ is the distance between the original object and the object

441     broken at an unnatural location, and $d_n$ is the distance between the original object and

442     the same object broken at a natural location. A positive part processing index implies

443     an effect similar to that seen in perception. The part processing index across layers of

444     the VGG-16 network is depicted in Figure 6B. We found that the index becomes

445     positive in the intermediate layers, but becomes negative in the subsequent layers

446     (conv4/conv5 onwards).

447     In Experiment 11B, we asked what happens to objects that can be decomposed

448     into two possible ways without introducing a break (Figure 6C). In visual search,

449     search between pairs of whole objects is better explained by breaking the object down

450     into its natural parts compared to its unnatural parts (Pramod and Arun, 2016b). To

451     capture this effect, we defined the natural part advantage as the difference in model

452     correlation (see Methods) between natural and unnatural parts. A positive value

453     indicates effects similar to perception. This natural part advantage is shown across

454     layers of the VGG-16 network in Figure 6D. It can be seen that there is little or no

455    advantage for natural parts in most layers except temporarily in the later layers

456    (conv5/fc).

457        Based on Experiments 11A & 11B, we conclude that the VGG-16 network

458    shows no systematic part processing.

459

460    **Experiment 12: Do deep networks show a global shape advantage?**

461        In perception a classic finding is that we see the forest before the trees, i.e. we

462    can detect global shape before local shape (Navon, 1977; Kimchi, 1994). We can

463    recast this effect into a statement about distances in perception: the distance between

464    two hierarchical stimuli differing only in global shape will be larger than the distance

465    between two such stimuli differing only in local shape. This is depicted schematically

466    in Figure 6E. To calculate a single measure for this effect, we defined a global

467    advantage index as $(d_{global} - d_{local})/(d_{global} + d_{local})$, where $d_{global}$ is the average distance

468    between all image pairs differing only in global shape and $d_{local}$ is the average distance

469    between all image pairs differing only in local shape. A positive global advantage index

470    implies an effect similar to perception. The global advantage index is depicted across

471    layers of the VGG-16 network in Figure 6F. While there is a slight global advantage in

472    the initial layers, the network representation swings rapidly in the later layers towards

473    the opposite extreme, which is a local advantage. Thus, deep networks see the trees
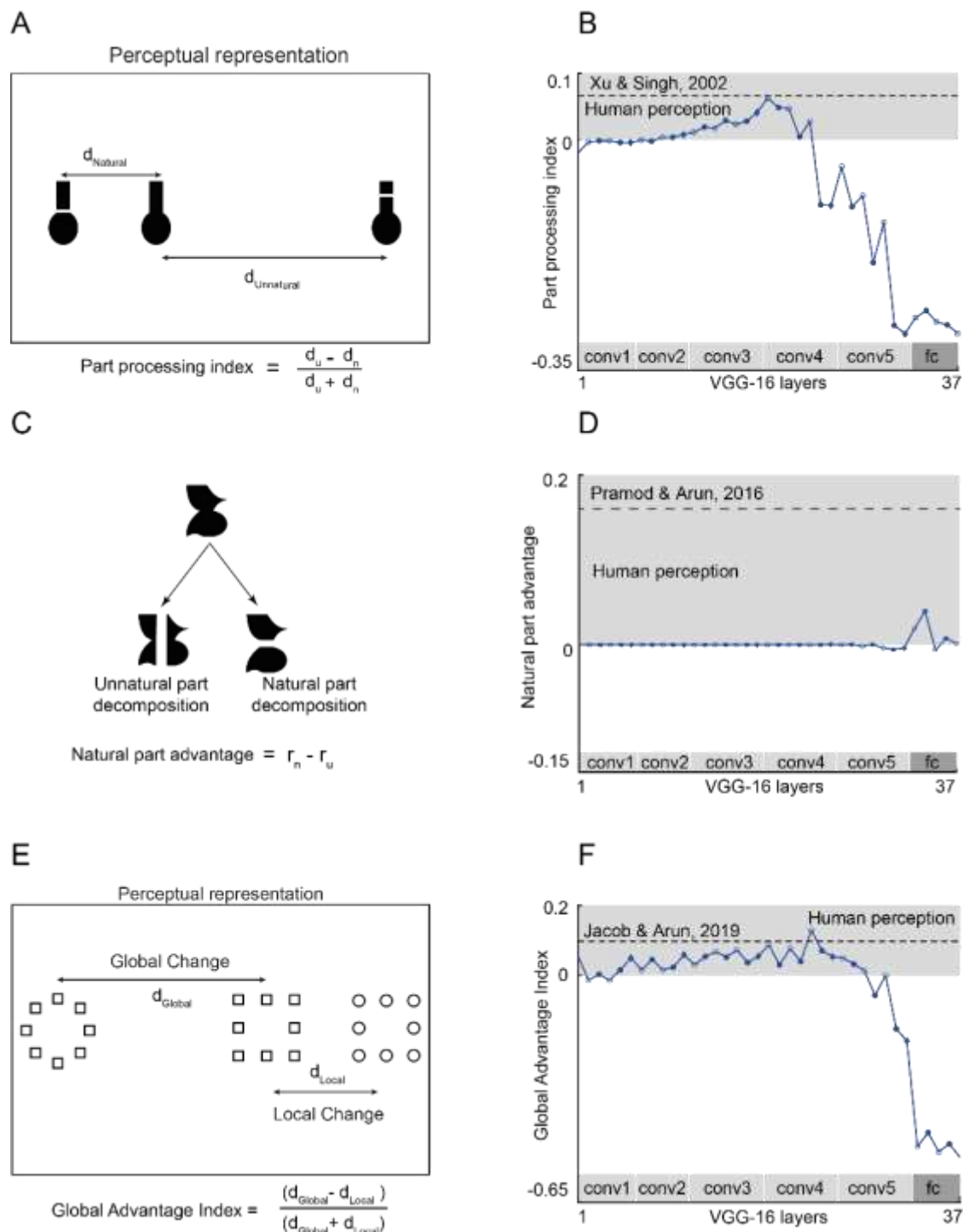
474    and not the forest.

**Figure 6: Part-whole relations in deep networks**

(A) Schematic showing the perceptual representation of objects with a break introduced either at natural or unnatural parts.

(B) Part processing index across layers of the VGG-16 network. The shaded gray bar represents effects similar to human perception, and the dotted line represents the effect size estimated from human visual search (Xu and Singh, 2002).

(C) Schematic showing how the same object can be broken into either natural or unnatural parts.

(D) Natural part advantage across layers of the VGG-16 network. The shaded gray bar represents effects similar to human perception, and the dotted line

487    represents the effect size estimated from human visual search (Pramod and
488    Arun, 2016b).
489 (E) Schematic showing the perceptual representation of hierarchical stimuli. The
490    left and middle images differ only in global shape whereas the middle and right
491    images differ only at the local level.
492 (F) Global Advantage index across layers of the VGG-16 network. The shaded gray
493    bar represents effects similar to human perception, and the dotted line
494    represents the effect size estimated from human visual search (Jacob and
495    Arun, 2019).

**DISCUSSION**

496

497    Because object representations in deep networks match coarsely with the

498    brain, it is widely believed that any remaining differences between brains and deep

499    networks must be of degree but not of kind. Here, we show that this is not always the

500    case: In some respects, deep networks do see the way we do. They exhibit perceptual

501    phenomena such as the Thatcher effect, mirror confusion, scene incongruence and

502    Weber's law. Their units show multiple object normalization, sparseness along multiple

503    dimensions and encode relative size. Yet in other ways, they don't see the way we do.

504    They fail to encode relational properties such as surface invariance, do not show

505    processing of 3d features, occlusions or natural parts and do not show a global

506    advantage.

507    These findings are important for several reasons. First, they describe the

508    similarities and differences between our vision and deep networks, and challenge the

509    prevailing belief that they can be treated as accurate models for biological vision.

510    Second, they show that object recognition training alone is sufficient to produce some

511    emergent properties but not others, thereby elucidating the computational problem of

512    vision itself. Finally, the missing properties could be incorporated as training or

513    architecture constraints on deep networks to yield better or more robust performance.

514    Below we discuss our findings in the context of the existing literature.

515    We begin by discussing some general concerns regarding our findings. First

516    and foremost, it could be argued that our results are based on testing with artificial

517    objects or images, and that it is unreasonable to expect deep networks to respond

518    sensibly to unnatural images. However, these concerns apply equally to humans as

519    well, who in fact do respond sensibly to these artificial displays despite no prior

520    exposure. Indeed, there is a long tradition in psychology and neuroscience of using

521    artificial images to elucidate visual processing (Rust and Movshon, 2005). Second, it

522    could be argued that deep networks could potentially be trained to report all the tested

523    properties. However, such a finding would only be circular if the network did indeed

524    exhibit the property it was trained for. We do note however that it would be interesting

525    if deep networks were unable to learn certain properties. Indeed, certain relational

526    properties have been reported as difficult to learn by computer vision algorithms

527    (Fleuret et al., 2011b), although this study did not evaluate deep neural networks. By

528    contrast, we consider our experiments to be far more interesting, since they reveal

529    that deep networks trained for specific other purposes show emergent properties that

530    they were not explicitly trained for, such as Weber's law.

