
MindSet: Vision. A toolbox for testing DNNs on key psychological experiments

Valerio Biscione^{1*} Dong Yin¹ Gaurav Malhotra² Marin Dujmovic³
Milton L. Montero⁴ Guillermo Puebla⁵ Federico Adolfi^{1,6}
Rachel F. Heaton^{7,8} John E. Hummel⁸ Benjamin D. Evans⁹ Karim Habashy¹
Jeffrey S. Bowers¹

¹School of Psychological Science, University of Bristol

²University at Albany, State University of New York

³School of Physiology, Pharmacology and Neuroscience, University of Bristol

⁴IT University of Copenhagen

⁵Instituto de Alta Investigación, Universidad de Tarapacá, Chile

⁶Ernst Strüngmann Institute for Neuroscience, Max-Planck Society, Germany

⁷Siebel Center for Design, University of Illinois Urbana-Champaign

⁸Department of Psychology, University of Illinois Urbana-Champaign

⁹School of Engineering and Informatics, University of Sussex

Abstract

1 Multiple benchmarks have been developed to assess the alignment between deep
2 neural networks (DNNs) and human vision. In almost all cases these benchmarks
3 are observational in the sense they are composed of behavioural and brain re-
4 sponses to naturalistic images that have not been manipulated to test hypotheses
5 regarding how DNNs or humans perceive and identify objects. Here we intro-
6 duce the toolbox *MindSet: Vision*, consisting of a collection of image datasets
7 and related scripts designed to test DNNs on 30 psychological findings. In all
8 experimental conditions, the stimuli are systematically manipulated to test spe-
9 cific hypotheses regarding human visual perception and object recognition. In
10 addition to providing pre-generated datasets of images, we provide code to regen-
11 erate these datasets, offering many configurable parameters which greatly extend
12 the dataset versatility for different research contexts, and code to facilitate the
13 testing of DNNs on these image datasets using three different methods (similar-
14 ity judgments, out-of-distribution classification, and decoder method), accessible
15 at <https://github.com/ValerioB88/mindset-vision>. We test ResNet-152
16 on each of these methods as an example of how the toolbox can be used.

17 1 Introduction

18 Deep neural networks (DNNs) provide the best solution for visual identification of objects short of
19 biological vision, and many researchers claim that DNNs are the best current models of human vision
20 and object recognition [70, 79, 120, 38, 53]. Key evidence in support of this claim comes from the
21 finding that DNNs perform the best on various behavioural and brain benchmarks. In the case of
22 behavioural benchmarks, models are assessed on how well they account for human (or macaque)

*corresponding author: valerio.biscione@gmail.com

23 errors in classifying a large set of objects [91, 108, 60], or how well they predict human similarity
24 judgements [84, 12]. In the case of brain benchmarks, models are assessed with regard to how well
25 they predict brain recordings (e.g., single-cell responses or fMRI data) in response to a set of objects
26 [96, 97, 60, 4]. The general assumption is that the better a model does at predicting the data, the more
27 similar the model is to biological vision. For instance, the Brain-Score benchmark is described as “a
28 composite of multiple neural and behavioural benchmarks that score any [artificial neural network]
29 on how similar it is to the brain’s mechanisms for core object recognition” [96]

30 A common feature of most benchmark studies is that they treat the to-be-predicted data as observa-
31 tional. That is, there is rarely an attempt to predict the impact of experimental manipulation designed
32 to test specific hypotheses about how human or machine vision works. Rather, observers perform
33 a single task over a set of images that satisfies some general criterion, such as objects presented in
34 isolation [69], in naturalistic contexts [84, 12, 4, 60], or on a range of arbitrary backgrounds [30, 96].
35 This approach is problematic because it is possible to make good predictions on these datasets even
36 when models identify objects in a qualitatively different way from monkeys or humans [41]. For
37 example, if the images contain multiple diagnostic cues for object classification (e.g., shape and
38 texture both predict object category), then good predictions might be driven by different features than
39 those that drive human object recognition – that is, predictions might be driven by confounds. For
40 example, a DNN that classifies objects by texture might still be able to predict brain activations in a
41 visual system that classifies objects by shape.

42 The standard way to rule out confounds in order to determine causal relations (e.g. inferring that
43 DNNs learn brain-like representations) is to carry out experiments designed to rule out confounds
44 as the basis of making good predictions. In fact, there is a large literature in psychology describing
45 experiments designed to test specific hypotheses about how human vision works, but surprisingly,
46 this literature is often ignored when modellers compare DNNs to biological vision. Bowers et
47 al.[28] reviewed a wide range of psychological phenomena that current DNNs either fail to capture
48 or that have yet to be considered. Furthermore, when researchers do consider the psychological
49 literature when making claims regarding DNN-human similarities, the models are rarely subject
50 to the kind of “severe” tests that are required to make any strong conclusions, that is tests that are
51 likely to challenge claims in case they are false. Instead, strong conclusions are often drawn based on
52 superficial similarities [26].

53 There are at least four (related) reasons for this. First, many researchers in computer science
54 and computational neuroscience may be unfamiliar with the rich set of experiments carried out
55 in psychology that manipulate independent variables to better understand human vision, memory,
56 language, etc. Those who are aware of these studies might find it challenging to engage with them,
57 as psychological datasets are not readily available in formats that the community is accustomed to
58 working with. Second, it is not always obvious how one should test a model against psychological
59 data. Hence, it may be easier to focus on improving performance on the current benchmarks, and this
60 may have discouraged researchers from exploring data from psychology. A third potential reason is
61 an overall skepticism towards psychological results, a sentiment that may reflect the well-documented
62 replication crisis in psychology [7]. Forth, there is a strong bias to look for DNN-human similarities
63 and downplay the differences [26], and severely testing on psychological data might not result in
64 similarities. However, characterizing these failures provides key insights into the ways DNNs need to
65 be improved when modelling biological vision.

66 Here we present *MindSet: Vision*, a toolbox aimed at facilitating testing DNNs on visual psychological
67 phenomena by addressing all the problems presented above: our main contribution is to provide a
68 large, easily accessible, parameterized, set of 30 image datasets (and related scripts to re-generate
69 and modify them) accounting for a wide array of well-replicated visual experiment and phenomena
70 reported in psychology. Our stimuli cover aspects of low and mid-level vision (including Gestalt
71 phenomena), visual illusions, and object recognition tasks. We provide a high-level descriptions of
72 the visual phenomena in the main text (Section 2) and more detailed descriptions in the Appendix
73 (A). To facilitate experimentation across a variety of scenario, each dataset can be easily regenerated
74 across different configurations (image size, background colour, stroke colour, number of samples,

75 etc.). To address the difficulty in testing DNNs on these stimuli, we provide scripts for using one (or
76 more) of three methods: Similarity Judgment Analysis, Decoder Approach, and Out-of-Distribution
77 classification (Section 3). We provide examples illustrating how to use these scripts with a classic
78 feed-forward CNN (ResNet-152), and an extensively documented code (Section 4).

79 With *MindSet: Vision*, we aim to bridge the gap between computational modeling and psychological
80 research, bringing experimental studies that manipulate independent variables to the forefront of
81 developing and evaluating of DNN models of human vision. We also hope this initiative will drive
82 further interest in other areas of human psychology, such as memory, language, and speech perception
83 when attempting to understand and replicate human-like intelligence in machines.

84 1.1 Related Work

85 Several recent studies share some similarities with our project: [119] introduced a new dataset
86 containing five types of visual illusions falling into two categories: color constancy and geometrical
87 illusions. The authors formulated four tasks specifically designed to examine the performance
88 of Visual Language Models, finding low alignment with human responses. [75] developed the
89 Good Gestalt datasets, consisting of six types of datasets covering several types of Gestalt grouping
90 principles, including Closure, Continuity, and Proximity, aimed at testing a Latent Noise Segmentation
91 Network. Similarly, [50] developed the model-vs-human benchmark that compares ANN-human
92 classification errors on various ‘out-of-distribution’ datasets composed of naturalistic images that
93 were modified in various ways, including low-level feature manipulations of contrast and spatial
94 frequency, as well as higher-level manipulations, such as generating silhouettes and sketches of
95 images. Evans et al. [44] used a dataset of silhouettes, line-drawings, and contours, to investigate
96 robustness to these stimuli in DNNs pretrained on CIFAR-10, and Baker et al. [10] employed a
97 dataset of line drawings and silhouettes from ImageNet classes to investigate model robustness to
98 local versus global features. In comparison to these works, we present a toolbox to test DNNs
99 on visual psychological effects, investigating not only a much richer set of visual phenomena, but
100 providing the code base to regenerate images in batches, changing the parameters, and testing each
101 one of them on a variety of methods.

102 2 Datasets

103 We have included datasets from experiments that characterize a wide range of visual phenomena,
104 ranging from low- to high-level vision. We grouped the datasets (indicated in **bold**) into 3 broad
105 categories (see following Sections) as illustrated in Figure 1. Each dataset comprised multiple
106 sub-conditions designed to test DNN-human similarities, and in some cases, image datasets used to
107 train decoders, as described in Section 3.3.

108 While most of the stimuli are created by us, in a few instances we incorporate stimuli from external
109 sources (when needed, permission was obtained from the authors). In all cases, the stimuli have
110 been integrated into a versatile framework which offers significant flexibility in adjusting parameters
111 such as image size, background, stroke colour, and more, to allow their application to a variety
112 of models and methodologies. Given the extensive range of datasets provided, we only offer a
113 brief summary for each in the article, and provide more details in Appendix A, including details
114 about the suggested way to test each dataset, and the expected result for model-human perceptual
115 alignment. All resources are open-source and freely available under the MIT license at <https://github.com/ValerioB88/mindset-vision>.
116

117 2.1 Low and Mid-Level Vision

118 A fundamental low-level vision phenomena is captured by **Weber’s Law** [112], which states that
119 the minimum physical change of a stimulus on some dimension (e.g., its size) that is noticeable to
120 an observer is a constant ratio of the original stimulus value on this dimension. For example, it is

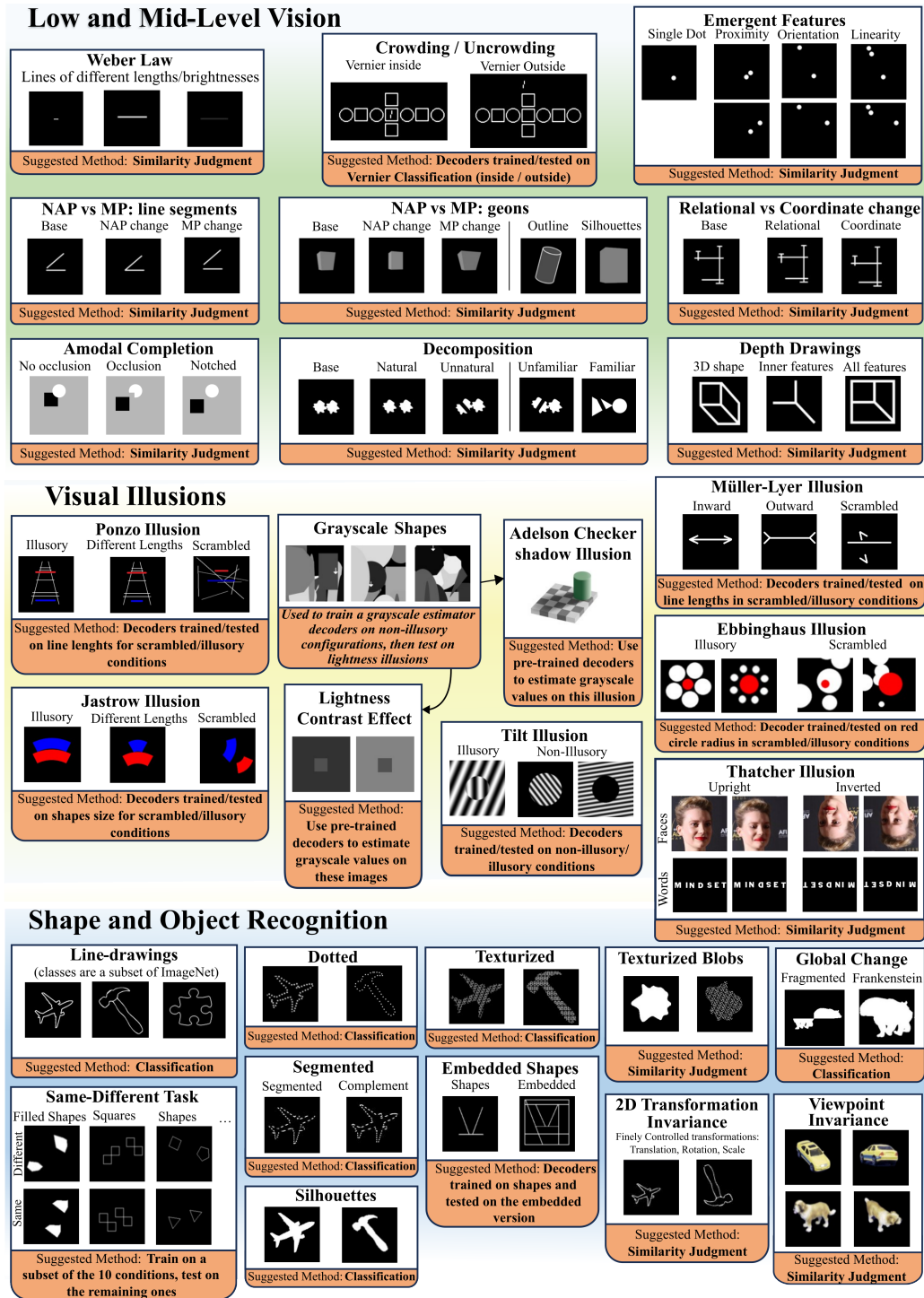


Figure 1: Comprehensive overview of the ‘MindSet: Vision’ datasets, arranged in three main categories. Each panel represents a distinct dataset, which is further divided into conditions. The images provide examples from these conditions, generated with default parameters.

121 equally easy to distinguish between line lengths of 1 and 2 cms and between 2 and 4 cms. We created
 122 a dataset that can be used to assess this relation for both line length and stimulus intensity.

123 Human perception is also sensitive to various **Emergent Features** in which simple image features
124 interact to generate “Gestalts” [85]. The dataset is comprised of a set of dots arranged in such a way
125 as to induce the emergent feature of proximity, orientation, and linearity [86, 22]. Another Gestalt
126 effect is manifest in the **Crowding/Uncrowding** phenomenon. In crowding, the ability to identify
127 an object is compromised by the presence of nearby objects or visual patterns, but in uncrowding,
128 object identification is improved when additional objects or visual patterns are added to the scene and
129 grouped such that they are segregated from the target. We adapt the dataset from [40] so that many
130 crowding conditions with several different shapes can be investigated.

131 Human perception is highly sensitive to non-accidental image features, that is features that are largely
132 invariant to changes in viewpoint when projected on to the retina [16], as opposite to accidental
133 features (e.g., degree of curvature) in which the projected image varies with viewpoint. Human vision
134 is known to be more sensitive to changes in images that alter non-accidental compared to accidental
135 properties [5, 6]. We use two datasets to examine model sensitivity to these features: one with **3D**
136 **geons** [5] and another with **2D line segments** (based on [71]). Similarly, we present a dataset to
137 compare **Relational Changes** with **Coordinate changes** between object parts. DNNs are commonly
138 insensitive to relational change, even after being explicitly trained on these relations [77], whereas
139 human perception is highly sensitive to relational changes [64].

140 To identify partly occluded objects, the human visual system groups contours and surfaces through
141 an amodal completion process [81]. The **Amodal Completion** dataset (based on [93]) enables the
142 investigation of these processes using images with shapes that are either occluded, unoccluded, or
143 “notched”. These latter shapes are unoccluded but notched in such a way as to maintain a high degree
144 of feature similarity with their occluded counterparts.

145 With the **Decomposition** dataset we provide a mean to test the extent to which DNNs group object
146 parts in a human-like fashion. We designed familiar and unfamiliar objects composed of two conjoined
147 parts that undergo what humans would perceive as a “natural” or “unnatural” break (inspired by
148 [65]). In the same work, [65] showed that VGG-16 trained on ImageNet did not possess human-like
149 sensitivity to images that could be interpreted as 3D-shapes by using a set of stimuli based on [42].
150 Accordingly, we reconstructed this **Depth Drawings** dataset.

151 2.2 Visual Illusions

152 Visual illusions are not mere curiosities, but often arise from adaptive perceptual processes [56].
153 Detailed computational models of multiple illusions have been advanced that provide theoretical
154 insights into the mechanisms that underlie them (e.g. [58]). We provide datasets exploring illusions
155 related to size perception, orientation, and lightness contrast.

156 Several illusions relate to size perception. In the **Müller-Lyer illusion**[35], arrow-like segments
157 at the ends of equal-length lines impact our perception of length. In the **Ponzo illusion**[88], two
158 equal-length horizontal lines cross a pair of converging lines. In this configuration, the top line looks
159 longer, an illusion often explained as related to the process of inferring depth. In the **Ebbinghaus**
160 **illusion**[2], the size of a circle is perceived differently depending on the size of surrounding circles.
161 Similarly, in the **Jastrow illusion**[66, 35], a specific arrangement of identical objects affects our
162 perception of their relative size. For all these illusions, we provide both an illusory condition and a
163 condition in which all elements of the original illusion are “scrambled up”, so that a decoder can be
164 trained to predict a specific feature (e.g. the size of the centre circle in the Ebbinghaus illusion) and
165 subsequently tested on illusory configurations.

166 In the **Tilt illusion**[52], the orientation of a central grating is perceived as being repulsed from or
167 attracted to the orientation of a surrounding grating. In this case, we provide conditions with either
168 central or background gratings (configurations which do not support the illusion in humans and could
169 be used for training a decoder), and a condition with both (eliciting the illusion in humans) for testing
170 the model. Another orientation illusion is the **Thatcher Effect** [105]. This is a phenomenon where
171 local changes in facial features (like inverted eyes or mouth) are less noticeable when the entire face is
172 upside down, highlighting our sensitivity to orientation in face perception. An interesting unresolved

173 issue is the extent to which this inversion effect is specific to faces [25, 115]. Together with a dataset
174 of faces and their Thatcherized version, we also include a dataset of **Thatcherized Words**, that is a
175 dataset of images containing words in which one or more letters are rotated by 180 degrees [115].

176 The **lightness contrast effect**[109] and the **Adelson Checker shadow**[1] illusions reveal how our
177 visual system perceives color and lightness based on context. We provide a **Grayscale Shapes**
178 **dataset** to train a decoder to output estimates of lightness at a given location of an image (indicated
179 by a small white arrow). After training, the network is presented with test images that induce illusions
180 and help assess whether DNNs show similar effects by pointing the arrow at the relevant parts of the
181 images (see Section C.2.6 for a detailed description of this approach).

182 It is important to note that there is no accepted account for some of the illusions described above.
183 However, even when we have no good understanding of the functional role or the mechanism that
184 drives an illusion, a DNN model of human vision should show similar effect. Indeed, understanding
185 the conditions under which DNNs show an illusion may advance our understanding of why the
186 phenomenon is observed in humans. There are now several articles exploring such illusions in various
187 types of DNNs trained in different ways, with some highlighting similarities (e.g., [14, 103, 111])
188 others reporting mixed or discrepant results (e.g., [55, 110, 119]); for a review of the relevant findings,
189 see [67].

190 2.3 Shape and Object Recognition

191 DNN object recognition is much more sensitive than human vision to distributional shifts from the
192 training set. For instance, humans can easily identify line drawings the first time they are exposed
193 to them [62], whereas DNNs perform poorly under these conditions [45] and need to be trained on
194 line drawings in order to recognize them at human levels [99]. We have included the **line drawing**
195 and **silhouettes** datasets (from [10, 8]) and also manipulated them in various ways to construct
196 additional datasets. The line drawings were converted into dotted contours (**Dotted line drawings**),
197 line segments (**Segments line drawings**) [18], or “texturized” (**Texturized line drawings**). The
198 texturized images are composed of oriented lines/characters applied to either/both the background
199 or/and the inside area of the line drawing. In all these cases, the resulting images are easily identifiable
200 by human observers due to various Gestalt rules that organize the image features into boundaries.
201 We also apply the same texturization technique outlined above on unfamiliar “blob”-like shapes
202 (**Texturized Unfamiliar dataset**). Human observers have no difficulty matching a novel “blob”
203 object to its texturized counterpart. In addition, we provide a dataset of fragmented images based on
204 [9] in which the global features of silhouettes or line drawing are modified by reflecting the top part of
205 an object along its vertical axis, leaving the local features mostly unchanged (**Global Modifications**
206 dataset). Human performance on these stimuli is greatly reduced but typically DNN performance
207 is largely unchanged, suggesting that human vision is more sensitive to global object structure and
208 DNN vision is more sensitive to local features.

209 The **Embedded Shapes Dataset** (inspired by [36]) provides another condition that greatly impacts
210 on human perception, by embedding geometric shapes within complex arrays of lines in ways that
211 camouflage the original shape. We include both the original images from [36] and a procedurally
212 generated dataset in which random polygons are embedded into a configuration that makes recognition
213 challenging for humans.

214 The human visual system supports object recognition following a wide variety of transformations
215 [104, 24]. Importantly, this extends to cases in which an object has only been viewed at one pose.
216 Previous works suggest a complex link between DNN pretraining and their object recognition
217 capabilities under object transformations [20, 21]. To test whether DNNs share these capacities, we
218 provide a dataset in which translations, plane rotations, and scale changes (**2D Transformations**)
219 are applied to line drawings. To test for **Viewpoint Invariance** (e.g. the ability to recognize an
220 object from a new viewpoint after a rotation in depth) we adapt the ETH-80 dataset, [32] allowing for
221 controlled variation in azimuth and inclination.

222 Finally, we provide a dataset to test whether DNNs possess the ability to solve a basic form of visual
223 reasoning task, namely, the **Same/Different** task. Drawing from [89], our dataset comprises images
224 composed of pairs of objects, which may be identical or different. These images are organized into
225 ten conditions that vary in their visual form, such as ‘filled polygons’, ‘open squares’, and ‘colored
226 shapes’. While humans effortlessly accomplish this task across all conditions without training, DNNs
227 often struggle when the training and test images come from different conditions.

228 **3 Testing methods.**

229 Each dataset is designed to align with at least one of three methods of testing, but other approaches
230 can be used as well. We discuss further possibilities in Appendix B.

231 **3.1 Out-of-Distribution Classification**

232 In this approach, a DNN pretrained on one dataset is tested on a new dataset composed of out-of-
233 distribution images taken from the trained classes (e.g., a DNN pre-trained on ImageNet is tested on
234 line drawings taken from the same categories). This approach is well suited for most of the Shape
235 and Object Recognition datasets that use images from ImageNet categories modified in such a way
236 that human observers have no trouble recognizing them, even without training. We provide scripts to
237 test a wide variety of vision models.

238 **3.2 Similarity Judgment Analysis**

239 This method involves assessing the pairwise similarity of activation patterns in DNNs (using a
240 Cosine Similarity or an Euclidean Distance metric) evoked by pairs of images and comparing these
241 similarities to human performance. This method has been used to assess how well DNNs capture
242 human similarity judgments [84] and response times to identify target stimuli from foils [22]. It is
243 often useful to carry out these analyses across multiple layers of DNNs given that some psychological
244 phenomena are known to manifest at earlier or later stages of visual processing. A DNN mimicking
245 human perception should show relevant similarity effects at the relevant layers. One key advantage of
246 this approach is that it can be applied to novel images that cannot be classified by a DNN.

247 To illustrate, we applied this method to the Texturized Unfamiliar dataset (Figure 2). The human
248 visual system groups elements in a scene by texture [13] and classify objects by their shapes [16].
249 Accordingly, texturized versions of the same shape should be judged as more similar than texturized
250 versions of different shapes. To explore if DNNs exhibit similar behaviour, we input pairs of images
251 into a ImageNet pre-trained ResNet-152 and, for each pair, we computed the Euclidean Distance
252 between their internal activations at every processing level. A human-like response is indicated by a
253 smaller distance for pairs of the same compared to different shapes. ResNet-152 exhibited a weak
254 manifestation of this pattern in the early layers, a reduced effect in the later layers, and no effect in
255 the output layer. By contrast, the human visual system supports similarity judgements on the basis of
256 shape-based representations that are computed following the early stages of visual processing.

257 **3.3 Decoder Method**

258 In this method a small, often single-layer, “decoder” network is attached to a layer of a frozen DNN
259 and trained on a task designed to reveal how the DNN encodes a specific type of information. For
260 instance, a frozen DNN might be presented with a set of images that contain a target object varying
261 in size, colour, and orientation, and a decoder is trained to output the value of one or more of these
262 properties at a given layer. We provide scripts for both classification and regression training, and
263 scripts to train and test a series of five decoders at varying levels of a ResNet-152 model. Although
264 these scripts are tailored to ResNet-152, they can easily be used as a template to streamline the
265 adaptation of this technique for different networks.

266 To illustrate, consider the Ebbinghaus Illusion. The Ebbinghaus dataset we provide consists of three
267 conditions: two illusory conditions in which a red centre circle (at different radii) is surrounded by

268 either small or large white circles (flankers) in a configuration that, in humans, induces a biased
 269 size estimation of the centre circle: the circle appears larger when surrounded by small flankers,
 270 everything else being equal. Another condition again contains a red centre circle of different sizes,
 271 but the surrounding circles are placed randomly on the canvas so that they would not elicit any
 272 illusion on a human observer. We use the latter condition to train decoders attached to a ImageNet
 273 pre-trained ResNet-152 model with frozen weights. The task consists of estimating the size of the
 274 centre circle. After training, we feed the illusory images to the decoders. For a network to exhibit the
 275 Ebbinghaus visual illusion, the size of the centre circle should be overestimated for small flankers and
 276 underestimated for big flankers. We did not find this pattern in ResNet-152 and, indeed, no significant
 277 difference across prediction errors for the different conditions was observed (result for one decoder
 278 shown in Figure 2).

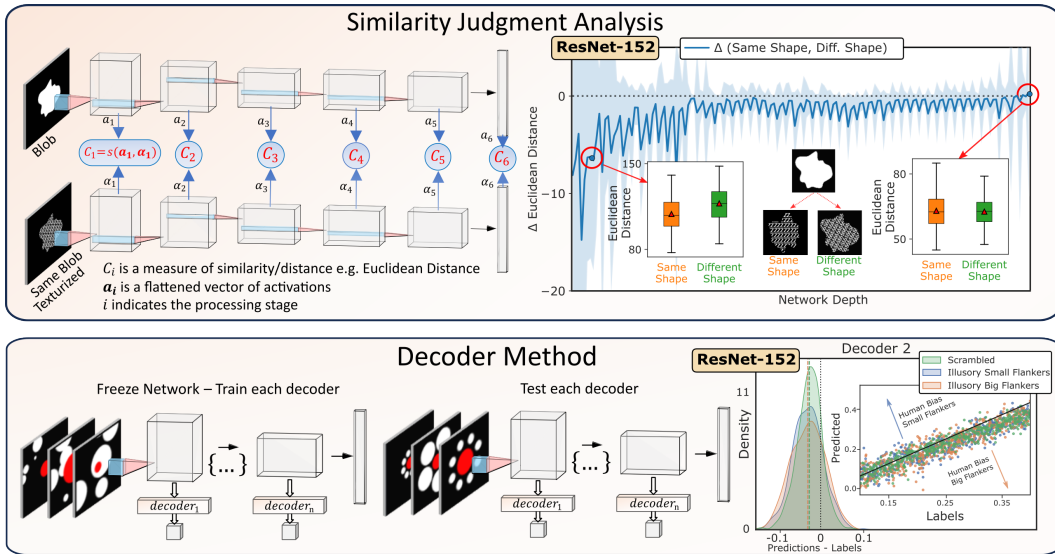


Figure 2: Depiction of two of the three proposed methods of evaluating DNNs in the context of two representative datasets. The first method, out-of-distribution classification, is not depicted here. The Similarity Judgment Analysis (top panel) involves feeding pairs of images to DNNs and comparing the elicited internal representations. We illustrate this method via the ‘Texturized Unfamiliar’ dataset, showing that the network possesses human-like responses in earlier layers which diminish in the later ones. The Decoder Method (bottom panel) involves training and testing a simple linear layer attached to different stages of a frozen network. In the given example, we assess the response to the Ebbinghaus illusion. Our findings indicate an absence of illusory perception. Both examples use an ImageNet pre-trained ResNet-152.

279 4 Code and Resources

280 We provide both ready-to-use datasets and scripts to generate them with varying parameters. Most of
 281 the ready-to-use datasets’ size span to around 5,000 images per condition, and larger dataset can easily
 282 be generated using the provided scripts. To separate code and configuration, each dataset generation
 283 script relies on a configuration file, consisting of a plain-text file in TOML format specifying all
 284 the available parameters for that dataset. Some parameters are used across most datasets, such as
 285 image size, background colour, and number of samples. Other parameters are dataset-specific, for
 286 example the size and distance of dots in the Dotted line drawing dataset. For convenience, the same
 287 configuration file can specify the configuration for multiple (or all) datasets, so that they can be
 288 generated in batches. The “default” configuration file we used to generate the ready-to-use versions is
 289 included, which can be used as a template. The output of each script is the dataset itself (with several
 290 sub-conditions depending on the dataset) together with a CSV annotation file, containing the path
 291 and parameters of each generated image.

292 We also provide the code and utilities to evaluate DNNs using the three methods noted above.
293 Each method is highly configurable through TOML files, with options including the type of data-
294 augmentation to apply, the network architecture, the metric to use for the similarity judgments and
295 more. Users have the flexibility to choose specific factors from this file for analysis, extending beyond
296 the factors that we deemed the most relevant for each task. For example, in the script testing the
297 Ebbinghaus Illusion used for Section 3.3, a decoder is trained to predict the normalized size of the
298 centre circle. However, a different research goal might involve predicting the size of the flankers.
299 This can be achieved by simply specifying the corresponding column ('NormSizeFlankers') in the
300 annotation file, without needing to re-generate the dataset or change the code.

301 Each method produces a pandas DataFrames [78] as its output, which can be independently analyzed.
302 Additionally, supplementary files containing simple tests and comparisons that serve as a springboard
303 for further and more detailed analysis are automatically generated. Comprehensive documentation
304 for every configurable option across all datasets and methods. Additionally, we offer guidance on the
305 general usage of various scripts and utilities through several examples and multiple README files
306 on the GitHub page.

307 **5 Limitations**

308 While *MindSet: Vision* offers a valuable resource for exploring visual psychological phenomena using
309 deep neural networks, there are several limitations to consider. Firstly, our focus is primarily on visual
310 tasks that do not involve high levels of reasoning and are not directly connected with other areas of
311 cognition such as language and memory. Secondly, the methodology for comparing DNN performance
312 to human participants often allows only for qualitative comparisons, as quantitative comparisons may
313 not be feasible with the current analysis methods. Lastly, while we have selected phenomena based
314 on well-replicated and famous visual experiments, there may be additional phenomena that are not
315 covered by our selection. These limitations underscore the need for further research and development
316 in the field of computational modeling of human vision to address these gaps and enhance the utility
317 of *MindSet: Vision* as a comprehensive toolbox for studying visual perception.

318 **6 Conclusion**

319 There is much interest in DNNs as models of human vision, but relatively little research is concerned
320 with how DNNs capture key psychological findings. When DNNs are tested against key psychological
321 findings, they often fail [27]. And when they do succeed, it is often because the DNNs have not been
322 severely tested [26]. In our view, to better characterize DNN-human alignment, and to build better
323 DNN models of human vision, it is necessary to systematically test models against key experiments
324 reported in psychology. The *MindSet: Vision* dataset is designed to facilitate this.

325 Currently it is quite common to rank models in term of how well they perform across several datasets
326 or tasks. For example, the Brain-Score benchmark [96] provides an overall leaderboard that scores any
327 DNN in terms of how good they are at explaining neural activity variance for core object recognition,
328 and the "model-vs-human" benchmark [50] ranks and scores models in terms of their behavioural
329 overlap with humans in identifying a range of out-of-distribution object datasets. We do not propose
330 to rank models in this way as each experiment in *MindSet: Vision* tests a specific hypothesis regarding
331 how DNNs and humans perceive and encode visual inputs. It makes little sense to provide a score
332 that averages across qualitatively different hypotheses. By making stimuli underlying psychological
333 experiments more accessible, easy to generate, configure, and modify, and by providing ready-to-use
334 scripts to test existing models, we hope that the *MindSet: Vision* toolbox encourages computational
335 modelling researcher to focus on testing their models on key experiments rather than competing on
336 observational datasets that do not support any conclusions regarding the mechanistic similarity of
337 DNNs and brains.