531      Our finding that deep networks trained for object recognition exhibit Weber's

532    law or encode relative size is puzzling at first glance. Why would the demands of

533    recognizing objects require sensitivity to relative changes? One possibility is that

534    object recognition requires a representation invariant to changes in size, position,

535    viewpoint and even illumination of objects in the image, which in turn requires

536    processing all object features relative to the surrounding features. This could be tested

537    by training deep networks on controlled sets of images with variations of one kind but

538    not the other. We note that there could be other visual task requirements that could

539    also give rise to Weber's law (Pardo-Vazquez et al., 2019).

540      Our finding that deep networks exhibit the Thatcher effect, mirror confusion and

541    scene incongruence are consistent with them being sensitive to image regularities in

542    scenes. In fact, the VGG-16 network may be over-reliant on scene context, because

543    it showed a larger drop in accuracy for incongruent scenes compared to humans

544    (Figure 2F). This is consistent with a previous study in which human scene

545    expectations benefited deep network performance (Katti et al., 2019). However, the

546    VGG-16 network did not exhibit 3d processing, occlusions and surface invariance,

547    suggesting that these properties might emerge only with additional task demands such

548    as extracting 3d shape from the image. Likewise, the absence of any part processing

549    or global advantage in the network suggests that these too are properties that might

550    emerge with additional task demands, such as part recognition or global form

551    recognition (Belongie et al., 2002).

552         We have found that deep networks do not show a global advantage effect but

553    instead seem to process local features. This finding is surprising considering that units

554    in later layers receive convergent inputs from the entire image. However, our finding

555    is consistent with reports of a bias towards local object texture in deep networks

556    (Geirhos et al., 2018a). It is also consistent with the large perturbations in classification

557    observed when new objects are added to a scene (Rosenfeld et al., 2018), which

558    presumably change the distribution of local features. Our finding that deep networks

559    experience large scene incongruence effects is therefore likely to be due to

560    mismatched local features rather than global features. Indeed, incorporating scene

561    expectations from humans (presumably driven by global features) can lead to

562    substantial improvements in object recognition (Katti et al., 2019). Finally, a reliance

563    on processing local features is probably what makes deep networks detect

564    incongruously large objects in scenes better than humans (Eckstein et al., 2017). We

565    speculate that training on global shape could make deep networks more robust and

566    human-like in their performance.

567         Comparing human vision and deep neural networks depends critically on the

568    choice of network architecture, learning algorithm, dataset and the task learned. An

569    important but neglected finding is that even a randomly initialized network can

570    generate features useful for certain tasks. For tasks like texture generation and

571 discrimination, even a randomly initialized network can yield reliable features (Jarrett

572 et al., 2009; Mongia et al., 2016). Likewise, a significant amount of explainable

573 variance in the early visual responses of MEG signals in humans can be predicted

574 using features extracted from a randomly initialized and untrained deep neural network

575 (Cichy et al., 2016). By contrast, we have found that most perceptual effects are absent

576 in randomly initialized networks, except for global advantage (Section S2). Thus,

577 training on object classification abolishes the global advantage and introduces

578 sensitivity to local features.

579 How can deep networks be made to match neural and perceptual

580 representations? There could be several ways of doing so. The first and perhaps most

581 promising direction would be to explicitly train deep networks to produce such

582 properties in addition to categorisation (Ruder, 2017). Another alternative would be to

583 train deep networks on tasks such as navigation or agent-object interaction rather than

584 (or in addition to) object recognition as this is ostensibly what humans also do (Haber

585 et al., 2018; Yang et al., 2019).

586 Finally, we note that deep networks are notorious for their susceptibility to

587 adversarial attacks. State-of-the-art deep neural networks have been shown to fail

588 catastrophically when input images are subjected to carefully constructed changes

589 that are barely perceivable to human eyes (Szegedy et al., 2013; Su et al., 2019).

590 Likewise, deep networks can give erroneous predictions on completely nonsensical

591 images (Nguyen et al., 2015). Finally, realistic multi-part 3D objects are consistently

592 misclassified by deep networks across viewpoint changes (Athalye and Carlini, 2018).

593 What could underlie such unusual behaviour? One possible reason could be the

594 tendency for deep networks to prioritize local features as described earlier. We

595 speculate that training deep networks to exhibit all the perceptual and neural

596    properties described in this study might not only improve their performance but also

597    make them more robust to adversarial attacks.

598          **METHODS**

599     *VGG-16 network architecture & training*

600          All experiments were performed on the VGG-16 network, a feedforward pre-

601     trained deep convolutional neural network trained for object classification on the

602     ImageNet dataset(Deng et al., 2009). We briefly describe the network architecture

603     here, and the reader is referred to the original paper for more details (Simonyan and

604     Zisserman, 2014). The input to the network is an RGB image of size 224 x 224 x 3

605     and the final output is a vector of confidence scores across 1000 categories.  We

606     subtracted the mean RGB value across all images (mean values across all images:

607     R=123.68, G=116.78, B=103.94). The image is passed through a stack of

608     convolutional filters (Figure 1C), where the initial layers have small receptive field (3x3)

609     and later layers are fully connected. A non-linear rectification (ReLu) operation is done

610     after each convolution operation. Five max-pooling layers are present to spatial pool

611     the maximum signal over a 2x2 window of neurons. We used the MATLAB-based

612     MatConvNet software platform (Vedaldi and Lenc, 2014) to extract features and do

613     the analysis. In addition to VGG-16, we also used VGG-face which has the same

614     architecture but trained instead on face identification (Parkhi et al., 2015).

615

616     *Feature Extraction*

617          For each image, we passed it as input into the network, stored the activations

618     of each layer as a column vector. Hence, a single image we will have 37 feature

619     vectors (one column vector from each layer). To calculate the distance between

620     images A and B, we calculated the Euclidean distance between the corresponding

621     activation vectors.

622     In all the experiments, we define specific measure based on distances and this

623     quantity is plotted across layers with a specific chain of symbols as shown in Figure

624     1C. Symbols used indicates the underlying mathematical operations done in that layer:

625     unfilled circle for convolution, filled circle for ReLu, diamond for maxpooling and

626     unfilled square for fully connected layers. Broadly, filled symbols denote linear

627     operations and unfilled ones indicate non-linear operations.

628

629     **Experiment 1: Thatcher effect**

630     The stimuli comprised 20 co-registered Indian faces (19 male, 1 female) from

631     the IISc Indian face dataset (Katti and Arun, 2019). All faces were grayscale, upright

632     and front-facing. To Thatcherize a face, we inverted the eyes and mouth while keeping

633     rest of the face intact. We implemented inversion by first registering facial landmarks

634     on frontal faces using an Active appearance model-based algorithm (Cootes et al.,

635     2001). Briefly, this method models face appearance as a two-dimensional mesh with

636     76 nodes, each node represents local visual properties of stereotyped locations such

637     as corners of eyes, nose, and mouth. We then defined bounding boxes for left and

638     right eye as well as mouth, by identifying landmarks that correspond to the four corners

639     of each box. We then locally inverted eye and mouth shape by replacing the top row

640     of eye or mouth image pixels by the last row and likewise repeating this procedure for

641     each pair of equidistant pixels rows above and below the middle of the local region.

642     The inversion procedure was implemented as a custom script written in MATLAB. The

643     full stimulus set is shown in Section S3.

644     To calculate a single measure for the Thatcher effect, we calculated a Thatcher

645     index defined as $\frac{d_{upright} - d_{inverted}}{d_{upright} + d_{inverted}}$ , where $d_{upright}$ is the distance between an normal

646     face and Thatcherized face in upright orientation and $d_{inverted}$ is the distance between

647    a normal face and Thatcherized face in inverted orientation. We estimated the

648    Thatcher index for humans from the similarity ratings reported from humans albeit on

649    a different set of faces (Bartlett and Searcy, 1993). We calculated Thatcher index after

650    converting the similarity rating (humans gave a rating between 1 to 7 on pair of images)

651    into a dissimilarity rating (dissimilarity rating = 7- similarity rating).

652

653    **Experiment 2: Mirror Confusion**

654         The stimuli consisted of 100 objects (50 naturalistic objects and 50 versions of

655    these objects made by rotating each one by 90°). This was done to avoid any effect

656    due to the objects own axis of elongation. We created a horizontal and vertical mirror

657    image of each object. We then gave as input the original image and the two mirror

658    images to the VGG-16 network and calculated for each layer the distance between the

659    object and two mirror images. The full stimulus set is shown in Section S3.

660         To calculate a single measure for mirror confusion, we defined a mirror

661    confusion index of the form $\frac{d_{horizontal} - d_{vertical}}{d_{horizontal} + d_{vertical}}$ , where $d_{horizontal}$ is the distance

662    between an object and its horizontal mirror image and $d_{vertical}$ is the distance between

663    an object and its vertical mirror image. We estimated the strength of mirror confusion

664    index in the brain using previously published data from monkey IT neurons

665    (Rollenhagen and Olson, 2000). Specifically, we took $d_{horizontal}$ to be the reported

666    average firing rate difference between the original objects and its horizontal mirror

667    image, and analogously for $d_{vertical}$.