338 Acknowledgments and Disclosure of Funding

339 This project has received funding from the European Research Council (ERC) under the European
340 Union’s Horizon 2020 research and innovation programme (grant agreement No 741134).

341 References

- 342 [1] E. H. Adelson. Checkershadow Illusion. *Perceptual Science Group*, 2005.
- 343 [2] Salvatore Aglioti, Joseph F.X. DeSouza, and Melvyn A. Goodale. Size-contrast illusions
344 deceive the eye but not the hand. *Current Biology*, 5(6):679–685, June 1995.
- 345 [3] Christian Agrillo, Michael J. Beran, and Audrey E. Parrish. Exploring the Jastrow Illusion
346 in Humans (*Homo sapiens*), Rhesus Monkeys (*Macaca mulatta*), and Capuchin Monkeys (*Sapajus apella*). *Perception*, 48(5):367–385, May 2019.
- 347 [4] Emily J. Allen, Ghislain St-Yves, Yihan Wu, Jesse L. Breedlove, Jacob S. Prince, Logan T.
348 Dowdle, Matthias Nau, Brad Caron, Franco Pestilli, Ian Charest, J. Benjamin Hutchinson,
349 Thomas Naselaris, and Kendrick Kay. A massive 7T fMRI dataset to bridge cognitive
350 neuroscience and artificial intelligence. *Nature neuroscience*, 25(1):116–126, January 2022.
- 351 [5] Ori Amir, Irving Biederman, and Kenneth J. Hayworth. Sensitivity to nonaccidental properties
352 across various shape dimensions. *Vision Research*, 62:35–43, June 2012.
- 353 [6] Ori Amir, Irving Biederman, Sarah B. Herald, Manan P. Shah, and Toben H. Mintz. Greater
354 sensitivity to nonaccidental than metric shape properties in preschool children. *Vision Research*,
355 97:83–88, 2014.
- 356 [7] Monya Baker. Over half of psychology studies fail reproducibility test. *Nature*, August 2015.
- 357 [8] Nicholas Baker and James H. Elder. Deep learning models fail to capture the configural nature
358 of human shape perception. *iScience*, 25(9), September 2022.
- 359 [9] Nicholas Baker and James H. Elder. Deep learning models fail to capture the configural nature
360 of human shape perception. *iScience*, 25(9):104913, September 2022.
- 361 [10] Nicholas Baker, Hongjing Lu, Gennady Erlikhman, and Philip J. Kellman. Deep convolutional
362 networks do not classify based on global object shape. *PLoS Computational Biology*, 14(12):1–
363 43, 2018.
- 364 [11] H. B. Barlow. Optic nerve impulses and Weber’s law. *Cold Spring Harbor symposia on
365 quantitative biology*, 30:539–546, 1965.
- 366 [12] Ruairidh M. Battleday, Joshua C. Peterson, and Thomas L. Griffiths. Capturing human
367 categorization of natural images by combining deep networks and cognitive models. *Nature
368 Communications 2020 11:1*, 11(1):1–14, October 2020.
- 369 [13] Jacob Beck, K. Prazdny, and Aziel Rosenfeld. A Theory of Textural Segmentation. In *Human
370 and Machine Vision*, pages 1–38. Elsevier, 1983.
- 371 [14] Ari Benjamin, Cheng Qiu, Ling-Qi Zhang, Konrad Kording, and Alan Stocker. Shared visual
372 illusions between humans and artificial neural networks. *2019 Conference on Cognitive
373 Computational Neuroscience*, August 2019.
- 374 [15] Evelin Bertin and Ramesh S. Bhatt. The Thatcher illusion and face processing in infancy.
375 *Developmental Science*, 7(4):431–436, 2004.
- 376 [16] Irving Biederman. Recognition-by-Components: A Theory of Human Image Understanding.
377 *Psychological Review*, M(2):115–147, 1987.
- 378

- 379 [17] Irving Biederman. Recognition-by-Components: A Theory of Human Image Understanding.
380 *Psychological Review*, 94(2):115–147, 1987.
- 381 [18] Irving Biederman and Eric E. Cooper. Priming contour-deleted images: Evidence for interme-
382 diate representations in visual object recognition. *Cognitive Psychology*, 23(3):393–419, July
383 1991.
- 384 [19] Irving Biederman and Ginny Ju. Surface versus edge-based determinants of visual recognition.
385 *Cognitive Psychology*, 20(1):38–64, January 1988.
- 386 [20] Valerio Biscione and Jeffrey S. Bowers. Convolutional Neural Networks Are Not Invariant to
387 Translation, but They Can Learn to Be. *Journal of Machine Learning Research*, 22(229):1–28,
388 2021.
- 389 [21] Valerio Biscione and Jeffrey S. Bowers. Learning online visual invariances for novel objects
390 via supervised and self-supervised training. *Neural Networks*, March 2022.
- 391 [22] Valerio Biscione and Jeffrey S. Bowers. Mixed Evidence for Gestalt Grouping in Deep Neural
392 Networks. *Computational Brain and Behavior*, 6(3):438–456, September 2023.
- 393 [23] Ryan Blything, Valerio Biscione, and Jeffrey Bowers. A case for robust translation tolerance in
394 humans and CNNs. A commentary on Han et al. *arXiv preprint arXiv: 2012.05950*, December
395 2020.
- 396 [24] Ryan Blything, Valerio Biscione, Ivan I Vankov, Casimir J H Ludwig, and Jeffrey S Bowers.
397 The human visual system and CNNs can both support robust online translation tolerance
398 following extreme displacements. *Journal of Vision*, 21(2):1–16, 2021.
- 399 [25] Luc Boutsen, Glyn W. Humphreys, Peter Praamstra, and Tracy Warbrick. Comparing neural
400 correlates of configural processing in faces and objects: An ERP study of the Thatcher illusion.
401 *NeuroImage*, 32(1):352–367, August 2006.
- 402 [26] Jeffrey S. Bowers, Gaurav Malhotra, Federico Adolffi, Marin Dujmović, Milton L. Montero,
403 Valerio Biscione, Guillermo Puebla, John H. Hummel, and Rachel F. Heaton. On the im-
404 portance of severely testing deep learning models of cognition. *Cognitive Systems Research*,
405 82:101158, December 2023.
- 406 [27] Jeffrey S. Bowers, Gaurav Malhotra, Marin Dujmović, Milton Llera Montero, Christian
407 Tsvetkov, Valerio Biscione, Guillermo Puebla, Federico Adolffi, John E. Hummel, Rachel F.
408 Heaton, Benjamin D. Evans, Jeffrey Mitchell, and Ryan Blything. Deep problems with neural
409 network models of human vision. *The Behavioral and brain sciences*, 46, December 2022.
- 410 [28] Jeffrey S. Bowers, Gaurav Malhotra, Marin Dujmović, Milton Llera Montero, Christian
411 Tsvetkov, Valerio Biscione, Guillermo Puebla, Federico Adolffi, John E. Hummel, Rachel F.
412 Heaton, Benjamin D. Evans, Jeffrey Mitchell, and Ryan Blything. Deep Problems with Neural
413 Network Models of Human Vision. *Behavioral and Brain Sciences*, pages 1–74, December
414 2023.
- 415 [29] Astrid Busch and Hermann J. Müller. The Ebbinghaus illusion modulates visual search for
416 size-defined targets: Evidence for preattentive processing of apparent object size. *Perception*
417 *& Psychophysics*, 66(3):475–495, April 2004.
- 418 [30] Charles F. Cadieu, Ha Hong, Daniel L.K. Yamins, Nicolas Pinto, Diego Ardila, Ethan A.
419 Solomon, Najib J. Majaj, and James J. DiCarlo. Deep Neural Networks Rival the Representa-
420 tion of Primate IT Cortex for Core Visual Object Recognition. *PLOS Computational Biology*,
421 10(12):e1003963, December 2014.
- 422 [31] Juan Chen, Irene Sperandio, Molly J. Henry, and Melvyn A. Goodale. Changing the Real
423 Viewing Distance Reveals the Temporal Evolution of Size Constancy in Visual Cortex. *Current*
424 *Biology*, 29(13):2237–2243.e4, July 2019.

- 425 [32] Kai-Xuan Chen, Jie-Yi Ren, Xiao-Jun Wu, and Josef Kittler. Covariance descriptors on a
426 Gaussian manifold and their application to image set classification. *Pattern Recognition*,
427 107:107463, November 2020.
- 428 [33] Colin W.G. Clifford. The tilt illusion: Phenomenology and functional implications. *Vision*
429 *Research*, 104:3–11, November 2014.
- 430 [34] Christoph D. Dahl, Nikos K. Logothetis, Heinrich H. Bülthoff, and Christian Wallraven. The
431 Thatcher illusion in humans and monkeys. *Proceedings of the Royal Society B: Biological*
432 *Sciences*, 277(1696):2973–2981, October 2010.
- 433 [35] R. H. Day and H. Knuth. The Contributions of F C Müller-Lyer.
434 <http://dx.doi.org/10.1068/p100126>, 10(2):126–146, April 1981.
- 435 [36] Lee de-Wit, Hanne Huygelier, Ruth Van der Hallen, Rebecca Chamberlain, and Johan Wage-
436 mans. Developing the Leuven Embedded Figures Test (L-EFT): Testing the stimulus features
437 that influence embedding. *PeerJ*, 5, 2017.
- 438 [37] Mostafa Dehghani, Josip Djolonga, Basil Mustafa, Piotr Padlewski, Jonathan Heek, Justin
439 Gilmer, Andreas Steiner, Mathilde Caron, Robert Geirhos, Ibrahim Alabdulmohsin, Rodolphe
440 Jenatton, Lucas Beyer, Michael Tschannen, Anurag Arnab, Xiao Wang, Carlos Riquelme,
441 Matthias Minderer, Joan Puigcerver, Utku Evci, Manoj Kumar, Sjoerd van Steenkiste,
442 Gamaleldin F. Elsayed, Aravindh Mahendran, Fisher Yu, Avital Oliver, Fantine Huot, Jasmijn
443 Bastings, Mark Patrick Collier, Alexey Gritsenko, Vighnesh Birodkar, Cristina Vasconcelos,
444 Yi Tay, Thomas Mensink, Alexander Kolesnikov, Filip Pavetić, Dustin Tran, Thomas Kipf,
445 Mario Lučić, Xiaohua Zhai, Daniel Keysers, Jeremiah Harmsen, and Neil Houlsby. Scaling
446 Vision Transformers to 22 Billion Parameters, February 2023.
- 447 [38] James J. Dicarlo, Daniel L.K. Yamins, Michael E. Ferguson, Evelina Fedorenko, Matthias
448 Bethge, Tyler Bonnen, and Martin Schrimpf. Let’s move forward: Image-computable models
449 and a common model evaluation scheme are prerequisites for a scientific understanding of
450 human vision. *The Behavioral and brain sciences*, 46, December 2023.
- 451 [39] A. Doerig, A. Bornet, O. H. Choung, and M. H. Herzog. Crowding reveals fundamental
452 differences in local vs. global processing in humans and machines. *Vision Research*, 167:39–
453 45, February 2020.
- 454 [40] Adrien Doerig, Alban Bornet, Ruth Rosenholtz, Gregory Francis, Aaron M. Clarke, and
455 Michael H. Herzog. Beyond Bouma’s window: How to explain global aspects of crowding?
456 *PLOS Computational Biology*, 15(5):e1006580, May 2019.
- 457 [41] Marin Dujmović, Jeffrey S Bowers, Federico Adolfi, and Gaurav Malhotra. The pitfalls of
458 measuring representational similarity using representational similarity analysis. *bioRxiv*, page
459 2022.04.05.487135, April 2022.
- 460 [42] James T. Enns and Ronald A. Rensink. Preattentive recovery of three-dimensional orientation
461 from line drawings. *Psychological review*, 98(3):335–351, 1991.
- 462 [43] Goker Erdogan and Robert A. Jacobs. Visual shape perception as Bayesian inference of 3D
463 object-centered shape representations. *Psychological Review*, 124(6):740–761, 2017.
- 464 [44] Benjamin D. Evans, Gaurav Malhotra, and Jeffrey S. Bowers. Biological convolutions improve
465 DNN robustness to noise and generalisation, September 2021.
- 466 [45] Benjamin D. Evans, Gaurav Malhotra, and Jeffrey S. Bowers. Biological convolutions improve
467 DNN robustness to noise and generalisation. *Neural Networks*, 148:96–110, April 2022.
- 468 [46] Alban Flachot, Arash Akbarinia, Heiko H. Schütt, Roland W. Fleming, Felix A. Wichmann,
469 and Karl R. Gegenfurtner. Deep neural models for color classification and color constancy.
470 *Journal of Vision*, 22(4):17–17, March 2022.

- 471 [47] Gregory Francis, Mauro Manassi, and Michael H. Herzog. Neural dynamics of grouping and
472 segmentation explain properties of visual crowding. *Psychological Review*, 124(4):483–504,
473 July 2017.
- 474 [48] V. H. Franz, F. Scharnowski, and K. R. Gegenfurtner. Illusion effects on grasping are tem-
475 porally constant not dynamic. *Journal of experimental psychology. Human perception and*
476 *performance*, 31(6):1359–1378, December 2005.
- 477 [49] Daniel K. Freeman, Gilberto Graña, and Christopher L. Passaglia. Retinal Ganglion Cell
478 Adaptation to Small Luminance Fluctuations. *Journal of Neurophysiology*, 104(2):704, August
479 2010.
- 480 [50] Robert Geirhos, Kantharaju Narayanappa, Benjamin Mitzkus, Tizian Thieringer, Matthias
481 Bethge, Felix A. Wichmann, and Wieland Brendel. Partial success in closing the gap between
482 human and machine vision. June 2021.
- 483 [51] Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A. Wichmann,
484 and Wieland Brendel. ImageNet-trained CNNs are biased towards texture; increasing shape
485 bias improves accuracy and robustness, November 2022.
- 486 [52] J. J. Gibson and M. Radner. Adaptation, after-effect and contrast in the perception of tilted
487 lines. *Journal of Experimental Psychology*, 20(5):453–467, May 1937.
- 488 [53] Tal Golan, Johnmark Taylor, Heiko Schütt, Benjamin Peters, Rowan P. Sommers, Katja
489 Seeliger, Adrien Doerig, Paul Linton, Talia Konkle, Marcel Van Gerven, Konrad Kording,
490 Blake Richards, Tim C. Kietzmann, Grace W. Lindsay, and Nikolaus Kriegeskorte. Deep
491 neural networks are not a single hypothesis but a language for expressing computational
492 hypotheses. *The Behavioral and brain sciences*, 46, December 2023.
- 493 [54] A. Gomez-Villa, A. Martín, J. Vazquez-Corral, M. Bertalmío, and J. Malo. Color illusions
494 also deceive CNNs for low-level vision tasks: Analysis and implications. *Vision Research*,
495 176:156–174, November 2020.
- 496 [55] Alexander Gomez-Villa, Adrián Martín, Javier Vazquez-Corral, and Marcelo Bertalmío. Con-
497 volutional Neural Networks Deceived by Visual Illusions. November 2018.
- 498 [56] Richard L. Gregory. Knowledge in perception and illusion. *Philosophical Transactions of the*
499 *Royal Society B: Biological Sciences*, 352(1358):1121, August 1997.
- 500 [57] Stephen Grossberg. The quantized geometry of visual space: The coherent computation of
501 depth, form, and lightness. *Behavioral and Brain Sciences*, 6(4):625–657, 1983.
- 502 [58] Stephen Grossberg. The visual world as illusion: The ones we know and the ones we don’t.
503 In *The Oxford Compendium of Visual Illusions*, pages 90–118. Oxford University Press, New
504 York, NY, US, 2017.
- 505 [59] Dongjun He, Ce Mo, Yizhou Wang, and Fang Fang. Position shifts of fMRI-based population
506 receptive fields in human visual cortex induced by Ponzo illusion. *Experimental brain research*,
507 233(12):3535–3541, December 2015.
- 508 [60] M. N. Hebart, O. Contier, L. Teichmann, A. H. Rockter, C. Y. Zheng, A. Kidder, A. Coriveau,
509 M. Vaziri-Pashkam, and C. I. Baker. THINGS-data, a multimodal collection of large-scale
510 datasets for investigating object representations in human brain and behavior. *eLife*, 12,
511 February 2023.
- 512 [61] Robert F. Hess and Anthony Hayes. Neural recruitment explains “weber’s law” of spatial
513 position. *Vision Research*, 33(12):1673–1684, August 1993.
- 514 [62] J. Hochberg and V. Brooks. Pictorial recognition as an unlearned ability: A study of one
515 child’s performance. *The American journal of psychology*, 75:624–628, 1962.

- 516 [63] D. D. Hoffman and W. A. Richards. Parts of recognition. *Cognition*, 18(1-3):65–96, December
517 1984.
- 518 [64] John E. Hummel and Brian J. Stankiewicz. Categorical relations in shape perception. *Spatial*
519 *Vision*, 10(3):201–236, 1996.
- 520 [65] Georgin Jacob, R. T. Pramod, Harish Katti, and S. P. Arun. Qualitative similarities and
521 differences in visual object representations between brains and deep networks. *Nature Com-*
522 *munications* 2021 12:1, 12(1):1–14, March 2021.
- 523 [66] Joseph Jastrow. Studies from the Laboratory of Experimental Psychology of the University of
524 Wisconsin. II. *The American Journal of Psychology*, 4(3):381–428, 1892.
- 525 [67] Nancy Kanwisher, Meenakshi Khosla, and Katharina Dobs. Using artificial neural networks to
526 ask ‘why’ questions of minds and brains. *Trends in Neurosciences*, 46(3):240–254, March
527 2023.
- 528 [68] Greet Kayaert, Irving Biederman, and Rufin Vogels. Shape Tuning in Macaque Inferior
529 Temporal Cortex. *Journal of Neuroscience*, 23(7):3016–3027, April 2003.
- 530 [69] Nikolaus Kriegeskorte, Marieke Mur, Douglas A. Ruff, Roozbeh Kiani, Jerzy Bodurka, Hos-
531 sein Esteky, Keiji Tanaka, and Peter A. Bandettini. Matching categorical object representations
532 in inferior temporal cortex of man and monkey. *Neuron*, 60(6):1126–1141, December 2008.
- 533 [70] Jonas Kubilius, Martin Schrimpf, Kohitij Kar, Rishi Rajalingham, Ha Hong, Najib J Majaj,
534 Elias B Issa, Pouya Bashivan, Jonathan Prescott-Roy, Kailyn Schmidt, Aran Nayebi, Daniel
535 Bear, Daniel L K Yamins, and James J Dicarlo. Brain-Like Object Recognition with High-
536 Performing Shallow Recurrent ANNs. *33rd Conference on Neural Information Processing*
537 *Systems (NeurIPS 2019)*, 2019.
- 538 [71] Jonas Kubilius, Charlotte Sleurs, and Johan Wagemans. Sensitivity to nonaccidental configu-
539 rations of two-line stimuli. *i-Perception*, 8(2):1–12, April 2017.
- 540 [72] Jonas Kubilius, Johan Wagemans, and Hans P. Op de Beeck. Emergence of perceptual Gestalts
541 in the human visual cortex: The case of the configural-superiority effect. *Psychological science*,
542 22(10):1296–1303, 2011.
- 543 [73] Drew Linsley, Junkyung Kim, Alekh Ashok, and Thomas Serre. Recurrent Neural Circuits for
544 Contour Detection. *Iclr*, pages 1–23, 2020.
- 545 [74] Zili Liu, David C. Knill, and Daniel Kersten. Object classification for human and ideal
546 observers. *Vision Research*, 35(4):549–568, 1995.
- 547 [75] Ben Lonnqvist, Zhengqing Wu, and Michael H. Herzog. Latent Noise Segmentation: How
548 Neural Noise Leads to the Emergence of Segmentation and Grouping, September 2023.
- 549 [76] S. P. MacEvoy and M. A. Paradiso. Lightness constancy in primary visual cortex. *Proceedings*
550 *of the National Academy of Sciences of the United States of America*, 98(15):8827–8831, July
551 2001.
- 552 [77] Gaurav Malhotra, Marin Dujmović, John Hummel, and Jeffrey S Bowers. Human shape
553 representations are not an emergent property of learning to classify objects. *bioRxiv*, page
554 2021.12.14.472546, August 2022.
- 555 [78] Wes McKinney. Data Structures for Statistical Computing in Python. In *Python in Science*
556 *Conference*, pages 56–61, Austin, Texas, 2010.
- 557 [79] Johannes Mehrer, Courtney J. Spoerer, Emer C. Jones, Nikolaus Kriegeskorte, and Tim C.
558 Kietzmann. An ecologically motivated image dataset for deep learning yields better models of
559 human vision. *Proceedings of the National Academy of Sciences*, 118(8), February 2021.

- 560 [80] Ken Nakayama, Zijiang J. He, and Shinsuke Shimojo. Visual Surface Representation: A
561 Critical Link between Lower-Level and Higher-Level Vision. *An Invitation to Cognitive*
562 *Science*, October 1995.
- 563 [81] Ken Nakayama and Shinsuke Shimojo. Experiencing and Perceiving Visual Surfaces. *Science*,
564 257(5075):1357–1363, September 1992.
- 565 [82] Aviad Ozana and Tzvi Ganel. A double dissociation between action and perception in bimanual
566 grasping: Evidence from the Ponzo and the Wundt–Jastrow illusions. *Scientific Reports 2020*
567 *10:1*, 10(1):1–10, September 2020.
- 568 [83] Kirsten R. Panton, David R. Badcock, and Johanna C. Badcock. A Metaanalysis of Perceptual
569 Organization in Schizophrenia, Schizotypy, and Other High-Risk Groups Based on Variants of
570 the Embedded Figures Task. *Frontiers in Psychology*, 7, 2016.
- 571 [84] Joshua C. Peterson, Joshua T. Abbott, and Thomas L. Griffiths. Evaluating (and improving) the
572 correspondence between deep neural networks and human representations. *Cognitive Science*,
573 42(8):2648–2669, June 2017.
- 574 [85] James R. Pomerantz and Mary C. Portillo. Grouping and Emergent Features in Vision: Toward
575 a Theory of Basic Gestalts. *Journal of Experimental Psychology: Human Perception and*
576 *Performance*, 37(5):1331–1349, October 2011.
- 577 [86] James R. Pomerantz and Mary C. Portillo. Grouping and emergent features in vision: Toward
578 a theory of basic Gestalts. *Journal of experimental psychology. Human perception and*
579 *performance*, 37(5):1331–1349, October 2011.
- 580 [87] James R. Pomerantz, Lawrence C. Sager, and Robert J. Stoeber. Perception of wholes and
581 of their component parts: Some configural superiority effects. *Journal of Experimental*
582 *Psychology: Human Perception and Performance*, 3(3):422–435, 1977.
- 583 [88] Mario Ponzo. *Intorno ad alcune illusioni nel campo delle sensazioni tattili, sull’illusione di*
584 *Aristotele e fenomeni analoghi*. Wilhelm Engelmann, 1910.
- 585 [89] Guillermo Puebla and Jeffrey S. Bowers. Can deep convolutional neural networks support
586 relational reasoning in the same-different task? *Journal of Vision*, 22(10):11, September 2022.
- 587 [90] R. D.S. Raizada and S. Grossberg. Context-sensitive binding by the laminar circuits of V1
588 and V2: A unified model of perceptual grouping, attention, and orientation contrast. *Visual*
589 *Cognition*, 8(3-5):431–466, 2001.
- 590 [91] Rishi Rajalingham, Elias B. Issa, Pouya Bashivan, Kohitij Kar, Kailyn Schmidt, and James J.
591 DiCarlo. Large-Scale, High-Resolution Comparison of the Core Visual Object Recognition
592 Behavior of Humans, Monkeys, and State-of-the-Art Deep Artificial Neural Networks. *Journal*
593 *of Neuroscience*, 38(33):7255–7269, August 2018.
- 594 [92] Einat Rashal, Aline F. Cretienoud, and Michael H. Herzog. Perceptual grouping leads to
595 objecthood effects in the Ebbinghaus illusion. *Journal of Vision*, 20(8):11–11, August 2020.
- 596 [93] Ronald A. Rensink and James T. Enns. Early completion of occluded objects. *Vision Research*,
597 38(15-16):2489–2505, August 1998.
- 598 [94] Irvin Rock. *An Introduction to Perception*. Macmillan, New York, NY, 1975.
- 599 [95] Toni P. Saarela, Bilge Sayim, Gerald Westheimer, and Michael H. Herzog. Global stimulus
600 configuration modulates crowding. *Journal of Vision*, 9(2):5–5, February 2009.
- 601 [96] Martin Schrimpf, Jonas Kubilius, Ha Hong, Najib J. Majaj, Rishi Rajalingham, Elias B. Issa,
602 Kohitij Kar, Pouya Bashivan, Jonathan Prescott-Roy, Kailyn Schmidt, Daniel L. K. Yamins,
603 and James J. DiCarlo. Brain-Score: Which Artificial Neural Network for Object Recognition
604 is most Brain-Like? *bioRxiv*, page 407007, September 2018.

- 605 [97] Martin Schrimpf, Jonas Kubilius, Michael J. Lee, N. Apurva Ratan Murty, Robert Ajemian,
606 and James J. DiCarlo. Integrative Benchmarking to Advance Neurally Mechanistic Models of
607 Human Intelligence. *Neuron*, 108(3):413–423, November 2020.
- 608 [98] Thomas Serre, Drew Linsley, and Junkyung Kim. What is the function of the orientation-tilt
609 illusion? *Journal of Vision*, 20(11):868, October 2020.
- 610 [99] Johannes J.D. Singer, Katja Seeliger, Tim C. Kietzmann, and Martin N. Hebart. From photos
611 to sketches - how humans and deep neural networks process objects across different levels of
612 visual abstraction. *Journal of Vision*, 22(2):4–4, February 2022.
- 613 [100] Chen Song and Geraint Rees. Intra-hemispheric integration underlies perception of tilt illusion.
614 *NeuroImage*, 175:80–90, July 2018.
- 615 [101] Valeria Anna Sovrano, Osvaldo da Pos, and Liliana Albertazzi. The Müller-Lyer illusion in
616 the teleost fish *Xenotoca eiseni*. *Animal cognition*, 19(1):123–132, January 2016.
- 617 [102] Irene Sperandio, Philippe A. Chouinard, and Melvyn A. Goodale. Retinotopic activity in
618 V1 reflects the perceived and not the retinal size of an afterimage. *Nature neuroscience*,
619 15(4):540–542, April 2012.
- 620 [103] Katherine R. Storrs, Barton L. Anderson, and Roland W. Fleming. Unsupervised learning
621 predicts human perception and misperception of gloss. *Nature Human Behaviour* 2021 5:10,
622 5(10):1402–1417, May 2021.
- 623 [104] Keiji Tanaka. Inferotemporal Cortex and Object Vision. *Annual Review of Neuroscience*,
624 19(1):109–139, March 1996.
- 625 [105] P. Thompson. Margaret Thatcher: A new illusion. *Perception*, 9(4):483–484, 1980.
- 626 [106] Katrien Torfs, Kathleen Vancleef, Christophe Lafosse, Johan Wagemans, and Lee De-Wit. The
627 Leuven Perceptual Organization Screening Test (L-POST), an online test to assess mid-level
628 visual perception. *Behavior Research Methods*, 46(2):472–487, November 2014.
- 629 [107] Michel Treisman. Noise and Weber’s law: The discrimination of brightness and other dimen-
630 sions. *Psychological Review*, 71(4):314–330, July 1964.
- 631 [108] Shikhar Tuli, Ishita Dasgupta, Erin Grant, and Thomas L. Griffiths. Are Convolutional Neural
632 Networks or Transformers more like human vision? *Proceedings of the 43rd Annual Meeting
633 of the Cognitive Science Society: Comparative Cognition: Animal Minds, CogSci 2021*, pages
634 1844–1850, May 2021.
- 635 [109] H. von Helmholtz. *Handbuch Der Physiologischen Optik*. Leipzig: Voss, 1867.
- 636 [110] Emily J Ward. Exploring perceptual illusions in deep neural networks. *Journal of Vision*,
637 19(10):34b–34b, September 2019.
- 638 [111] Eiji Watanabe, Akiyoshi Kitaoka, Kiwako Sakamoto, Masaki Yasugi, and Kenta Tanaka.
639 Illusory Motion Reproduced by Deep Neural Networks Trained for Prediction. *Frontiers in
640 Psychology*, 9, March 2018.
- 641 [112] E. H. Weber. Der Tastsinn und das Gemeingefühl. *Handwörterbuch der Physiologie*, 3:481–
642 588, 1983.
- 643 [113] Gerald Westheimer. Simultaneous orientation contrast for lines in the human fovea. *Vision
644 Research*, 30(11):1913–1921, January 1990.
- 645 [114] Robert L. Whitwell, Mehul A. Garach, Melvyn A. Goodale, and Irene Sperandio. Looking at
646 the Ebbinghaus illusion: Differences in neurocomputational requirements, not gaze-mediated
647 attention, explain a classic perception-action dissociation. *Philosophical Transactions of the
648 Royal Society B: Biological Sciences*, 378(1869):20210459, December 2022.

- 649 [115] Yetta K Wong, Elyssa Twedt, David Sheinberg, and Isabel Gauthier. Does Thompson's
650 Thatcher Effect Reflect a Face-Specific Mechanism? *Perception*, 39(8):1125–1141, August
651 2010.
- 652 [116] Yetta K Wong, Elyssa Twedt, David Sheinberg, and Isabel Gauthier. Does Thompson's
653 Thatcher Effect Reflect a Face-Specific Mechanism? *Perception*, 39(8):1125–1141, August
654 2010.
- 655 [117] Yaoda Xu and Manish Singh. Early computation of part structure: Evidence from visual
656 search. *Perception & psychophysics*, 64(7):1039–1054, 2002.
- 657 [118] Hongtao Zhang, Zhen Li, and Shinichi Yoshida. Müller-Lyer illusion is Replicated by Higher
658 Layer of Pre-trained Deep Neural Network for Object Recognition Müller-Lyer illusion is
659 Replicated by Higher Layer of Pre-trained Deep Neural Network for Object Recognition. In
660 *The 10th International Symposium on Computational Intelligence and Industrial Applications*
661 (*ISCIIA2022*, Zhang2022mullerlyer, 2022).
- 662 [119] Yichi Zhang, Jiayi Pan, Yuchen Zhou, Rui Pan, and Joyce Chai. Grounding Visual Illusions in
663 Language: Do Vision-Language Models Perceive Illusions Like Humans? October 2023.
- 664 [120] Chengxu Zhuang, Siming Yan, Aran Nayebi, Martin Schrimpf, Michael C. Frank, James J.
665 DiCarlo, and Daniel L. K. Yamins. Unsupervised Neural Network Models of the Ventral
666 Visual Stream. *bioRxiv*, page 2020.06.16.155556, 2020.