668

669    **Experiment 3: Scene Incongruence**

670         The stimuli consisted of 40 objects which was taken from previous studies: 17

671    objects were from the Davenport study (Davenport and Potter, 2004) and the

672    remaining 33 from the Munneke study (Munneke et al., 2013). We discarded a few

673    objects from each set since they did not have a matching category in the ImageNet

674    database. Each object was embedded against a congruent and an incongruent

675    background. The full stimulus set is shown in Section S3.

676        We measured the classification accuracy (Top-1 and Top-5) of the VGG-16

677    network for the objects pasted onto congruent and incongruent scenes. The final layer

678    of VGG-16 (Layer 38) returns a probability score for all 1000 categories in the

679    ImageNet database. The top-1 accuracy is calculated as the average accuracy with

680    which the object class with the highest probability matches the ground-truth object

681    label. The top-5 accuracy is calculated as the average accuracy with which the ground-

682    truth object is present among the object classes with the top 5 probability values. We

683    report the human (object naming) accuracy on the same dataset from previous studies

684    (Davenport and Potter, 2004; Munneke et al., 2013).

685

686    **Experiment 4: Multiple object normalization**

687        The stimuli consisted of forty-nine natural grayscale images and placed them

688    at three different locations in the image (Figure 4A). We have 147 (49 x 3 = 147)

689    singletons and randomly selected 200 pairs and 200 triplet composites. We extracted

690    features for all images (singletons, pairs and triplets) from every layer of the CNN. We

691    selected a unit for further analysis only if the unit responded differently to at least one

692    of the images in all the three positions. We then plotted the sum of activations of

693    selected units in a layer to the singleton images against the activation for the

694    corresponding pairs (or triplets). The slope of this scatterplot across layers was used

695    to infer the nature of normalization in CNNs – a slope of 0.5 for pairs and 0.33 for

696     triplets indicated divisive normalization matching that observed in high-level visual

697     cortex. The full stimulus set is shown in Section S3.

698

699     **Experiment 5: Selectivity along multiple dimensions**

700     Here we used the stimuli used in a previous study to assess the selectivity of

701     IT neurons along multiple dimensions (Zhivago & Arun, 2016). These stimuli consisted

702     of 8 reference shapes (Figure 3D; top row) and created intermediate parametric

703     morphs between pairs of these shapes (Figure 3D; example morph between camel to

704     cat). In addition, to compare texture and shape selectivity, we used 128 natural

705     textures and 128 silhouette shapes. The full stimulus set is shown in Section S3.

706     As before we calculated the activation of every layer of the VGG-16 network to

707     each of the above stimuli are input. Visually active neurons were selected by finding

708     units with a non-zero variance across this stimulus set. We found the visually active

709     neurons for each set separately and we selected a unit for further analysis only if that

710     unit is visually active for both sets. The response of each unit was normalized between

711     0 and 1. We then calculated the sparseness of each unit across different stimulus sets:

712     the reference set, the four morphlines, shape set and texture set. For a given stimulus

713     set with responses $r_1, r_2, r_3, \ldots r_n$, where n is the number of stimuli, the sparseness is

714     defined as follows: $S = (1 - \frac{\left(\frac{\Sigma r_i}{n}\right)^2}{\Sigma \frac{r_i^2}{n}}) / \left(1 - \frac{1}{n}\right)$  (Vinje et al., 2000; Zhivago and Arun,

715     2016). We then calculated the correlation across neurons between the sparseness on

716     one stimulus set versus another stimulus set.

717

718

719

**Experiment 6: Weber's Law**

720

721     To test for the presence of Weber's law in the deep network, we manipulated

722     two elementary features: length and intensity. We created two sets of images, one

723     varying in the length and the other varying in brightness of a rectangular bar. We chose

724     the values of lengths and brightness's such that the feature difference computed on

725     pairs of images spanned a wide range both in terms of absolute as well as relative

726     differences. The full stimulus set is shown in Section S3.

727     For each layer of the neural network, we extracted the pattern of activations for

728     each image and then computed the pairwise activation pattern dissimilarity for all pairs

729     of images. We then computed the correlation between pattern dissimilarities and

730     actual feature differences (i.e, difference between the actual lengths or brightnesses

731     of the rectangular bars in the two images). We computed this correlation for both

732     absolute (denoted by $r_{abs}$) as well as for relative feature differences (denoted as $r_{rel}$).

733     A positive value for the difference between the two correlation coefficients ($r_{rel} - r_{abs}$)

734     indicated that the Weber's law was present in a given layer of the neural network.

735     We also analysed deep networks for the presence of Weber's law for image

736     intensity, but found highly inconsistent and variable effects. Specifically, the pre-

737     trained VGG-16 network showed Weber's law for low image intensity levels but not for

738     high intensity levels.

739

**Experiment 7: Relative Size**

740

741     We used the stimuli used in a previous study to test whether units in the VGG-

742     16 network encode relative size (Vighneshvel and Arun, 2015). This stimulus set

743     consisted of 24 tetrads. A sample tetrad is shown in Figure 2D, with the stimuli

744    arranged such that images that elicit similar activity are closer. The full stimulus set is

745    shown in Section S3.

746         In our previous study (Vighneshvel and Arun, 2015), only a small fraction of

747    neurons (around 7%) encoded relative size. To identify similar neurons in the deep

748    network, we first identified the visually responsive units by taking all units with a non-

749    zero variance across the stimuli. For each unit and each tetrad, we calculated a

750    measure of size interactions of the form abs(r11 + r22 – r12 – r21), where r11 is the

751    response to both parts at size 1, r12 is the response to part 1 at size 1 and part 2 at

752    size 2 etc. We then selected the top 7% of all tetrads with the largest interaction effect.

753    Note that this step of selection does not guarantee the direction of the relative size

754    effect. For the selected tetrads, we calculated the relative size index, defined as $\frac{d_1 - d_2}{d_1 + d_2}$

755    where $d_1$ and $d_2$ are distances between the incongruent and congruent stimuli

756    respectively.

757

758    **Experiment 8: Decouple pattern shape from surface shape**

759         The stimuli consisted of six patterns superimposed on four surfaces. Each

760    pattern-surface pair was used to create a tetrad of stimuli as depicted in Figure 4E.

761    The full stimulus set consisted of 24 tetrads, which were a subset of those tested in

762    our previous study (Ratan Murty and Arun, 2017). The full stimulus set is shown in

763    Section S3.

764         In each VGG-16 layer, we selected visually responsive neurons and normalized

765    their responses across all stimuli as before in Experiment 5. We then selected the top

766    9% of all tetrads with an interaction effect calculated as before, as with the previous

767    study (Ratan Murty and Arun, 2017).  For all the selected tetrads we calculated the

768    surface invariance index, defined as $\frac{d_1 - d_2}{d_1 + d_2}$ where $d_1$ and $d_2$ are distances between

769    incongruent and congruent stimuli.

770

771    **Experiment 9: 3D processing**

772    We investigated 3D processing in the VGG-16 network by comparing line

773    drawing stimuli used in a previous perceptual study (Enns and Rensink, 1991). We

774    compared three pairs of shapes: cuboid, cube and frustum of square in isometric view

775    with the corresponding Y junctions. The full stimulus set is shown in Section S3. For

776    each shape, we calculated three distances between equivalent shape pairs with the

777    same feature difference (Figure 5A). We calculated a 3D processing index, defined as

778    $\frac{d_1 - d_2}{d_1 + d_2}$ where $d_1$ and $d_2$ are distances between the 3D shape and control conditions

779    respectively.

780

781    **Experiment 10: Occlusions**

782    We recreated the stimulus set used in a previous study (Rensink and Enns,

783    1998) as depicted in Figure 5C. The full stimulus set is shown in Section S3. As before

784    we compared the distance between two pairs of shapes: a pair that differed in

785    occlusion status (occluded vs unoccluded, or two images that differed in their order of

786    occlusion), and an equivalent 2D feature control containing the same feature

787    difference. We then calculated an occlusion index defined as $\frac{d_1 - d_2}{d_1 + d_2}$ where $d_1$ and $d_2$

788    are the distances between the occluded and control conditions respectively.

789

790

791

**Experiment 11: Object Parts**

We performed two experiments to investigate part processing in deep networks. In Experiment 11A, we tested what happens when a break is introduced into an object at a natural cut or an unnatural cut (Xu and Singh, 2002). The full stimulus set is shown in Section S3. The critical comparison is shown in Figure 6A. For each layer of the CNN, we extracted features for the three objects and computed the distance of the intact object with each of the broken objects ($d_n$ and $d_u$ denote distances to the broken objects with natural and unnatural parts respectively). We then computed a part processing index, defined as $\frac{d_u - d_n}{d_u + d_n}$.