667 **Checklist**

- 668 1. For all authors...
- 669 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
670 contributions and scope? [Yes]
- 671 (b) Did you describe the limitations of your work? [Yes] Yes, Section 5.
- 672 (c) Did you discuss any potential negative societal impacts of your work? [N/A] *We do*
673 *not believe this work could have any negative societal impact.*
- 674 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
675 them? [Yes]
- 676 2. If you are including theoretical results...
- 677 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 678 (b) Did you include complete proofs of all theoretical results? [N/A]
- 679 3. If you ran experiments (e.g. for benchmarks)... *The two experiments we ran are only pre-*
680 *sented as illustrations of how the datasets could be tested with our suggested methodologies,*
681 *and are not intended as benchmarks.*
- 682 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
683 mental results (either in the supplemental material or as a URL)? [Yes] Provided code
684 in GitHub repo contains details instruction to replicate the exemplary results presented
685 in this work.
- 686 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
687 were chosen)? [Yes] We used .toml file which contains parameters for each dataset
688 generation process and each methodology, so that full replicability is ensured.
- 689 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
690 ments multiple times)? [Yes] see 2.
- 691 (d) Did you include the total amount of compute and the type of resources used (e.g., type
692 of GPUs, internal cluster, or cloud provider)? [N/A]
- 693 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 694 (a) If your work uses existing assets, did you cite the creators? [Yes] Creators are cited
695 within the main text, and in the codebase in each file using their assets. All used assets
696 are either publicly available under CC or we obtained the explicit permission from the
697 authors.
- 698 (b) Did you mention the license of the assets? [Yes]
- 699 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
700 *We provide Kaggle links for two different versions of the datasets, and a GitHub repo*
701 *to regenerate the datasets.*
- 702 (d) Did you discuss whether and how consent was obtained from people whose data you’re
703 using/curating? [Yes] *in Section 2.*
- 704 (e) Did you discuss whether the data you are using/curating contains personally identifiable
705 information or offensive content? [No] *data are programmatically generated, and*
706 *their material is not offensive nor contains identifiable information.*
- 707 5. If you used crowdsourcing or conducted research with human subjects... *The datasets are*
708 *ideally used to test psychological phenomena observed and studied in human subjects, but*
709 *we did not directly conduct any research with humans participants.*
- 710 (a) Did you include the full text of instructions given to participants and screenshots, if
711 applicable? [N/A]
- 712 (b) Did you describe any potential participant risks, with links to Institutional Review
713 Board (IRB) approvals, if applicable? [N/A]
- 714 (c) Did you include the estimated hourly wage paid to participants and the total amount
715 spent on participant compensation? [N/A]

716 **Appendices**

717 **A General Dataset Info**

718 **A.1 Pre-Generated Datasets**

719 In the pre-generated dataset, we use 224x224 pixel images, and a variable number of samples
720 depending on the task and condition (see section below on the number of samples). However, image
721 size and dataset sizes are parameters that the user can easily modify if needed. The images in the
722 datasets all have 3 channels (RGB). For almost all datasets, the user can specify the background color
723 (either a uniform value, or request a different RGB value for every image), whether to use antialiasing,
724 and the size of the item in the image relative to the whole canvas.

725 **A.2 Data Augmentation**

726 We do not apply any affine transformation or other data augmentation techniques during the dataset
727 generation phase. For instance, in the majority of Shape and Object Recognition datasets, a sample
728 typically comprises a modified line drawing or silhouette centrally positioned on a canvas without
729 rotation. We deliberately avoid creating replicas of the same sample with additional transformations.
730 This approach prevents unnecessary expansion of the dataset’s size, as most popular deep learning
731 libraries allow for an easy application of data augmentation. Furthermore, our testing methods
732 allow for the application of affine transformation during testing through the configuration file, again
733 avoiding the need to generate pre-augmented datasets.

734 **A.3 Number of Samples and procedural generations**

735 Some datasets contain a fixed and limited number of samples. For example, the “NAP vs MP: line
736 segments” dataset, recreating the stimuli used in [72], contains 3 conditions with 26 items each. This
737 can be potentially expanded by changing the line stroke, the background color, or by applying dataset
738 augmentation separately. For other datasets much larger samples with extremely low probability of
739 repetition are easily constructed. For example, most Visual Illusions contain a “scrambled” condition
740 in which the elements of the illusions are presented “scrambled up” on the canvas, with varying
741 positions and orientation of each different element. A virtually limitless number of samples can be
742 generated for these conditions, which is important since they are often used for training decoders. The
743 pre-generated dataset typically includes approximately 5,000 samples for these conditions. For users
744 requiring larger sample sizes, they can generate the dataset by changing the configuration argument
745 relative to sample size (e.g., num_samples_scrambled) depending on their needs.

746 **B Other Testing Methods**

747 The datasets featured in MindSet: Vision are suitable for a range of experimental approaches beyond
748 the ones suggested in Section 3. For example, the relational reasoning capabilities needed to solve the
749 Same/Different task could be tested by training a network (which does not need to be pre-trained) on
750 one or a few of the ten available conditions in the Same/Different dataset, and then testing the network
751 on the remaining conditions (as in [89]). In a similar vein, the 2D transformations and Viewpoint
752 datasets could be approached by training a network on certain transformations/viewpoints and then
753 testing it on the others, as demonstrated in [20].

754 Another method that extends beyond our provided scripts is to input images into Multimodal Large
755 Language Models (LLMs) and query them about what they see [119]. Assuming that the language
756 output of the LLM provides a reliable window into its perceptual processes, this approach allows for
757 an interactive examination of the LLM’s understanding of the images.

758 Our preliminary investigations with GPT-4 reveal that while these models are proficient at recognizing
759 silhouettes and line drawings, they struggle with textured representations of familiar objects. For

760 example, a textured image of a banana was misidentified as either a crescent moon or a pair of
 761 scissors, and an airplane was mistaken for a butterfly. Moreover, we found instances in which GPT-4
 762 is influenced from images it has previously processed, that is prior exposure to an image can lead
 763 the model to incorrectly identify later, differently textured images as the same as the initially viewed
 764 object. For instance, when the model is first presented with the silhouette of a banana, followed by a
 765 textured depiction of an airplane, it sometimes erroneously classifies the airplane as a banana.

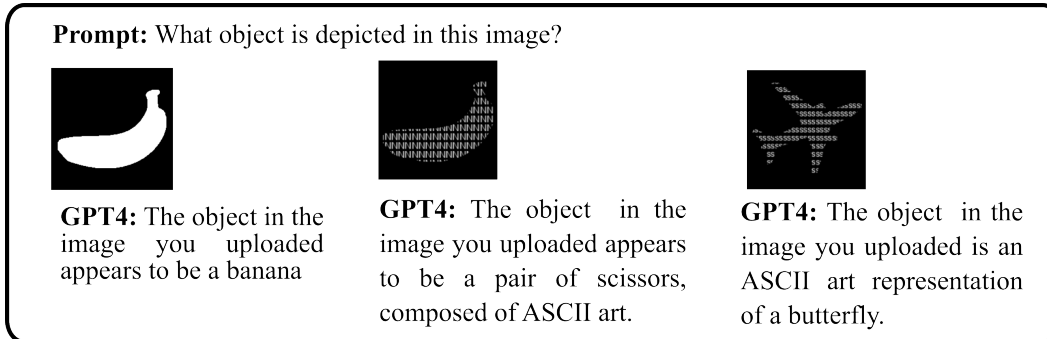


Figure 3: Samples of GPT4 responses after being prompted by different images from our silhouettes and texturized datasets, each time spawning a new conversation. The model accuracy drops significantly with texturized images.

766 C Detailed Datasets Information

767 Below we provide more details about the various conditions included in MindSet: Vision and their
 768 relevance to understanding human vision. We also highlight the most important parameters for each
 769 dataset. For a complete list, refer to the generated HTML page² which additionally contains several
 770 samples for each condition and each dataset when generated using default parameters.

771 C.1 Low and Mid-level vision

772 There is no sharp dividing line between early and middle visual processing, but early vision extracts
 773 low level feature information (i.e., color and luminance contrasts, the orientation of bar and edge
 774 segments, and contour segments) from the retinal image. By contrast, mid level vision encodes more
 775 abstract aspects of shape, such as surfaces and parts of objects, in an increasingly viewpoint invariant
 776 manner. This is where representations of 2D and 3D shapes, material properties, and coherence of the
 777 substances and surfaces in the world are computed. That is, mid-level vision builds representations of
 778 the distal world from the proximal stimulus. This is an ill-posed problem, and accordingly, various
 779 heuristics (such as Gestalt perceptual grouping cues) are employed to provide the best estimate of the
 780 distal world. Illusions are striking examples of failures to correctly encode the distal world from the
 781 proximal stimulus. The experiments we include in MindSet: Vision largely focus on mid-level vision.

782 C.1.1 Weber Law

783 **Psychological Significance.** The Weber Law (or Weber-Fechner Law) quantifies the psychophysical
 784 relation between changes in the world and changes in perception. The law states that the minimum
 785 physical change of a stimulus on some dimension (e.g., its size) that is perceptible to an observer is a
 786 constant ratio of the original stimulus value on this dimension. For example, it is easy to distinguish
 787 between line lengths of 1 and 2 cm, but difficult to distinguish between lines of length 100 and 101
 788 cm, despite the fixed difference in length (1 cm). To make the latter distinction equally salient, the
 789 two stimuli should be 100 and 200 cm (a fixed ratio of 2). Although Weber’s Law breaks down at
 790 extreme values, this relationship applies to a wide range of dimensions, from weight, length, size,

²<https://bit.ly/mindsetvision-datasets>

791 brightness, and even numbers. Weber’s Law reflects the more general observation that perception is
 792 often based on relative rather than absolute encoding of stimulus dimensions. Importantly, Weber’s
 793 Law is often manifest in early visual areas [61], and indeed, in some cases, at the level of the retina
 794 [11, 49, 107].

795 Jacob et al. [65] reported that the convolutional DNN VGG-16 showed a human-like Weber Law
 796 effect when encoding line-lengths. However, the authors only observed Weber’s Law for line lengths
 797 in the late convolutional layers of the network, did not assess whether discrimination was a constant
 798 ratio of the original stimulus (they employed a weaker test), and failed to observe a reliable effect for
 799 image intensity.

800 **Dataset.** Images in this dataset are composed of a simple horizontal white line with varying length
 801 and brightness values. Configurable parameters include line width, min/max values for length and
 802 brightness. To assess DNNs sensitivity to Weber’s Law, a similarity judgment analysis assesses
 803 whether the relative change in the perception of these stimuli (as measured by the level of unit
 804 activation in the inner layers of a pre-trained DNN) adheres to a logarithmic relationship with the
 805 stimulus strength (e.g. line length).

806 C.1.2 Crowding / Uncrowding

807 **Psychological Significance.** Our ability to identify objects is impaired by the presence of nearby
 808 objects and shapes, a phenomenon called crowding. At the same time, in some conditions, the
 809 inclusion of additional surrounding objects makes the identification of the target easier, a phenomenon
 810 called uncrowding. This is illustrated in Figure 4, in which participants are asked to perform a
 811 vernier discrimination task by deciding whether the top vertical line from a pair of vertical lines
 812 is shifted to the left or right. When these lines are surrounded by a square rather than presented
 813 by themselves performance is impaired. However, the inclusion of additional squares dramatically
 814 improves performance. This is thought to reflect a Gestalt process in which the squares are grouped
 815 together and then processed separately from the vernier [95]. Standard DNNs are unable to explain
 816 uncrowding [39, 47], but the DNNs inspired by the LAMINART model of Grossberg and colleagues
 817 [90] designed to support grouping processes can capture some aspects of uncrowding.

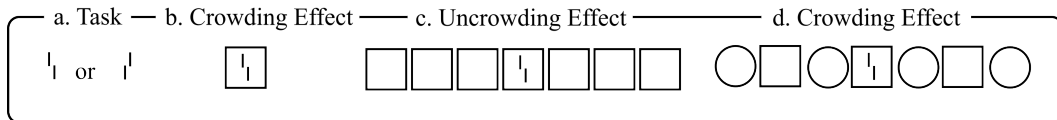


Figure 4: Illustration of the Crowding and Uncrowding effect. a. Observers perform a vernier discrimination task. A standard approach consists of measuring the vernier offset for which observers correctly discriminate in 75% of the trials. With the vernier alone, the offset is quite small. b. When a square is added the performance drastically drops (that is, the threshold-offset increases). This is the classic crowding effect. c. Adding more flankers increases performance again. This is referred to as uncrowding. d. The magnitude of crowding and uncrowding effects is contingent upon both short-range and long-range spatial interactions between visual elements. Furthermore, the specific characteristics and spatial positioning of flanker stimuli play a crucial role in modulating these effects. For example, the performance drops again for the depicted pattern.

818 **Dataset.** Based on [39], with code adapted with authors’ permission. Images are composed of a
 819 ‘vernier’ stimulus (two parallel line segments with some offset) placed either inside or outside a set of
 820 random flankers (squares, circles, hexagons, octagons, stars, diamonds). Each configuration has from
 821 1 to 7 columns and from 1 to 3 rows of flankers with a variety of same/different shape patterns used.
 822 The vernier can be left/right oriented. The suggested method for this dataset (as per [39]) consists
 823 of attaching a decoder at several stages of a pre-trained DNN. The decoder is trained and tested
 824 on a classification task to discriminate between left/right types of vernier but, significantly, during
 825 training, the vernier and the flankers were non-overlapping, whereas during test, the vernier was often
 826 placed inside one of the shapes, allowing the measuring of (un)crowding effect through change in
 827 classification accuracy across test conditions. A model with human-like visual characteristics should

828 match human perception with regards to both crowding and uncrowding effects, following the pattern
 829 in [39, 40]. Users can specify whether the size of the flankers varies or is fixed across samples.

830 C.1.3 Emergent features

831 **Psychological Significance.** Emergent features provide a compelling example of “the whole is
 832 different than the sum of its parts”. Pomerantz and colleagues [85, 87] relied on a simple visual
 833 search paradigm where participants were asked to identify a target amongst foils. They devised
 834 several different types of target and foil stimuli, but the simplest were composed of dot patterns as
 835 depicted below. Participants viewed a set of 4 panels, each of which contained a single dot. Three of
 836 these panels were identical (dots were in the same location) and one outlier panel (where the dot was
 837 in a different location). The task was to identify the outlier panel as quickly as possible. In the single
 838 dot condition, the outlier was simply the panel with a dot in a unique position. In the critical emergent
 839 feature condition(s), a dot (or more) was added to the single dot images as context. The context dot(s)
 840 was in the same location in all panels. Because these added dot(s) were identical in all four panels,
 841 there were no new features that could be used to facilitate the identification of the outlier other than
 842 configural “emergent” features. For example, in the top row of Figure 5, the extra dot (depicted in the
 843 middle column) produces the emergent feature of “orientation”, and in the bottom row, the extra dot
 844 produces the emergent feature of proximity. The critical finding was that participants could identify
 845 the location of the outlier panel more quickly in the emergent compared to the baseline condition.
 846 That is, the “whole” was more discriminable than the sum of its parts.

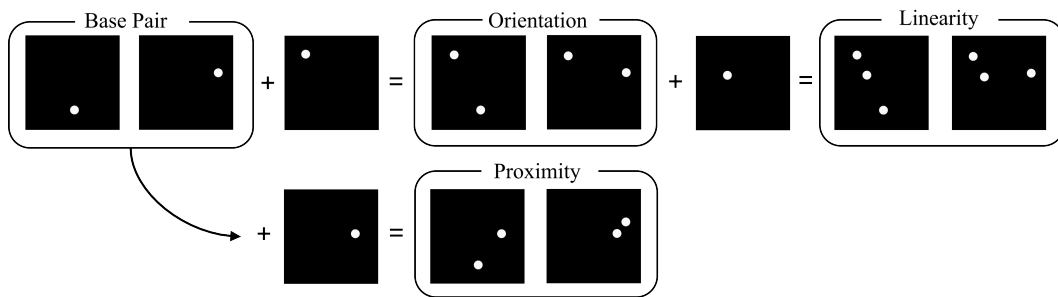


Figure 5: Schematic of the generation procedure for producing a set of dotted stimuli. Starting with a pair of images in which the only discriminant feature is the location of a dot (Base Pair), an additional dot is added, yielding the Emergent Feature of proximity or orientation. The Emergent Feature of linearity is obtained by adding a dot to the orientation pair. Notice that the added dot is the same to both elements of the pair so it does not add on its own any discriminative features, but it generates additional features in relation with the surrounding dots.

847 Biscione and Bowers [22] carried out a series of studies assessing whether DNNs were sensitive to a
 848 range of emergent features that facilitated human performance, including testing DNNs on the dot
 849 stimuli illustrated in Figure 5. We observed that DNNs did show some sensitivity to some of the
 850 emergent features, but only at the later layers of the network. This is problematic given that these
 851 emergent features are thought to be computed relatively early in the visual system, such that they
 852 support rapid “pop out” search.

853 **Dataset.** Adapted from [22]. The dataset consisted of sets of paired images. Each set includes four
 854 conditions: a base condition (single dots), and composite conditions (orientation, proximity, and
 855 linearity). The ‘single dots’ condition consists of paired images in which each image contains a
 856 single dot placed at a different location. In the composite conditions, one or more dots are added to
 857 both images of the base condition, in the same locations, in such a way that it would elicit different
 858 emergent properties when combined with the original single dots. In the orientation and proximity
 859 conditions, the added dot results in different orientation/proximity features. In the linearity condition
 860 (generated by adding a dot to the orientation condition), the added dot would either be placed on a
 861 straight line with the other two dots or on a different path. Each dot was constrained to be located
 862 at a distance of at least 20 pixels from one another, and 40 pixels from the border. By computing

863 the difference in similarity scores between each composite condition and the base condition, we can
 864 compute how much each emergent feature impairs/facilitates distinction of the additional dots. For
 865 example, if the average ‘orientation’ pair is found to be easier to distinguish (through a similarity
 866 analysis of the internal activations of the network) than the pairs from the ‘single dot’ condition,
 867 then we can infer that the network is sensitive to orientation (as the additional dot in the orientation
 868 condition was not-diagnostic, e.g. the same for both images in each pair). The same comparison with
 869 the ‘single dot’ pairs can be performed for the proximity and linearity conditions. The overall pattern
 870 of similarity scores should match human results, in which the highest effect is obtained through the
 871 feature of proximity, followed by linearity, and then orientation [85, 22].

872 C.1.4 Decomposition

873 **Psychological Significance.** The visual system represents objects in terms of their parts, separating
 874 regions at points of deep concavity [63]. Perceptually, searching for an object broken into its natural
 875 parts among a set of unsegmented versions of the same object is significantly more challenging than
 876 locating the same object when it is segmented at points that do not correspond to its natural divisions.
 877 In other words, a segmentation at natural points preserves the basic parts which make up the object
 878 and therefore make the segmented version more similar to the uncut object when compared to an
 879 ‘unnatural’ segmentation. There is good evidence that this occurs relatively early in visual processing
 880 [117]. To assess whether DNNs encode objects into parts in a similar manner, Jacob et al. [65]
 881 compared the internal representations of a base object composed of two parts to two segmentations
 882 of the object, one natural and one unnatural. The assumption is that a natural segmentation of the
 883 image will be encoded in a more similar way to the whole object (the segmented images maintain the
 884 integrity of parts that compose the complete object). However, they reported that the VGG-16 did not
 885 show this pattern, suggesting that DNNs do not encode objects by their parts, or at least, not in a way
 886 similar to humans.

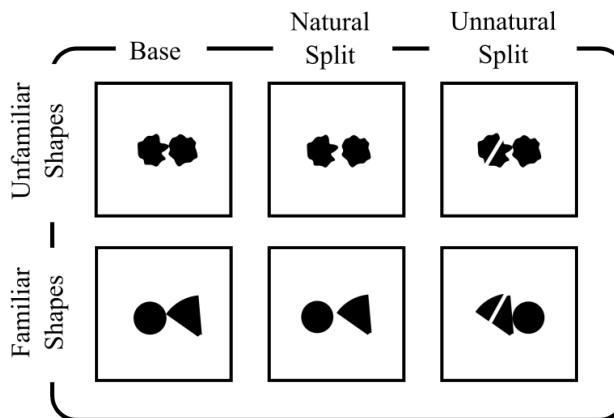


Figure 6: The dataset features base images depicting two objects in contact at a single point. It includes two variations: natural and unnatural splits. In natural splits, the objects are separated, while in unnatural splits, the division occurs within an object itself. Identifying differences between base images and unnatural splits is simpler than distinguishing between base and natural splits. The dataset presents examples with both familiar and unfamiliar shapes, showcasing the diversity in object recognition challenges.

887 **Dataset.** The dataset consists of a variation of the images used in [65]: instead of a single object
 888 composed of two parts, we used two objects joined at a single point of contact. There are three
 889 ‘split’ conditions and two ‘familiarity’ conditions. The ‘split’ conditions are: ‘no split’ in which two
 890 parts are touching at one point but not overlapping; ‘natural split’, in which two parts are separated;
 891 ‘unnatural split’ in which the two parts are touching each other as in the ‘no split’ condition, but one
 892 of the parts is ‘cut’ and separated from the rest. The items are silhouettes uniformly coloured on
 893 a uniform background, and they can be either familiar or unfamiliar shapes. The familiar shapes

894 consist of the following objects: circle, square, rectangle, triangle, heptagon, and a 50-degree arc
 895 segment; the unfamiliar shapes consist of blob-like objects. Within each familiar/unfamiliar condition,
 896 all possible combinations of two shapes are used (e.g. a triangle with a rectangle). Configuration
 897 parameters include the distance between pieces in the ‘unnatural split’ and ‘natural split’ conditions,
 898 the colour of items, and the number of different blob-like objects to use for the unfamiliar condition.
 899 Following the test from [65], similarity judgments between pairs composed of base samples and
 900 natural/unnatural splits can be computed for an ImageNet pre-trained network. To match human
 901 perception, the natural split samples should have internal representations that are closer to the base
 902 samples than the unnatural split samples. This should apply regardless of whether the shapes are
 903 familiar or unfamiliar.

904 **C.1.5 Encoding relations between object parts**

905 **Psychological Significance.** Humans not only encode objects in terms of their parts, but also the
 906 relations between parts which are essential for object recognition [16]. Early evidence for this was
 907 reported by [64] who trained participants to identify a small set of artificial stimuli in which they
 908 could easily manipulate relations between parts. Two types of changes were introduced to create foils
 909 for the base stimuli. First, a coordinate change in which relations between parts were maintained
 910 but the position of a part of the object was changed. And second, a relational change in which there
 911 was a categorical change in relations between object parts. They reported that participants were
 912 much more likely to mistake foil objects for the base object when the relations between object parts
 913 were maintained than when the relations changed (coordinate vs relational change in Figure 7). By
 914 contrast, [77] showed that two standard convolutional networks are completely insensitive to these
 915 relational features, treating Relational and Coordinate foils equally similar to the Basis objects.

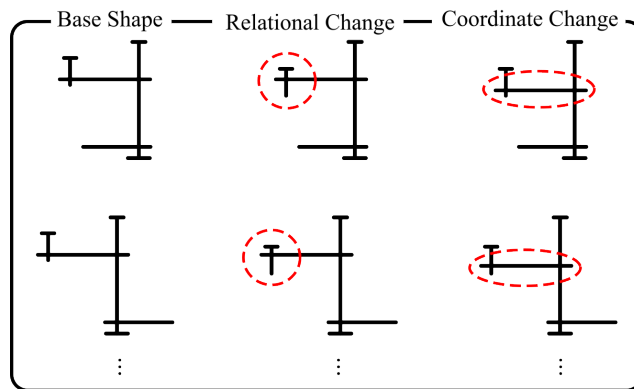


Figure 7: Reproduction of stimuli used in [64]. Starting with a base shape, the Relational Change variant was created by moving one part of the base object up or down (red dashed circle). The move was chosen to change the categorical above/below relation between the circled part and the part to which it is attached. The coordinate change variant was created by moving the whole horizontal (red circled) segment up or down, together with the part moved in the relational change. This resulted in no categorical relations change. Therefore, the perceived difference between a base and its relational-change pair is greater than the perceived difference between the corresponding base-coordinate change pair.

916 **Dataset.** We recreated images originally contained in [64], Experiment 5, using white strokes on a
 917 uniform background. To compare to human perception, similarity judgments can be computed from
 918 pre-trained DNNs by sequentially inputting pairs of images composed of a base shape and either their
 919 corresponding coordinate or relational change. A pattern that mirrors human perception would result
 920 in greater similarity between the base shapes and their coordinate modifications foils as opposed to
 921 their relational change foils.

922 **C.1.6 Encoding of 3D shapes**

923 **Psychological Significance.** The human visual system builds 3D representations of images for the
924 sake of object recognition [43, 74], and some perceptual illusions of size, such as the Ponzo illusion
925 described below, are thought to be a by-product of computing depth information. By contrast, there
926 is little evidence that DNNs infer 3D structure from the 2D images they process. For example, [65]
927 tested VGG-16 on three pairs of objects developed by [42]: a pair of objects composed of three
928 segmented lines (base pair in Figure 8) are transformed in two different ways (V1 and V2), each
929 time adding the same configuration to both elements of the base pair. Humans were assessed in how
930 quickly they could discriminate the two V1 images and the two V2 images. Discrimination was
931 highly improved for the V2 pair, but not for the V1 pair, most likely the result of enhanced 3D cues
932 in the V2 stimuli.

933 In contrast, Jacob et al. [65] obtained no evidence that VGG-16 was better at discriminating the base
934 pair, suggesting a failure to encode their 3D structure.

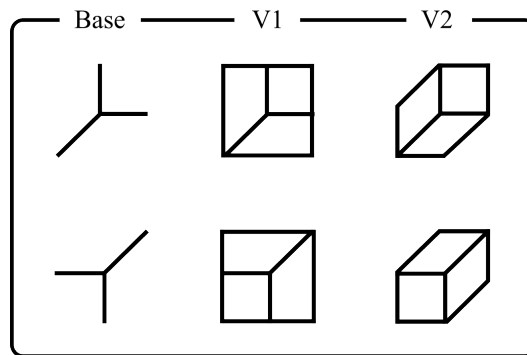


Figure 8: Illustration of the 3D Drawing dataset stimuli. One where segmented lines (Base shapes) are augmented with contextual features to clearly form distinguishable 3D shapes (V2), and another where the additions do not contribute as strongly to depth perception, making shape discrimination challenging (V1). Importantly, the identical contextual features applied to each pair highlight that enhanced discrimination stems solely from the perceived depth, rather than the features themselves.

935 **Dataset.** We recreated stimuli appearing in [42], using white strokes on a uniform background. Using
936 the similarity judgment method on pre-trained DNNs, a perception akin to humans would result in a
937 significantly lower similarity for the V2 pair compared to both the base and V1 images.

938 **C.1.7 Amodal completion**

939 **Psychological Significance:** The visual system needs to identify partly occluded objects in the 3D
940 world. A key part of the solution for humans is an amodal completion process in which a surface
941 representation of the occluded object is completed behind the occluder. This process is called amodal
942 because the visual system builds complete surface forms of occluded objects without generating
943 a visible experience of the missing shape. Amodal completion occurs early in the visual system,
944 perhaps as early as V1 [81]. Various compelling perceptual effects are associated with amodal
945 completion (for review, see [80]). Here we include the materials of [93] who showed that humans
946 quickly and automatically encode the shape of partially occluded objects in a visual search task.
947 Amongst the various conditions in their experiments, two illustrate the point most clearly. In the
948 'Target - Notched square' and 'Target - notched circle' conditions, participants searched for a notched
949 black square or notched white circle, respectively, among full black square and white circle distractors.
950 None of the objects overlapped in this condition. In the 'Occlusion' condition participants were again
951 searching for notched squares and circles, but in this case notched squares and circles touched to
952 give the impression of occlusion. Search was significantly faster in the 'Notched' condition when
953 compared to the 'Occlusion' condition. This is because in the 'Occlusion' condition the notched
954 squares were perceived as full squares occluded by white disks due to amodal completion. This made

955 the notched black squares much more difficult to find among full black squares. The same was true
 956 for the notched white circles in the occlusion condition. This pattern of results suggests that the
 957 notched square in the 'Occlusion' condition was encoded as a square early in visual processing (fast
 958 visual search is typically characterized as pre-attentive). Jacob et al. [65] reported that the DNN
 959 VGG-16 network pre-trained on ImageNet failed to show any evidence for amodal completion with
 960 these stimuli.

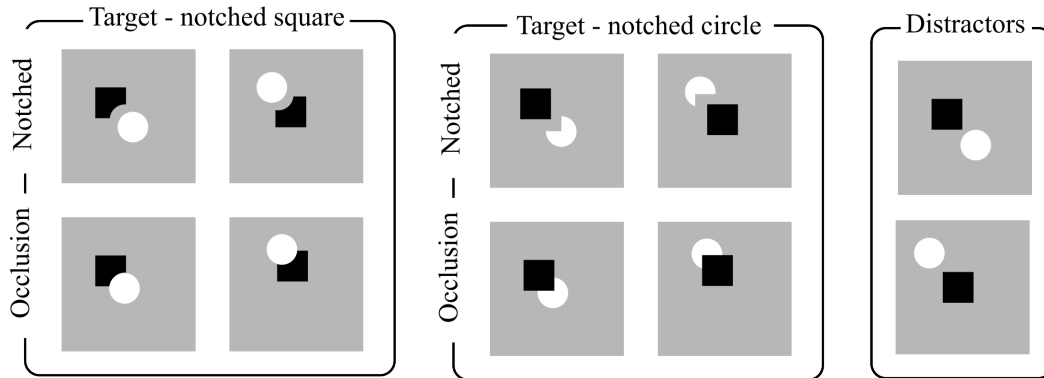


Figure 9: Illustration of the Occluded Shape stimuli as used in [93]. The experiment compares three conditions: a baseline condition with squares and disks with no occlusion, an occlusion condition in which one object obscures part of the other, and a notched condition in which the occluded part of the object in occlusion condition is removed. This turns the notch into a distinctive feature and effectively creates a differently perceived shape despite the visible parts being the same as in the occlusion condition. The finding that notched circles and squares are more easily identified than occluded circles and squares is taken to reflect amodal completion.

961 **Dataset.** We generated samples that look like the stimuli used in [93]. We generated samples for the
 962 distractors ('unoccluded'), 'occlusion', and 'notched' conditions, with either the square occluding
 963 the circle or vice versa. The occluding shape is placed at a variety of degrees from the occluded
 964 shape. Each occluded image has a corresponding notched image (that is, using the same shape
 965 configurations) so that they can be directly compared. The unoccluded condition is generated by
 966 using a non-occluded sample in which the occluding shape is moved radially away from the occluded
 967 shape, maintaining the same orientation.

968 To align with human similarity judgments as measured in [93], a DNN should yield internal activations
 969 that result in higher similarity scores for distractor (unoccluded) versus occluded samples (where
 970 amodal completion generates representations of the full shape of the notched stimuli), compared to
 971 those for (distractor) unoccluded versus notched samples.

972 C.1.8 Non-accidental and Metric Properties for Geons and line stimuli

973 **Psychological Significance:** Human object recognition is highly sensitive to non-accidental properties
 974 (NAPs) of an object, that is, visual features that are invariant over rotations in depth. NAPs are
 975 hypothesized to be critical for representing object parts such as Geons [16]. For example, curvature
 976 (as opposed to a straight line) is a NAP because a curved object in the 3D world will project a curved
 977 image on a 2D retina when viewed from most orientations apart from rare "accidental" viewpoints,
 978 as when a curve projects a straight contour. NAPs are distinguished from metric properties (MP),
 979 features of objects that change continuously with variations over depth orientation when projected on
 980 the retina. For example, a curved object in the world will project different degrees of curvature on the
 981 retina depending on its orientation to the viewer. A variety of research highlights how human vision
 982 is more sensitive to changes in images that alter NAPs (e.g., a change from a curve to straight line)
 983 compared to MPs (e.g., changes in degree of curvature) [16]. [71] provided evidence that several
 984 DNNs are also more sensitive to image manipulations that alter NAPs, although the effects were most

985 pronounced in later layers of the networks whereas sensitivity to NAP is thought to occur relatively
986 early in human visual processing to encode object parts.

987 **Datasets.** We have included images of both **2D line segments** based on [71], and **3D Geon stimuli**
988 originally used in [68] and obtained from ³ to assess the degree in which DNNs are sensitive to NAP
989 vs MP changes. In the case of the Geon stimuli, we have provided a version with shade (as in [72]), a
990 version in which no shades are present and the outline is highlighted, and a version in which only
991 the silhouettes are shown, as one concern with the shaded version is that the similarity judgements
992 produced by DNNs may reflect differences in shades as opposed to shape. For each Geon or line
993 segment, a feature dimension (such as the curvature of a Geon) is altered from a singular value (e.g.
994 straight contour with 0 curvature) to two different values (e.g. slightly curved or very curved). The
995 ‘reference’ condition includes items with the intermediate feature value; in this example, the slightly
996 curved geon. The ‘MP change’ condition consists of items with a greater non-singular value; in
997 this case, the greater curvature geon. Finally, the ‘NAP change’ condition includes items with the
998 singular value; the straight contour geon from this example. A human-like similarity judgment would
999 correspond to higher similarity between the reference object to the MP variants than the NAP variants
1000 (that is, NAP changes are easier to discriminate). [71] provides a more detailed description of human
1001 performance through reaction times that can be directly compared to similarity judgments in DNNs
1002 (where higher reaction times correspond to lower similarity).