In Experiment 11B, we asked whether whole object dissimilarities computed on CNN feature representations could be understood as a linear combination of dissimilarities between their natural or unnatural part decompositions as reported previously for visual search (Pramod and Arun, 2016b). We considered seven whole objects that could be broken down into either natural or unnatural parts and recombined the parts to form other objects. That is, we created two sets each containing 49 objects made either from natural or unnatural parts of the original seven objects. The full stimulus set is shown in Section S3. We then selected 492 pairs of objects from each set (including all 21 pairs from the common set) and calculated the feature distances from each layer of the CNN. We fit a part summation model to explain pairwise whole-object distances as a function of pairwise part relations, as described previously (Pramod and Arun, 2016b). We then compared model performance on the 21 pairwise distances between the common objects. We denoted by $r_{natural}$ the correlation between observed and predicted distances assuming natural part decomposition and by $r_{unnatural}$ the model correlation assuming unnatural part

816    decomposition. The natural part advantage was computed as ($r_{natural} - r_{unnatural}$). The

817    same measure was computed for human perception.

818

819    **Experiment 12: Global Shape Advantage**

820         We created a set of 49 hierarchical stimuli by combining seven shapes at global

821    scale and the same seven shapes at the local scale (Jacob and Arun, 2019). The full

822    stimulus set is shown in Section S3. We extracted features from all layers of CNNs

823    and calculated the Euclidean distance between all pairs of hierarchical shapes. We

824    calculated the global distance as the mean distance between image pairs having only

825    global change.   Similarly, we calculated the local distance as the mean distance

826    between image pairs having only local change. A sample global/local change pair is

827    shown in Figure 6E. We calculated a global advantage index as $\frac{d_{Global} - d_{Local}}{d_{Global} + d_{Local}}$.

**REFERENCES**

Athalye A, Carlini N (2018) On the Robustness of the CVPR 2018 White-Box Adversarial Example Defenses. arXiv:2–3.

Baker N, Lu H, Erlikhman G, Kellman PJ (2018) Deep convolutional networks do not classify based on global object shape Einhäuser W, ed. PLOS Comput Biol 14:e1006613.

Bartlett JC, Searcy J (1993) Inversion and configuration of faces. Cogn Psychol 25:281–316.

Belongie S, Malik J, Puzicha J (2002) Shape matching and object recognition using shape contexts. IEEE Trans Pattern Anal Mach Intell 24:509–522.

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.

Cichy RM, Pantazis D, Oliva A (2014) Resolving human object recognition in space and time. Nat Neurosci 17:455–462.

Cootes TF, Edwards GJ, Taylor CJ (2001) Active appearance models. IEEE Trans Pattern Anal Mach Intell 23:681–685.

Davenport JL, Potter MC (2004) Scene consistency in object and background perception. Psychol Sci 15:559–564.

Deng J, Dong W, Socher R, Li L-J, Kai Li, Li Fei-Fei (2009) ImageNet: A large-scale hierarchical image database. In: 2009 IEEE Conference on Computer Vision and Pattern Recognition, pp 248–255. IEEE.

Dodge S, Karam L (2019) Human and DNN Classification Performance on Images With Quality Distortions. ACM Trans Appl Percept 16:1–17.

Eckstein MP, Koehler K, Welbourne LE, Akbas E (2017) Humans, but Not Deep Neural Networks, Often Miss Giant Targets in Scenes. Curr Biol 27:2827-2832.e3.

Enns JT, Rensink R a (1991) Preattentive recovery of three-dimensional orientation from line drawings. Psychol Rev 98:335–351.

Enns JT, Rensink RA (1990) Sensitivity to Three-Dimensional Orientation in Visual Search. Psychol Sci 1:323–326.

Fleuret F, Li T, Dubout C, Wampler EK, Yantis S, Geman D (2011a) Comparing machines and humans on a visual categorization test. Proc Natl Acad Sci U S A 108:17621–17625.

Fleuret F, Li T, Dubout C, Wampler EK, Yantis S, Geman D (2011b) Comparing machines and humans on a visual categorization test. Proc Natl Acad Sci U S A 108:17621–17625.

Geirhos R, Rubisch P, Michaelis C, Bethge M, Wichmann FA, Brendel W (2018a) ImageNet-trained CNNs are biased towards texture; increasing shape bias improves accuracy and robustness. arXiv:1–22.

Geirhos R, Temme CRM, Rauber J, Schütt HH, Bethge M, Wichmann FA (2018b) Generalisation in humans and deep neural networks. arXiv.

Haber N, Mrowca D, Fei-Fei L, Yamins DLK (2018) Learning to Play with Intrinsically-Motivated Self-Aware Agents. Adv Neural Inf Process Syst 2018-Decem:8388–8399.

Jacob G, Arun SP (2019) How the forest interacts with the trees: Multiscale shape integration explains global and local processing. bioRxiv:777110.

Jarrett K, Kavukcuoglu K, Ranzato MA, LeCun Y (2009) What is the best multi-stage architecture for object recognition? In: 2009 IEEE 12th International Conference on Computer Vision, pp 2146–2153. IEEE.

878  Kar K, Kubilius J, Schmidt K, Issa EB, DiCarlo JJ (2019a) Evidence that recurrent
879      circuits are critical to the ventral stream's execution of core object recognition
880      behavior. Nat Neurosci 22:974–983.
881  Kar K, Kubilius J, Schmidt K, Issa EB, DiCarlo JJ (2019b) Evidence that recurrent
882      circuits are critical to the ventral stream's execution of core object recognition
883      behavior. Nat Neurosci 22:974–983.
884  Katti H, Arun SP (2019) Are you from North or South India? A hard face-classification
885      task reveals systematic representational differences between humans and
886      machines. J Vis 19:1.
887  Katti H, Peelen M V, Arun SP (2017) How do targets, nontargets, and scene context
888      influence real-world object detection? Atten Percept Psychophys 79:2021–2036.
889  Katti H, Peelen M V, Arun SP (2019) Machine vision benefits from human contextual
890      expectations. Sci Rep 9:2112.
891  Kietzmann TC, Spoerer CJ, Sörensen LKA, Cichy RM, Hauk O, Kriegeskorte N (2019)
892      Recurrence is required to capture the representational dynamics of the human
893      visual system. Proc Natl Acad Sci:201905544.
894  Kimchi R (1994) The role of wholistic/configural properties versus global properties in
895      visual form perception. Perception 23:489–504.
896  Mitchell S (1988) Tao te ching: A new English version. Harper Collins.
897  Mongia M, Kumar K, Erraqabi A, Bengio Y (2016) On Random Weights for Texture
898      Generation in One Layer Neural Networks. arXiv.
899  Munneke J, Brentari V, Peelen M V. (2013) The influence of scene context on object
900      recognition is independent of attentional focus. Front Psychol 4:552.
901  Navon D (1977) Forest before trees: The precedence of global features in visual
902      perception. Cogn Psychol 9:353–383.
903  Nguyen A, Yosinski J, Clune J (2015) Deep neural networks are easily fooled: High
904      confidence predictions for unrecognizable images. In: 2015 IEEE Conference on
905      Computer Vision and Pattern Recognition (CVPR), pp 427–436. IEEE.
906  Pardo-Vazquez JL, Castiñeiras-de Saa JR, Valente M, Damião I, Costa T, Vicente MI,
907      Mendonça AG, Mainen ZF, Renart A (2019) The mechanistic foundation of
908      Weber's law. Nat Neurosci 22:1493–1502.
909  Parkhi OM, Vedaldi A, Zisserman A (2015) Deep Face Recognition. In: Procedings of
910      the British Machine Vision Conference 2015, pp 41.1-41.12. British Machine
911      Vision Association.
912  Pramod RT, Arun SP (2014) Features in visual search combine linearly. J Vis 14:1–
913      20.
914  Pramod RT, Arun SP (2016a) Do computational models differ systematically from
915      human object perception? Proc IEEE Comput Soc Conf Comput Vis Pattern
916      Recognit 2016-Decem:1601–1609.
917  Pramod RT, Arun SP (2016b) Object attributes combine additively in visual search. J
918      Vis 16:8.
919  Rajalingham R, Issa EB, Bashivan P, Kar K, Schmidt K, DiCarlo JJ (2018) Large-scale,
920      high-resolution comparison of the core visual object recognition behavior of
921      humans, monkeys, and state-of-the-art deep artificial neural networks. J
922      Neurosci.
923  Ratan Murty NA, Arun SP (2017) Seeing a straight line on a curved surface:
924      decoupling of patterns from surfaces by single IT neurons. J Neurophysiol
925      117:104–116.
926  Rensink RA, Enns JT (1998) Early completion of occluded objects. Vision Res
927      38:2489–2505.

928 Rollenhagen JE, Olson CR (2000) Mirror-Image Confusion in Single Neurons of the
929        Macaque Inferotemporal Cortex. Science (80- ) 287:1506–1509.
930 Rosenfeld A, Zemel R, Tsotsos JK (2018) The Elephant in the Room. ArXiv.
931 Ruder S (2017) An Overview of Multi-Task Learning in Deep Neural Networks. arXiv.
932 Rust NC, Movshon JA (2005) In praise of artifice. Nat Neurosci 8:1647–1650.
933 Serre T (2019) Deep Learning: The Good, the Bad, and the Ugly. Annu Rev Vis Sci
934        5:399–426.
935 Simonyan K, Zisserman A (2014) Very Deep Convolutional Networks for Large-Scale
936        Image Recognition. arXiv:1–14.
937 Sinz FH, Pitkow X, Reimer J, Bethge M, Tolias AS (2019) Engineering a Less Artificial
938        Intelligence. Neuron 103:967–979.
939 Su J, Vargas DV, Sakurai K (2019) One Pixel Attack for Fooling Deep Neural
940        Networks. IEEE Trans Evol Comput 23:828–841.
941 Szegedy C, Zaremba W, Sutskever I (2013) Intriguing properties of neural networks.
942        arXiv Prepr arXiv:1–10.
943 Thompson P (1980) Margaret Thatcher: a new illusion. Perception 9:483–484.
944 Vedaldi A, Lenc K (2014) MatConvNet - Convolutional Neural Networks for MATLAB.
945        arXiv.
946 Vighneshvel T, Arun SP (2015) Coding of relative size in monkey inferotemporal
947        cortex. J Neurophysiol 113:2173–2179.
948 Vinje WE, Gallant JL, Vinje WE, Gallant JL (2000) Sparse coding and decorrelation in
949        primary visual cortex during natural vision. Science (80- ) 287:1273–1276.
950 Xu Y, Singh M (2002) Early computation of part structure: Evidence from visual search.
951        Percept Psychophys 64:1039–1054.
952 Yang J, Ren Z, Xu M, Chen X, Crandall D, Parikh D, Batra D (2019) Embodied Visual
953        Recognition. arXiv.
954 Zhivago KA, Arun SP (2016) Selective IT neurons are selective along many
955        dimensions. J Neurophysiol 115:1512–1520.
956 Zoccolan D, Cox DD, DiCarlo JJ (2005) Multiple Object Response Normalization in
957        Monkey Inferotemporal Cortex. J Neurosci 25:8150–8164.
958
959

## ACKNOWLEDGEMENTS

**SUPPLEMENTARY MATERIAL**
**for**

**"Do deep neural networks see the way we do?"**
**Georgin Jacob, RT Pramod, Harish Katti & SP Arun**


**CONTENTS**

# SECTION S1. RESULTS WITH OTHER FEED FORWARD NETWORKS

1  　　　The results in the main text were based on testing a specific feedforward network,
2  namely VGG-16. Here, we investigated other feedforward network architectures for the
3  presence of the same perceptual and neural phenomena. We did not test the recurrent
4  networks since unfolding recurrent networks over time make them equivalent to a deep
5  feedforward network (LeCun, Bengio and Hinton, 2015; Liang and Hu, 2015).
6
7  **Methods**
8  　　　We selected four popular pre-trained feedforward networks, all trained on the
9  ImageNet ILSVRC challenge data (Deng *et al.*, 2009; Russakovsky *et al.*, 2014). We
10 selected architectures that are shallower and deeper than VGG-16, to investigate whether
11 the depth of the network influences the emergence of the perceptual and neural
12 properties. All networks were implemented using MatConvNet framework in MATLAB,
13 and their performance is summarized in Table S1.
14
15 *Network 1: AlexNet.* This network won the ILSVRC 2012 challenge by a large margin
16 (Krizhevsky, Alex, Ilya Sutskever, 2012). The network consists of five convolutional layers
17 and three fully-connected layers. Drop-out technique is used fully connected layers to
18 reduce overfitting. The architecture of this network is shallower compared to VGG-16.
19
20 *Network 2: GoogLeNet.* This network follows the inception architecture which is well
21 known for better utilization of computing resources inside the network. This network won
22 the classification track of the ILSVRC 2014 challenge (Szegedy *et al.*, 2015).
23
24 *Network 3: ResNet-50.* ResNet-50 is a shallower variant of the ResNet-152 detailed
25 below.
26
27 *Network 4: ResNet-152.* The network uses a residual learning principle which make them
28 capable of training deeper networks without the problem of vanishing gradients (He *et al.*,
29 2016). The ResNet architecture won three tracks (classification, detection and
30 localization) of the ILSVRC 2015 challenge and two tracks (detection and segmentation)
31 of the COCO 2015 challenge.
32

| Name of network | Performance | |
|---|---|---|
| | Top-1 error (%) | Top-5 error (%) |
| AlexNet | 42.6 | 19.6 |
| VGG-16 | 28.5 | 9.9 |
| GoogLeNet | 34.2 | 12.9 |
| ResNet-50 | 24.6 | 7.7 |
| ResNet-152 | 23.0 | 6.7 |

33 **Table S1. Performance of deep networks on the ILSVRC 2012 validation dataset**
34 **(accuracy reported from MatConvNet website, accessed on 27th November 2019).**
35
36
37

## Results

*Experiment 1: Thatcher effect.* The Thatcher index for each network across layers is shown in Figure S1A. It can be seen that the Thatcher index is negative for all networks in their final layers except for GoogLeNet which showed a small positive level in the final layers. For the networks with higher classification performance (GoogLeNet, ResNet-50, ResNet-152), we observed an interesting pattern whereby the Thatcher index is positive in the intermediate layers. This is true even for the VGG-16 network (Figure 2B).

*Experiment 2: Mirror Confusion.* The mirror confusion index for each network is shown in Figure S1B. All networks exhibited an increasing mirror confusion index across layers, just as we observed for VGG-16 (Figure 2D).

*Experiment 3: Scene incongruence.* The classification accuracy for objects in congruent and incongruent scenes is shown for each network in Figure S1C. It can be seen that the deeper architectures show smaller incongruence effects.

*Experiment 4: Multiple object normalization.* The normalization slope for pairs and triplets for all networks is shown in Figure S1C. It can be seen that there is increased normalization in the later layers in all networks.

*Experiment 5: Selectivity across multiple dimensions.* The correlation between sparseness of units in each layer for textures and shapes is shown in Figure S1E. It can be seen that all networks show an increasing trend in later layers. We obtained qualitatively similar results for comparing sparseness on the reference shape set and morph lines (not shown for brevity).

*Experiment 6. Weber's law.* The Weber's law measure (difference in correlation for relative vs absolute length) for all networks is shown in Figure S1F. It can be seen that the Weber's law arises in the later layers for all the networks.

*Experiment 7. Relative size.* The relative size effect for each network across layers is shown in Figure S2A. It can be seen that the relative size effect is extremely weak and variable across networks, and never approaches the levels observed in the brain (relative size index = 0.39).

*Experiment 8: Decoupling patterns from surfaces.* The surface invariance index for each network across layers is shown in Figure S2B. It can be seen that the index is consistently negative for all networks, as observed for VGG-16.

*Experiment 9: 3D processing.* The 3D processing indices (for Condition 1 & 2) for each network across layers is shown in Figure S2C. It can be seen that the 3D processing indices are generally negative for all networks, and even if the index is positive, the levels are much smaller than observed in humans.

83    *Experiment 10: Occlusions.* The occlusion indices for each network across layers is
84    shown in Figure S2D. It can be seen that both indices are consistently negative across
85    layers, as observed for VGG-16.
86
87    *Experiment 11: Object parts.* The natural part advantage for Experiment 11B is shown for
88    each network across layers in Figure S2E. It can be seen that the natural part advantage
89    is highly variable across networks, with GoogLeNet showing levels comparable to
90    humans in the later layers.
91
92    *Experiment 12: Global shape advantage.* The global advantage index for each network
93    across layers is shown in Figure S2F. Across all networks, there is a slight global
94    advantage in the intermediate layers, which reverses into a local advantage in the later
95    layers. Thus, it appears that all the feedforward networks are using local features for
96    classification.
97