1003 C.2 Visual Illusions

1004 There is now a growing number of articles exploring various illusions in various different types of
1005 DNNs trained in different ways. Some of these highlight similarities with human perception (e.g.,
1006 [14, 103, 111]) while others report mixed or discrepant results (e.g., [55, 54, 110, 119]). For a review
1007 of various findings see [67]). The conditions in which DNNs ‘experience’ human-like illusions may
1008 provide insights into why humans experience these phenomena. For instance, [55] found that various
1009 brightness and color illusions can be induced in CNNs trained for image denoising, image deblurring,
1010 and computational color constancy. They argue that these illusions are a byproduct of biological
1011 processes designed to improve efficiency of low-level visual processes. In addition, Storrs et al.
1012 [103] found that unsupervised (as opposed to supervised) learning led DNNs to factorize images into
1013 encoding of reflectance and illumination that resulted in a human-like perceptual illusion of gloss.
1014 Here we consider several classic size, lightness, and orientation illusions.

1015 C.2.1 Müller-Lyer illusion

1016 **Psychological Significance.** The Müller-Lyer illusion is perhaps the most famous of all illusions.
1017 There is no agreed-upon explanation of the effect but the fact that it is observed across species,
1018 including fish [101], suggests that it reflects something basic about the architecture of the visual
1019 system rather than the training environment. Ward [110] reported that VGG-19 showed the illusion
1020 (to a rough approximation), although they only reported the effect at the final stage of the network. A
1021 similar result was reported by [118] who reported more robust effects in the higher levels of VGG-19
1022 and ResNet-101.

1023 **Dataset.** The Müller-Lyer illusion stimuli were generated in one of two ‘illusory’ configurations
1024 (with inward or outward ‘fins’) or in a ‘scrambled’ configuration. In the latter, the fins are arranged
1025 randomly in the canvas, separated from the line segment. In all three conditions, we vary the
1026 line length, the position of the line, and the angle of the fins. A method to test whether a DNN
1027 is susceptible to this illusion involves training a set of decoders to predict the line length in the
1028 scrambled condition set. These decoders are then tested on the illusory conditions. A human-like
1029 response would be evidenced by a consistent pattern of both overestimating the line length in the
1030 outward illusory condition and underestimating it in the inward illusory condition. Additionally, the
1031 illusory effect should be larger for more acute fin angles.

³<https://geon.usc.edu/ori>

1032 C.2.2 Ponzo illusion

1033 **Psychological Significance.** In this classic illusion, two identical horizontal lines cross a pair of
1034 converging lines, a configuration similar to railway tracks. In this configuration, the top line looks
1035 longer. The standard explanation is that the visual system assumes that the converging lines are
1036 receding in depth and that the upper horizontal line is further away. Given that the two lines project
1037 the same length on the retina, the visual system assumes that the upper line must be longer. That is,
1038 the illusion is a by-product of the visual system attempting to compute size constancy. Interestingly,
1039 there is good evidence that the Ponzo illusion [59], and related size illusions [102], alter the activation
1040 in V1, although this may reflect top-down activation from higher-level visual areas [31]. [110] failed
1041 to observe a Ponzo effect in VGG-19.

1042 **Dataset.** Two target lines (red and blue) are placed across a railway track pattern. In the illusory
1043 condition, the target lines have the same length (varying across samples). In the scrambled condition,
1044 the target lines have different length, are still placed horizontally one on top of the other, but all the
1045 other segments are randomly placed across the canvas. We include a third condition in which the
1046 railway track pattern is used with target lines which differ in length. The railway track pattern for the
1047 illusory and different lengths conditions is composed of converging segments (with a varying degree
1048 of convergence), and horizontal segments (randomly placed at different horizontal positions). For all
1049 three conditions, the user can specify the number of horizontal segments to use.

1050 The suggested way to test whether DNNs perceive the Ponzo Illusion consists of training a set of
1051 decoders on the scrambled condition, to predict either the length of the target lines or a function of the
1052 length (for example, the difference between the top and the bottom line lengths). Then the decoders
1053 can be tested on the illusory condition. The different lengths condition could be used as a further
1054 way of analysing the decoders response. To match human perception, a decoder should overestimate
1055 the length of the top target line (or underestimate the length of the bottom line, or output a positive
1056 difference in top minus bottom line length, depending on the training setup) in the illusory condition
1057 (where the two target lines have the same length).

1058 C.2.3 Ebbinghaus (or Titchener) illusion.

1059 **Psychological Significance:** In this classic illusion, the perceived size of a central circle is altered
1060 by the size of surrounding circles. There is evidence that the illusion distorts the perception of size
1061 but not action ([2, 114]; but see [48]) and there is evidence that this illusion is mediated by relatively
1062 low-level (preattentive) vision [29]. Again, there are different explanations for the phenomena [92].
1063 [110] failed to observe this effect in VGG-19.

1064 **Dataset.** A red target circle is surrounded by a fixed number of white circles (flankers) on a uniform
1065 background. In the two illusory conditions ('big' and 'small' flankers) the flankers surround the target
1066 circle, and they all have the same size within each sample. In the scrambled condition the target circle
1067 is placed in the center, but white circles with random sizes are randomly placed on the canvas. Across
1068 illusory samples, we varied the radii of the flankers, the radius of the target circle, the displacement
1069 of the flankers around the target. To measure illusory effects in DNNs, decoders can be trained on
1070 estimating the circle size or radius in the scrambled condition, and tested on the big/small flankers
1071 condition. Human-like perception should induce overestimating in the small flankers condition and
1072 under-estimating in the big flankers condition (see example in Figure 2).

1073 C.2.4 Jastrow Illusion

1074 **Psychological Significance:** In the Jastrow Illusion [66], two identical-sized curved segments are
1075 perceived as different sizes when one is placed above the other in certain configurations. There are
1076 multiple explanations for the phenomenon, but perhaps the simplest explanation is that it is a form of
1077 a contrast effect. The length of the concave edge of the upper object in a Jastrow configuration is
1078 much shorter than the convex edge of the bottom object, and this contrast drives the perception of
1079 size when the edges are closely aligned [94]. Rhesus monkeys do not appear to be affected [3], nor

1080 do humans when assessed on grasping behavior [82]. As far as we are aware, no one has reported
1081 whether DNNs show a similar pattern.

1082 **Dataset.** We used a red and a blue arc shape, either one on top of the other at the centre of the
1083 canvas ('illusory' and 'different lengths' conditions) or randomly placed in the canvas with a random
1084 orientation ('scrambled' condition). In the scrambled and different lengths conditions the two shapes
1085 have different sizes. The size is the same (thus eliciting the illusion) in the illusory condition. To
1086 estimate DNNs susceptibility to the illusion the same approach as the Ponzo Illusion can be used.

1087 C.2.5 Tilt illusion

1088 **Psychological Significance:** In the tilt illusion, a central grating's orientation is perceived as being
1089 repulsed from or attracted to the orientation of a surrounding grating. A wide variety of mechanistic
1090 accounts of the illusion have been proposed (for review see [33]), and it is argued to be an adaptive
1091 feature rather than a bug of a visual system optimized for contour detection [98]. There is evidence
1092 that the illusion reflects processes in V1 [100]. Linsley, et al. [73] reported that a recurrent DNN
1093 optimized for contour detection produces a tilt illusion.

1094 **Dataset.** We provide one illusory condition, in which an oriented grating pattern is presented within
1095 a circular mask ('center grating') and a differently oriented grating is placed as the background
1096 ('context' grating); and two non-illusory conditions: one in which the background is uniformly
1097 colored and only a center mask contains the oriented grating pattern; and vice versa. The samples
1098 are varied in their orientation and spatial frequency of the gratings, and in the size of the central
1099 grating. Our suggested approach to test whether a DNN perceives the tilt illusion is to train a decoder
1100 to estimate the orientation of the center grating, and test it on the illusory condition to check whether
1101 the presence of a context affects performance. In particular, the decoder should present the largest
1102 repulsive bias at around 20° and an attractive bias at around 70° - 80° . Plus, the attractive effect should
1103 be much smaller than the repelling effect, and larger for matching center-surround gratings spatial
1104 frequencies. [113].

1105 C.2.6 Lightness Illusions

1106 Lightness refers to our perception of the reflective surface of an object (a stable property of an object)
1107 whereas brightness is a measure of the amount of light reflected from an object, something that is
1108 affected by both reflectance as well as the lighting source. We include two famous illusions related to
1109 lightness: the Lightness Contrast Illusion and the Adelson Checker Shadow Illusion. However, to
1110 facilitate testing for these and other lightness-related effects, we created an additional dataset called
1111 '**Grayscale shapes**'. The purpose of this dataset is not to elicit any illusion in humans or in DNNs
1112 but to train a network (or, with our suggested method, a decoder attached to a network) to output the
1113 grayscale value of a target pixel.

1114 **Grayscale shapes Dataset.** Each image is composed of 20 overlapping items amongst the following
1115 types of shapes (circle, circle sector, circle segment, ellipse, rectangle with straight and rounded
1116 corners, heptagon, irregular polygon composed of a random number of edges from 3 to 10). Position,
1117 dimension, orientation, and grayscale colour value are randomized for each shape. We place 20 items
1118 to be sure that most space in the canvas is filled by an item, but that only a few of them are fully
1119 visible. This results in a chaotic canvas with many different shapes with varying grayscale colours
1120 but with coherent patterns (as opposed to, for example, having each pixel of a different random
1121 grayscale value). In order to target a specific pixel to be predicted by the decoder, a small white
1122 vertical arrow (the 'marker') of fixed size is placed randomly on the canvas. The arrow points to the
1123 pixel whose value can be used for prediction. Notice that while the images are commonly normalized
1124 from -1 to 1 before being fed into the network, the targeted pixel value to predict is in the 0-255 range.
1125 Once a trained decoder reaches the desired level of accuracy, it can be tested on other configurations
1126 by simply adding the white arrow 'marker' into any image. We call this network with the decoder
1127 attached the **color-picker**. We can then test whether an illusory configuration impacts performance
1128 of the color picker by placing a white arrow marker at several points of the illusory image and check

1129 whether the output is biased in a human-like fashion. This is the approach we use in the Lightness
1130 Contrast Effect and Adelson Checker Shadow Illusion.

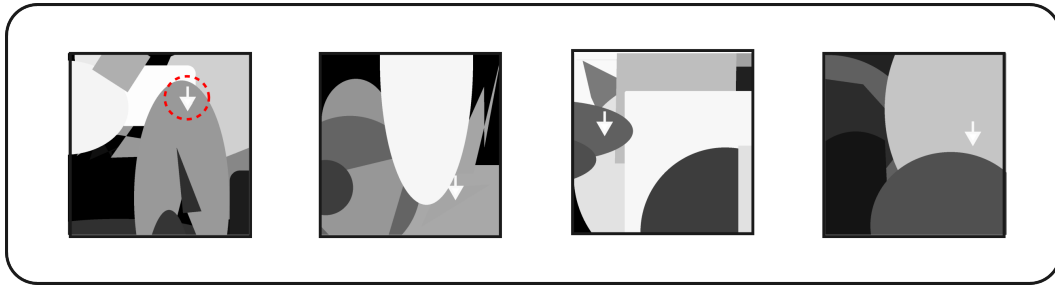


Figure 10: Samples of five images from the grayscale dataset, used to train a color-picker, as detailed in the text.

1130

1131 C.2.7 Lightness Contrast effects

1132 **Psychological Significance:** In the Lightness Contrast effects our perception of two identical central
1133 gray patches is altered by their surround, such that a patch surrounded by a dark background is
1134 perceived as lighter, and the patch surrounded by a light background is perceived as darker. The
1135 standard explanation of this is that lightness perception is the product of the relative brightness
1136 of surfaces across a boundary given that this ratio will remain constant regardless of the general
1137 illumination, allowing lightness (and color) constancy. However, in the lightness contrast context,
1138 the mechanism designed to produce lightness constancy results in the central grey squares being
1139 perceived differently. These computations are thought to occur in the primary visual cortex [76].
1140 Some DNNs can achieve color constancy [46] and other forms of constancy [103] under some
1141 conditions, although there remain questions as to whether this is achieved in a human-like manner.
1142 There are also computational theories of the lightness contrast effect [57], but we are not aware of
1143 any demonstrations that DNNs support this effect.

1144 **Dataset.** The dataset consists of the standard Lightness Contrast configuration: square within a
1145 uniform canvas of different grayscale values. The user can specify the grayscale value of the center
1146 square, which is kept fixed, while the value of the background is varied. Importantly, each sample
1147 is replicated many times with the white arrow marker placed at different locations in the canvas. A
1148 color picker network can then be queried for the grayscale value at different locations in order to
1149 measure whether the perceived value of the central square is affected by its surroundings.

1150 C.2.8 Adelson checker shadow illusion

1151 **Psychological Significance:** In this classic illusion [1], two squares on a checkerboard are perceived
1152 to have different reflectance due to one being in shadow and the other in light, despite being the same
1153 brightness. This phenomenon is explained with the ability of the human brain to perceive reflectance
1154 of a surface as invariant under variation of brightness. In this illusion, the patches inside and outside
1155 the shadow reflect the same brightness, and accordingly, the visual system assumes the patch in the
1156 shadow must be lighter.

1157 **Dataset.** This dataset simply consists of the Adelson Checker Shadow illusory image replicated
1158 many times, grayscaled, with a white arrow systematically placed at different locations of the canvas,
1159 covering the whole checkerboard. A color-picker network (that is a decoder trained to predict the
1160 value of a marked pixel, e.g. trained on the Grayscale shape dataset, see Section C.2.6) is queried at
1161 all locations. Critically, the color-picker will show illusory perception if the pixels in the two target
1162 patches are seen as two different colours. In particular the pixels of the unshaded patch should be
1163 seen as darker than the shaded patch by the network.

1164 C.2.9 Thatcher Illusion

1165 In this illusion the eyes and mouth of upright faces are inverted to produce a grotesque image of
1166 a person. The distinction between a normal face and a distorted face is highly salient. However,
1167 when faces are upside down, the distortions are much less salient. This effect was originally reported
1168 on images of Margret Thatcher, thus the name of the illusion. This effect is sometimes claimed
1169 to be more dramatic for faces compared to other categories, lending support to the hypothesis that
1170 face processing is special [25]. Other researchers claim that similar effects are observed for other
1171 types of objects [115]. Jacobs et al. [65] demonstrated that CNNs trained to identify faces exhibit a
1172 Thatcher-like effect, though they did not assess whether this effect extends to other object categories
1173 [65]. The extent to which this effect is specific to faces or can be generalized to non-face objects
1174 remains a subject of debate (e.g., see [116, 15, 34]). To test DNN sensitivity to the Thatcher Effect for
1175 both faces and non-face dataset, we provide both a Thatcherized dataset of faces and a Thatcherized
1176 dataset of words (in which individual letters are rotated).

1177 **Face Dataset.** We provide a small dataset of celebrity faces using a subset of CelebA⁴, but the user
1178 can specify any folder containing images of faces. Each image is resized according to parameters
1179 specified by the user and then reoriented into both an upright and a 180-degree inverted configuration.
1180 Furthermore, it is either 'Thatcherized' or remains unaltered. To 'Thatcherize' an image we compute
1181 landmarks of the eyes and mouth, compute the bounding rectangle for each, and rotate them around
1182 their centre of mass. Blurring on the edge is applied to minimize artefacts. To assess the susceptibility
1183 of DNNs to the Thatcher effect in faces, we propose a similarity judgment analysis. This involves
1184 comparing the perceived similarity between each upright face and its Thatcherized counterpart, as
1185 well as each inverted face with its Thatcherized version. To align with human perception, the latter
1186 comparison is expected to yield a higher similarity score than the former.

1187 **Word Dataset.** We employ a collection of 1000 English words or artificially generated sequences of
1188 random letters. All entries are uniformly presented in uppercase, covering a range from 3 to 8 letters
1189 in length. Following [116], to simulate the Thatcher Effect for words, we rotate one or more letters
1190 by 180 degrees. To increase variability, each word is displayed in one of ten different fonts, with
1191 variable font sizes, and includes jitter for each letter. The configurable parameters include the number
1192 of words, the exact or range of letter counts per word, the number or range of letters to be rotated, the
1193 font size, the level of jitter, and whether to use random strings or English words.

1194 C.3 Shape and object recognition

1195 A key feature of human vision is that we identify objects largely based on their shape. For example,
1196 we can easily identify line drawings of objects with no colour and texture [19]. To measure shape
1197 bias in DNNs there is now a benchmark that tests models on "style transfer" images composed of the
1198 shape of one category and the texture of another [51]. Many DNNs, including DNNs that perform at
1199 the top of the leaderboard on Brain-Score, rely primarily on non-shape features, as they classify the
1200 images based on their texture rather than shape. More recent DNNs trained on much larger datasets
1201 have started to show a more human-like shape bias [37], but there are many additional attributes of
1202 human shape perception that need to be accounted for by any DNN model of human vision.