98
99
100 **Figure S1: Experiments 1-6 with other feed-forward networks.** Each column
101 represents a deep network (from left to right: AlexNet, GoogLeNet, ResNet-50 and Reset-
102 152) and each row represents an experiment (A-F: experiment 1-6). In each subplot an
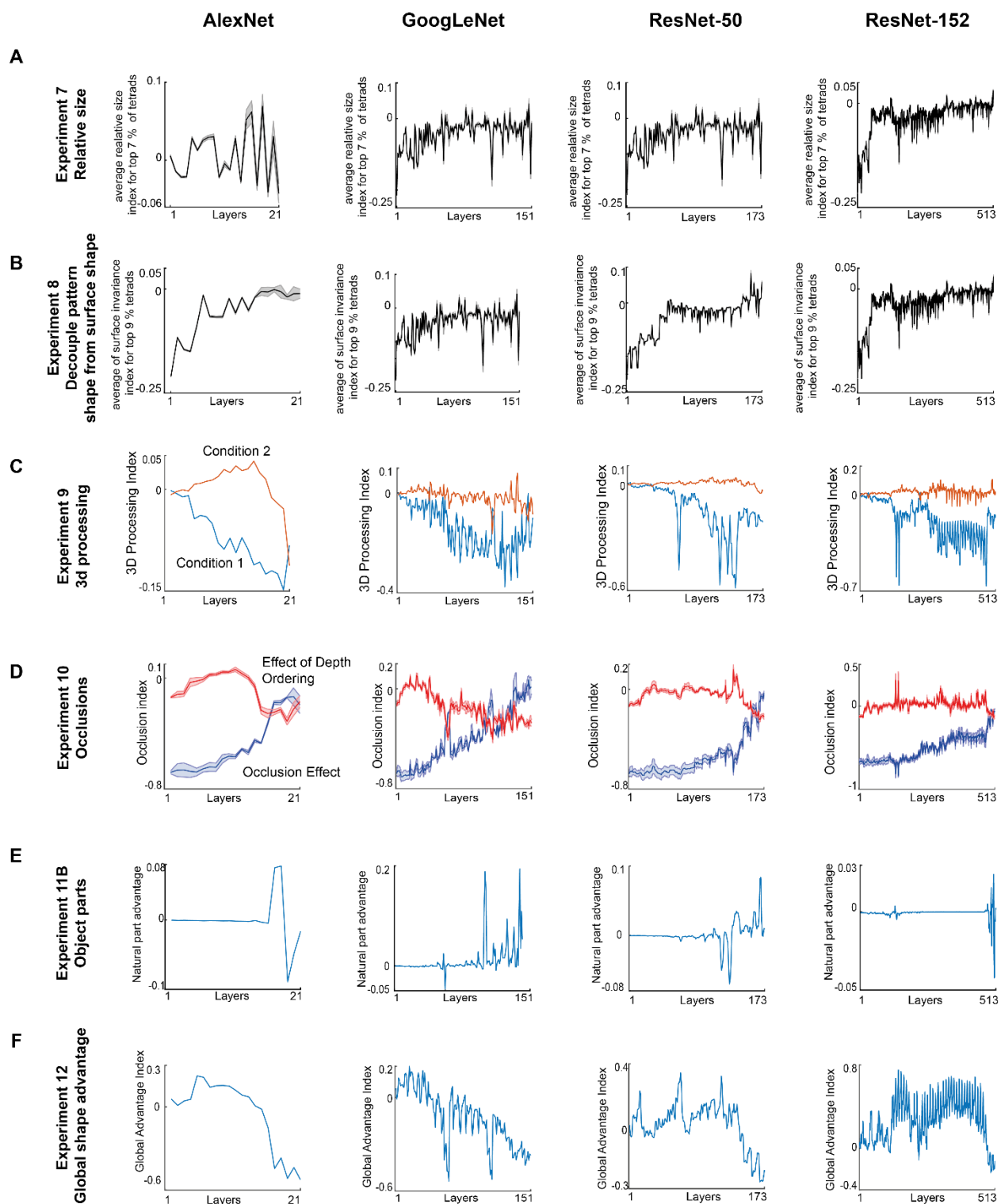103 experiment specific index or measure is plotted across layers.
104

**Figure S2: Experiments 7-12 for other feedforward networks.** All conventions as in Figure S1.

# SECTION S2. RESULTS WITH RANDOMLY INITIALIZED NETWORKS

108        We have shown that deep neural networks trained for object-recognition show
109 some perceptual phenomena but not others, but we wondered whether any of these
110 phenomena can be observed even in an untrained network, simply as a consequence of
111 the architecture. To address this issue, we repeated Experiments 1-12 on an untrained
112 VGG-16 architecture with randomly initialized weights.
113
114 **Methods**
115        We estimated the probability distribution functions of the weights in each layer of
116 the pre-trained VGG-16 network and picked the weights randomly from this estimated
117 distribution. This process was done layer-wise to obtain the randomly initialized network.
118 We then performed all Experiments on this network except for Experiment 3 (scene
119 incongruence) since this experiment requires object classification. We then performed all
120 data analyses exactly as before.
121
122 **Results**
123        The results for the random network and the pre-trained VGG-16 network are
124 shown together in Figure S3. In most cases (Experiments 1, 2, 6-11), we observed no
125 specific trend for the random network towards or away from human/neural levels, which
126 is what would be expected since the random network has no specific bias or training. We
127 observed systematic and interesting differences for the other experiments, which we
128 discuss in detail below.
129
130 *Experiment 4: Multiple Object normalization.* The divisive normalization slope for the
131 random network is shown in Figure S3C-D. Here, the random network showed perfect
132 divisive normalization, in that the net response to AB is exactly the average of the
133 responses to A & B separately. To investigate this puzzling observation further, we looked
134 at the unit activations in the random network. We found that the activations for any pair
135 of natural images was highly correlated (correlation of layer-37, mean ± sd: $r = 0.98 \pm$
136 $0.01$), suggesting that these units were not very selective for images. This means that
137 every image activates the network in the same way. As a result, the response to AB and
138 the response to A & B separately would be identical, giving rise to a perfect slope of 0.5
139 in the relationship between the response AB and the sum of responses A + B. Thus, the
140 divisive normalization observed in the random network is a trivial outcome of its lack of
141 image selectivity.
142
143 *Experiment 5: Selectivity across multiple dimensions.* Here too the random network
144 shows a high correlation between shape and texture selectivity (Figure S3E). We suspect
145 that this too is a consequence of the very low selectivity for images in the random network,
146 whereby some units have zero selectivity (and therefore respond equally to all images)
147 and others have weak selectivity (and therefore show slight differences in the response
148 across images).
149
150 *Experiment 12: Global shape advantage*. Here we observed an interesting pattern: the
151 random network showed a global advantage that increased across layers. This is likely

152    due to increased pooling in the higher layers, but interestingly the pre-trained VGG-16
153    network shows the opposite pattern. Thus, it appears that object classification training
154    abolishes the global advantage that is intrinsic to the network architecture. We speculate
155    that this local advantage might arise because of the demands of distinguishing between
156    highly similar categories present in the ImageNet dataset (e.g. there are 90 categories of
157    dogs among the total of 1000 categories in ImageNet). Testing this possibility will require
158    training the VGG-16 architecture on highly distinctive object classes.
159
160

**Figure S3. Experiments 1-12 for randomly initialized networks.** Results for Experiments 1-12 (except for #3) are shown. In each panel, the corresponding experiment-specific index is plotted for the randomly initialized VGG-16 (*red*) with the pre-trained VGG-16 (*blue*). All other conventions are as in the main text.

167

A

**Experiment 1: Thatcher effect**

Faces are removed due to bioRxiv policy

B

**Experiment-2 : Mirror confusion**



168

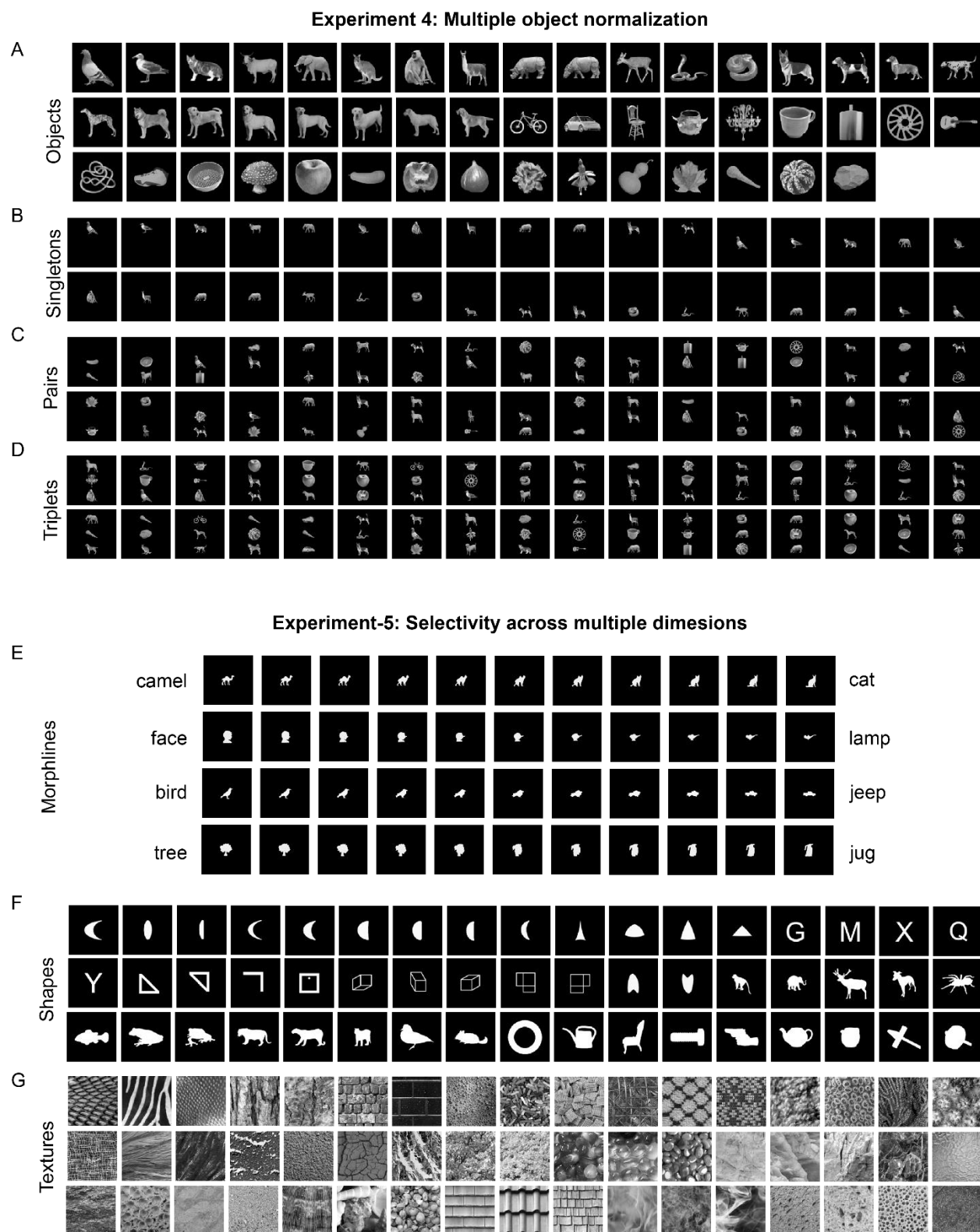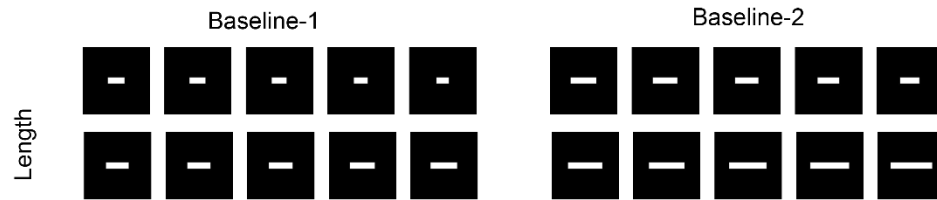**Figure S4. Stimulus set used for Experiments 1 & 2.**

**Experiment 3 : Scene incogruence**

Davenport & Potter, 2004



Munneke et. al., 2013



169
170 **Figure S5. Stimulus set used for Experiment 3**. Each pair of images depicts an object
171 against a congruent and incongruent background. Stimulus set reproduced with consent.
172
173

**Experiment 4: Multiple object normalization**

A

Objects



B

Singletons



C

Pairs



D

Triplets



**Experiment-5: Selectivity across multiple dimesions**

E

Morphlines

| camel |  | cat |
| face | | lamp |
| bird | | jeep |
| tree | | jug |

F

Shapes



G

Textures



174
175 **Figure S6. Stimulus set used for Experiments 4 & 5**
176 Plots A-D shows the images used for experiment 4 and plots E-G shows the images
177 used in experiment-5.

178     A) 49 natural images used for experiment-4
179     B) Singletons are formed by placing the 49 objects either at top, middle or bottom.
180        34 selected singletons are shown from 147 singletons used in this experiment.
181     C) 34 selected pairs are shown from 200 pairs used in experiment-4.
182     D) 34 selected triplets are shown from 200 triplets used in experiment-4.
183     E) Shows all four morphlines used in experiment-5.
184     F) Shown 51 selected shapes of the total set of 120 shapes.
185     G) Shown 51 selected textures out of the total set of 120 textures.
186

**A**

**Experiment-6 : Weber's law**

Baseline-1                    Baseline-2

Length

**B**

**Experiment-7 : Relative size**

**C**

**Experiment-8 : Decouple pattern shape from surface shape**
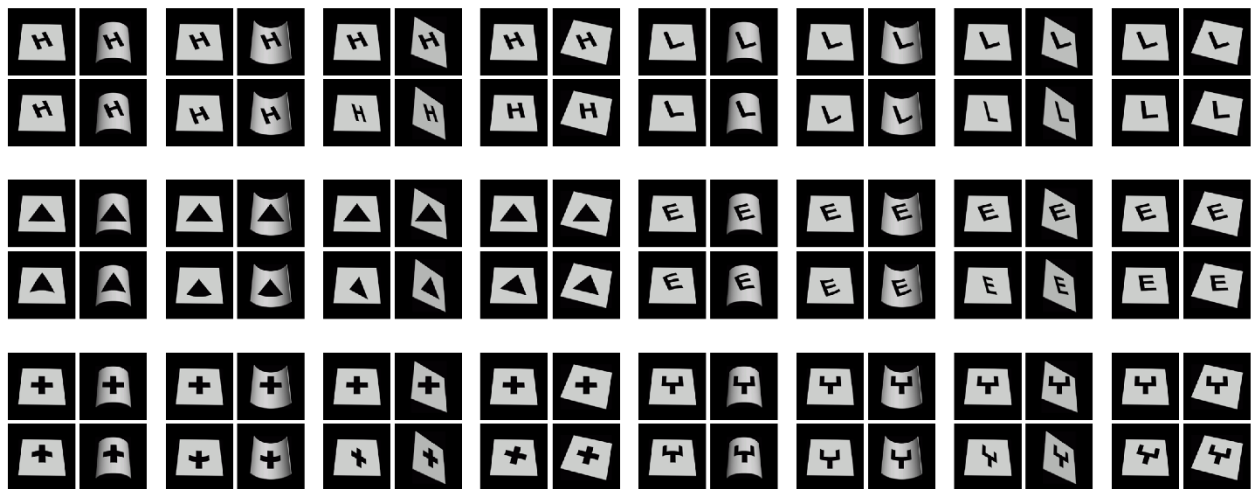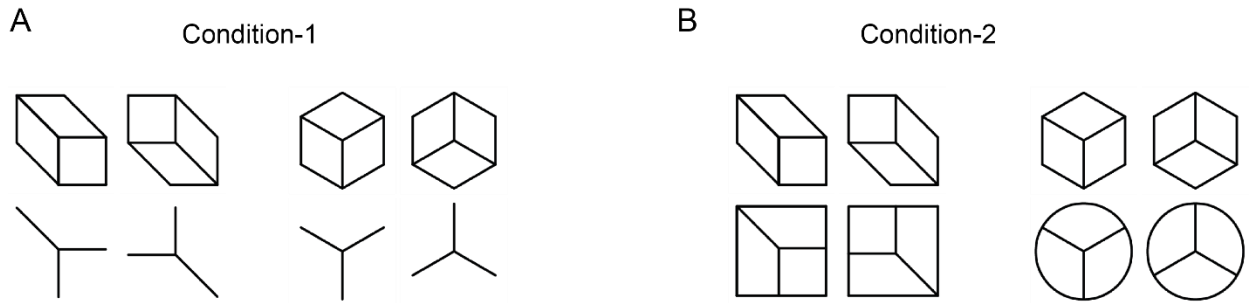
187

**Figure S7. Stimulus set used for Experiments 6-8**

188    A) Shows the images used in Weber's law experiment. There are two baselines and
189        each column is an image pairs which has equal length difference from the
190        baseline.
191    B) Shows 24 tetrads used in Relative Size experiment.
192    C) Shows 24 tetrad used in surface invariance experiments. Tetrads are made by
193        transforming eight shapes onto five different surfaces.

**Experiment 9: 3d processing**

A     Condition-1            B     Condition-2



**Experiment 10: Occlusions**

C                         D

Occlusion                      Depth Ordering
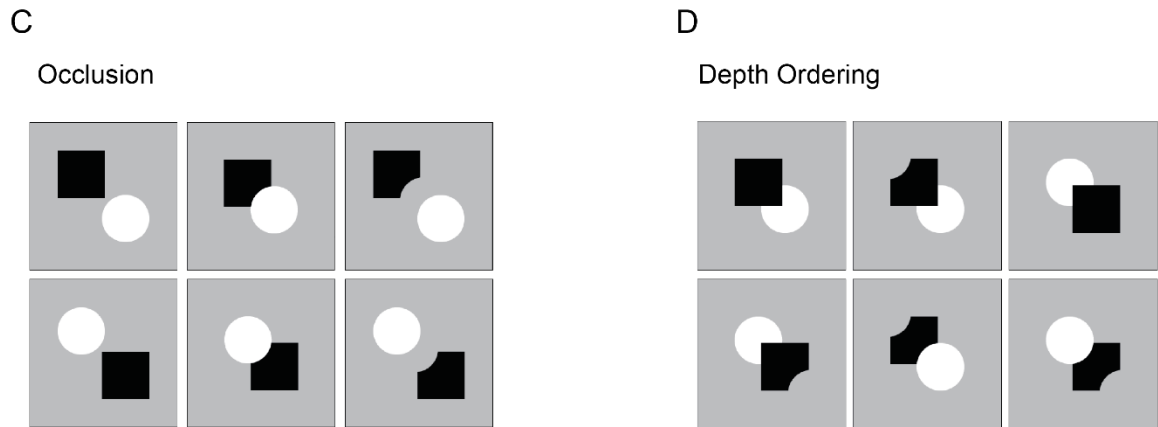


194

**Figure S8. Stimulus set used for Experiments 9-10**

    A) Two sets of images used to compare the 3D perception in the CNNs. Images in top row have 3D effect whereas the images in the bottom row have an equivalent same feature difference without a perceived 3D difference.

    B) Two sets of images used to compare the 3D perception in the CNNs. Images in top row have 3D effect whereas the images in the bottom row has the same feature difference and an additional common outer shape but not the 3D perception.

    C) Each row shows a set of images used for testing basic occlusion effect.

    D) Each row shows a set of images used for testing depth ordering effect.
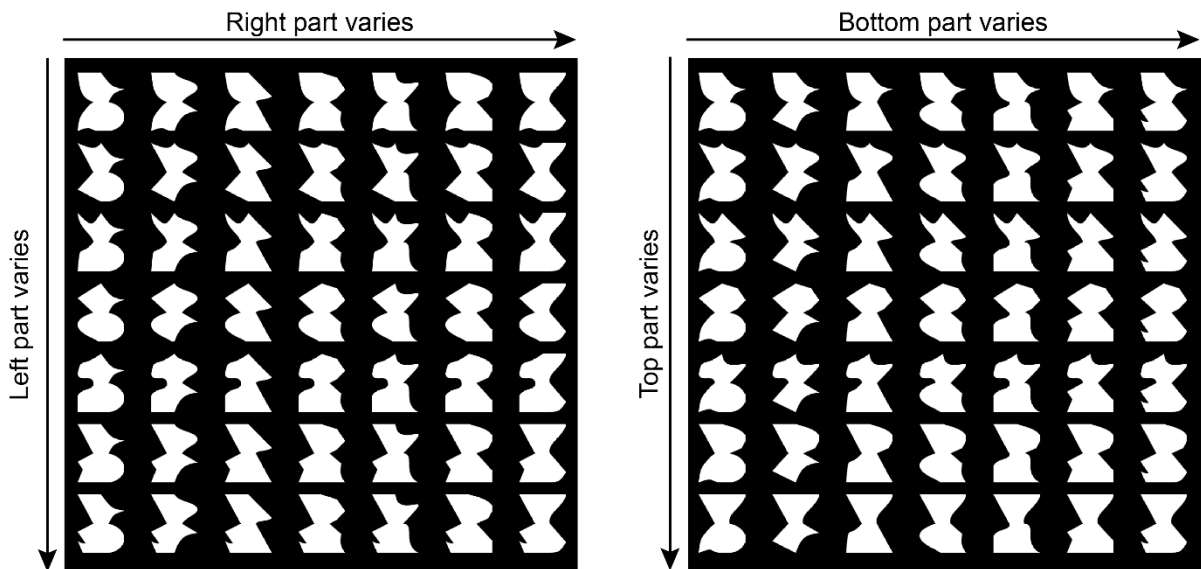
195

A

**Experiment 11: Object parts**



Xu & Singh, 2002

B

**Unnatural Cut**                    **Natural Cut**

Right part varies                    Bottom part varies



Left part varies                    Top part varies

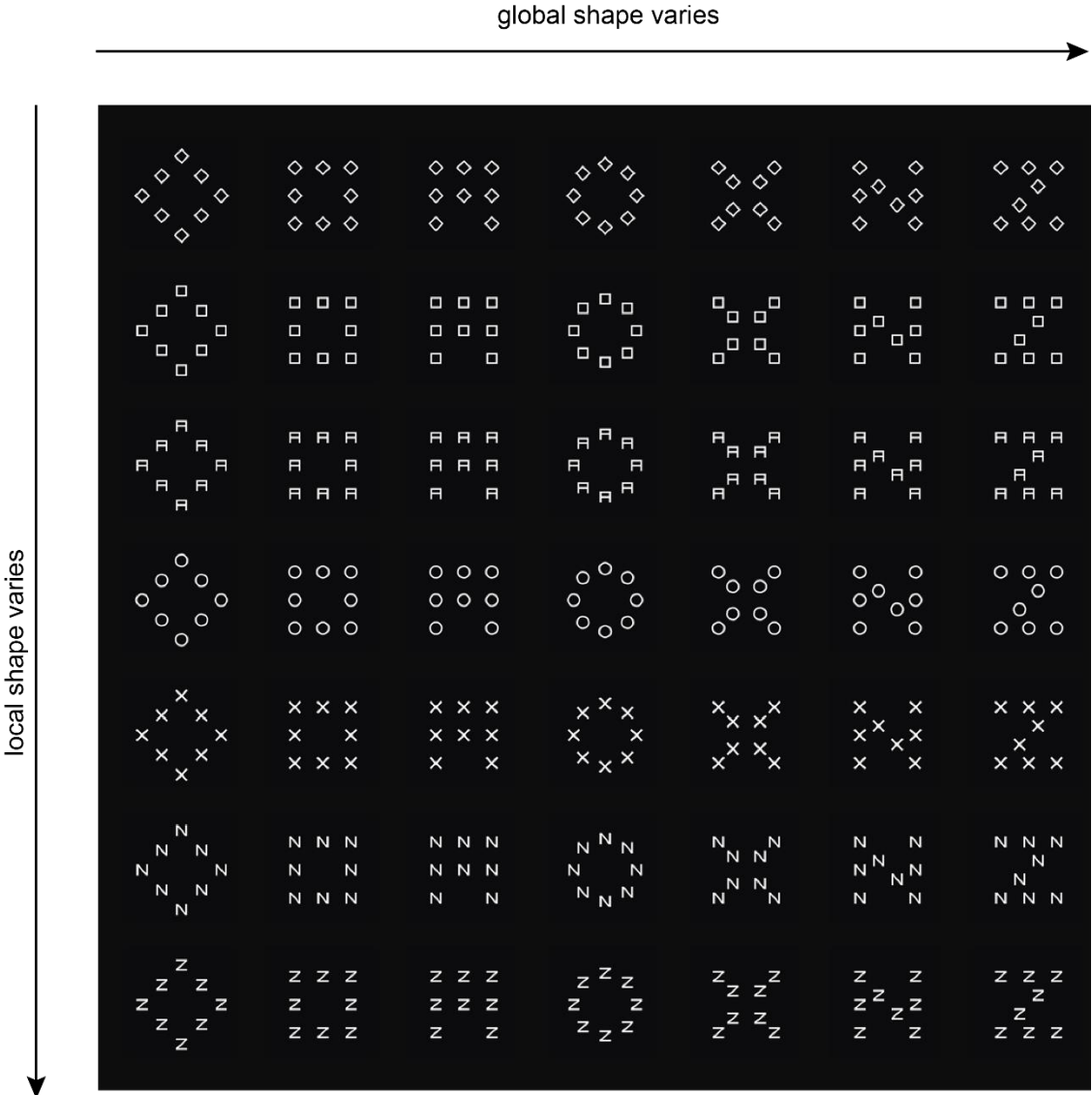Objects common to both sets



Pramod & Arun, 2016

196
197 **Figure S9. Stimulus set used for Experiments 11**
    A) Shows the images used to check part processing in experiment-11
    B) Shows the images used to check part advantage in experiment-11. The seven
       shapes in the diagonal position of both Unnatural and natural set are the same.

**Experiment 12: Global shape advantage**

global shape varies



Jacob & Arun, 2019

198    **Figure S10. Stimulus set used for Experiments 12.** A total of 49 images used to check
199    global advantage. These images are formed by all combinations of seven shapes are
200    global and local scales.
201

# REFERENCES

Deng, J. *et al.* (2009) 'ImageNet: A large-scale hierarchical image database', in *2009 IEEE Conference on Computer Vision and Pattern Recognition.* IEEE, pp. 248–255. doi: 10.1109/CVPR.2009.5206848.

He, K. *et al.* (2016) 'Deep Residual Learning for Image Recognition', in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, pp. 770–778. doi: 10.1109/CVPR.2016.90.

Krizhevsky, Alex, Ilya Sutskever, and G. E. H. (2012) 'ImageNet Classification with Deep Convolutional Neural Networks', *Advances in Neural Information Processing Systems 25 (NIPS 2012)*, pp. 1–9.

LeCun, Y., Bengio, Y. and Hinton, G. (2015) 'Deep learning.', *Nature*, 521(7553), pp. 436–44. doi: 10.1038/nature14539.

Liang, M. and Hu, X. (2015) 'Recurrent convolutional neural network for object recognition', *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*. IEEE, 07-12-June(Figure 1), pp. 3367–3375. doi: 10.1109/CVPR.2015.7298958.

Russakovsky, O. *et al.* (2014) 'ImageNet Large Scale Visual Recognition Challenge', *International Journal of Computer Vision*. Springer US, 115(3), pp. 211–252. doi: 10.1007/s11263-015-0816-y.

Szegedy, C. *et al.* (2015) 'Going deeper with convolutions', in *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, pp. 1–9. doi: 10.1109/CVPR.2015.7298594.