1203 C.3.1 Identifying line drawings, dotted line drawings, silhouettes, and image segments

1204 **Psychological Significance:** Humans can often identify line drawings of objects as quickly and
1205 accurately as photographs, highlighting the importance of shape for object identification [19]. Inter-
1206 estingly, a child who had never previously been exposed to line drawings can readily identify them,
1207 showing that there is no need to be trained on line drawings to identify them [62]. By contrast, DNNs
1208 need to be trained on line drawings in order to recognize them at human levels [99]. Similarly, humans
1209 can easily identify silhouettes of objects, whereas DNNs again struggle (although interestingly, they
1210 do better with silhouettes compared to line drawings; [10]).

⁴<https://mmlab.ie.cuhk.edu.hk/projects/CelebA.html>

1211 In addition, by exploiting the Gestalt principle of good continuation, humans can recognise line
1212 drawings with modified local features when the global shape is left intact. In our experiments, we
1213 modified line drawings in three ways: by replacing the continuous line with dots; by replacing
1214 continuous lines with segments; and by texturizing them. These images are easily identifiable, and
1215 accordingly, DNNs should be able to identify these out-of-distribution images.

1216 **Line drawings dataset.** We use the line-drawing stimuli from [10], consisting of 36 classes from
1217 ImageNet (one line-drawing per class). The line-drawings are white stroke on a uniform canvas
1218 (black by default). We used this dataset to build the **Dotted line drawings and Image Segments**
1219 **datasets.** In the former case, the user can specify the dot size and the distance between dots. In the
1220 latter case, we have generated complementary images that have complementary segments removed
1221 (see Figure 11). That is, each line segment in one image is absent in the other, and together the image
1222 is complete. These stimuli are generated by overlapping a grid on the line drawing and deleting
1223 complementary sections. The user can specify the grid orientation, the distance between each grid
1224 row and column, and thickness of each cell. Participants find these images trivial to identify, and
1225 accordingly, DNNs should also. Importantly, humans find complementary images like these hard
1226 to distinguish, and indeed, complementary images produce equivalent priming to repeated images,
1227 highlighting how the visual system treats them as equivalent [17]. This would also be the case if
1228 complementary dots were removed for the dotted line drawings. Thus a second approach to compare
1229 humans to DNNs is through a similarity judgment analysis across complementary images, which
1230 should return very high similarity value in some hidden layers of DNN.

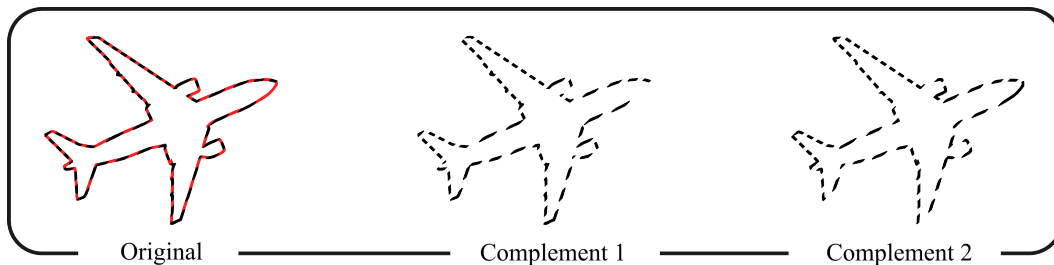


Figure 11: Example of the complementary segment dataset. Each linedrawing results in two images with complementary segments removed. The resulting samples are very difficult to discriminate for humans [17].

1231 There are many datasets of line drawings available, and the user can specify any folder containing line
1232 drawings, to generate a dotted/segmented dataset. The line drawings are expected to be composed of a
1233 black stroke on a white background, but can otherwise be of any shape, and the line drawing folder is
1234 expected to contain sub-folders for each class (e.g. ‘airplane’) which can contain multiple line drawing
1235 instances for that class (this follows the standard structure used in many deep learning libraries for
1236 image classification, e.g. the ImageFolder dataset in PyTorch). Our script will automatically convert
1237 the image to a white stroke over a uniform background (black by default).

1238 **Silhouette dataset.** For the Silhouette dataset, we use samples from [8] (9 classes from ImageNet,
1239 each containing 40 samples). As before, the user can specify any folder containing silhouettes. Altern-
1240 atively, the user can also specify a folder containing line-drawings (following the same constraints
1241 as above), which will be converted into silhouettes. Again, humans find these images easy to identify,
1242 so models should as well.

1243 C.3.2 Identifying familiar and unfamiliar images defined by texture boundaries

1244 **Psychological Significance:** The human visual system can group elements in scenes based on texture,
1245 with texture regions defined by the similarity of their elements. This is an example of a Gestalt
1246 principle (Similarity) contributing to object recognition [13]. One way this manifests is through the
1247 ability to identify familiar (**texturized objects**) and perceive unfamiliar (**texturized unfamiliar**) by
1248 their texture.

1249 **Datasets: Familiar and Unfamiliar Texturized Objects.** We provide a dataset of familiar texturized
1250 objects by using line drawings from [10] as base items. For unfamiliar shapes, we generated
1251 silhouettes of blob-like objects. For the pre-generated datasets, the texturization consists of masking
1252 the internal contour of a line drawing/silhouette with a pattern of a repeated character with a
1253 randomized font size, rotated by a random degree. The character is randomly selected from letters,
1254 digits, or punctuation, and we kept the background uniformly colored. When generating images, the
1255 user can also specify the texturization of the background as well, although we have found that doing
1256 so will turn object recognition from trivial to challenging, depending on the selected character.

1257 The same approach is used for the unfamiliar shapes. In this case, the user can specify the number of
1258 blobs to generate and texturize. For both familiar and unfamiliar datasets, the user can specify how
1259 many texturization samples to generate for each input image.

1260 To measure alignment with human visual perception, the different datasets require a different approach.
1261 For familiar shapes, DNNs can be tested by simply assessing classification accuracy. For unfamiliar
1262 shapes, a similarity analysis can be carried out. For example, a DNN should find a blob and its
1263 texturized counterpart more similar than a blob and a differently texturized blob (see example in
1264 Figure 2).

1265 C.3.3 Identifying Embedded Shapes

1266 **Psychological Significance:** The Embedded Figures Test (EFT, [36]) is a widely utilized tool in
1267 research exploring individual differences in perception, with a particular emphasis on studies of
1268 autism spectrum disorder, and as a measure of local versus global perceptual style [83]. Subsequently,
1269 [36] developed a set of stimuli in which several Gestalt grouping principles were manipulated in
1270 order to create increasingly difficult matching to sample tasks. They found that the principle of
1271 good continuation (operationalized in terms of the number of continued lines from the original
1272 shape) impacted performance the most. Each target shape was integrated into four distinct contexts,
1273 each exhibiting a progressive increase in the number of lines extending from the target shape into
1274 its surroundings. The higher the number of lines extending the shape, the lower the performance,
1275 highlighting human susceptibility to camouflage and the role of Gestalt organisation principles in
1276 camouflage.

1277 **Dataset.** We used the dataset from [36] who developed simple stimuli in which background lines
1278 camouflaged geometric shapes to various extent (Figure 12). Importantly, different embeddings have
1279 different levels of continued lines from the original shape, which strongly affects human performance.
1280 Furthermore, we developed our version by generating 5 irregular polygons, embedding them in a
1281 set of lines, some random and some extending directly from the polygon’s edges (similarly to the
1282 original dataset). Many camouflaging samples can then be procedurally generated from each polygon.
1283 Training decoders to classify the simple geometric forms provides one way to assess the impact
1284 of embedding shapes on DNNs. Decoders would be trained on simple shapes (either our irregular
1285 polygons or the original shapes from [36]) and would then be tested on the embedded version. A
1286 DNN with a human-like perceptual system should show reduced ability to identify the shapes, with
1287 the level of impairment being a function of the amount of lines originating from the polygon (as
1288 in [36]). Notice in this case, human alignment requires a degradation of performance after image
1289 alteration.

1290 C.3.4 Sensitivity to Global Shapes

1291 **Psychological Significance:** Human object recognition relies more heavily on global shape repre-
1292 sentations than on local features, whereas there is evidence that DNNs rely more heavily on local
1293 features [10], even when trained to have a shape bias [8]. In [8], humans and DNNs were presented
1294 with silhouette stimuli in their normal format, fragmented, or in a ‘Frankenstein’ format where most
1295 of the local features are preserved but the overall configuration of the image was distorted. That is,
1296 the authors modified global shape while maintaining most of the local features. In particular, the
1297 ‘fragmented’ condition (see Figure 13) divides the shape into two distinct, yet adjacent, entities while

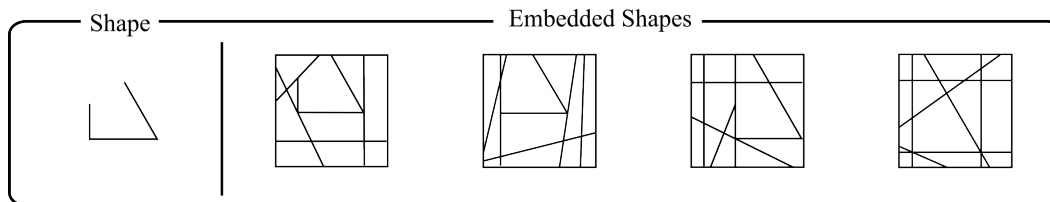


Figure 12: Illustration of one set of items from the embedded shape dataset. These are the stimuli recreated based on [106]. A basic shape is camouflaged using a variety of horizontal and vertical line, and extending the segments composing the shape. We also provide a variation of this dataset in which, given a set of polygons, camouflaged versions are procedurally generated.

1298 preserving the local characteristics of the object. The “Frankenstein” scenario involves adjusting the
 1299 upper section back into alignment with the lower half, so that the bottom and top halves are mirror
 1300 reversed. This method keeps the object intact as a single entity. Human performance was much
 1301 reduced in both the Fragmented and Frankenstein conditions, but DNNs performed similarly in the
 1302 Whole and Frankenstein conditions, highlighting the importance of the local features and the lack of
 1303 weighting for more global features in driving their performance. Attempts to train networks to focus
 1304 on the more global aspects of the images failed.

1305 **Dataset.** We provide both the dataset extracted directly from [8] and a a version in which the
 1306 fragmented and Frankenstein versions are automatically generated from any silhouettes or line
 1307 drawing samples. The [8] dataset contains 9 classes from ImageNet, each containing 40 samples. A
 1308 network with visual capabilities aligned with a humans’ should suffer from performance degradation
 1309 in both fragmented and Frankenstein condition, which can be measured through classification accuracy
 1310 (as usual, with a network pretrained on ImageNet or some other image dataset).

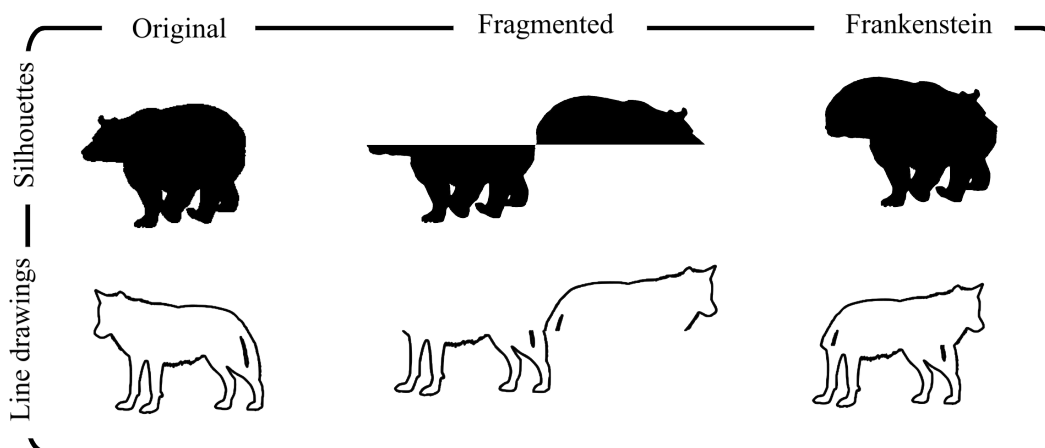


Figure 13: Example of a stimulus and its transformed version, following [8].

1311 C.3.5 Invariance to Object Transformation

1312 **Psychological Significance.** Humans possess the remarkable ability to recognize objects despite
 1313 the different retinal images the objects project depending on changes in size, orientation, lighting,
 1314 and placement [104]. This is performed on-line: an object, once seen in a new or altered form, can
 1315 typically be recognized instantly in subsequent exposures at different angles, without further training.
 1316 This capability applies irrespective of the object’s familiarity [24, 23]. Previous work [20, 21] found
 1317 that none of the 7 tested classic visual DNNs possess on-line invariance architecturally (that is:
 1318 training an object at a viewpoint would not automatically support object recognition at a different
 1319 viewpoint, even for simple affine transformations such as translation). However, this ability could be
 1320 induced by pre-training the model on the specific transformation of interest, and this would transfer

1321 to unfamiliar classes. For example, a network pretrained to classify images in which the objects are
1322 randomly rotated, develop on-line invariance to rotation, even for novel objects and novel classes.
1323 This partially extends even in transformation of viewpoint, in which the object is rotated in depth.

1324 **Datasets: 2D Affine Transformations and Viewpoint Transformations.** We provide a separate
1325 dataset for affine transformations (rotation in picture plane, translation, scale, and shear) and viewpoint
1326 invariance (rotation in depth). The configurable parameters allow for a fine-grained analysis of the
1327 effect of each transformation. For each transformation dimension, the user can chose one or multiple
1328 ranges of training and testing. As in the previous datasets, the user can specify any folder containing
1329 line-drawings or silhouettes (or any image with a clear contour on a white background). For the
1330 viewpoint invariance dataset, we use the ETH-80 dataset [32]⁵, which contains 8 categories (apples,
1331 cars, cows, cups, dogs, horses, pears and tomatoes), each consisting of 10 object instances, and each
1332 object captured from 41 different viewpoints. For the pre-generated dataset, we avoided including
1333 views straight from the top. The configurable parameters allow the user to generated dataset only
1334 within a specific azimuth and inclination range.

1335 For both the 2D transformation and viewpoint datasets, there are several ways to test whether a DNNs
1336 possess online invariance to transformation. First, a DNNs (not necessarily pretrained) could be
1337 trained on un-transformed images (e.g. with the object always in the center, unrotated, unscaled, or
1338 from a standard viewpoint). It could then be tested on various transformations of these objects. This
1339 could be used to establish whether the network is architecturally invariant to some transformations.
1340 Several pre-training steps could be used to test how the training environment affects performance.
1341 Another approach avoids training on the target classes (either because we want to test a pretrained
1342 network without altering its weights, or because we want to test the network on unfamiliar classes):
1343 a similarity judgments analysis is performed on transformed versions of the same object, and is
1344 compared to the similarity of different objects. A human-like DNN will have internal activations
1345 that are more similar for same objects across transformations, than for different objects. This is the
1346 approach used in [20] and [21].

1347 C.3.6 Same/Different Task

1348 **Psychological Significance.** Human shape representations not only support object recognition but
1349 also a wide variety of additional functions, including visual reasoning. Perhaps the simplest form of
1350 visual reasoning is tested in the same/different task – judging whether two shapes are identical apart
1351 from their spatial location. Although DNNs can solve the same/different task when training and test
1352 images are highly similar to one another, performance drops when training/test images are dissimilar
1353 [89]. By contrast, humans can make same/different judgements for any visual patterns as long as they
1354 are perceptible.

1355 **Same/Different Dataset.** The dataset was extracted from [89]. This dataset is composed of 10
1356 conditions. Each image consists of two items placed randomly on the canvas. The two items can be
1357 either the same shape or a different shape and cannot overlap. Each condition consists of a different
1358 type of item used. By default, the items are composed of white strokes with no fill on a black
1359 background. See Figure 14 for a summary of all the conditions.

1360 The suggested testing methodology for this condition is slightly different than all other methods, and
1361 consists of training a DNN (not necessarily pre-trained) on a subset of conditions, and testing it on a
1362 different subbset (as in [89]).

⁵<https://github.com/chenchkx/ETH-80/tree/master>

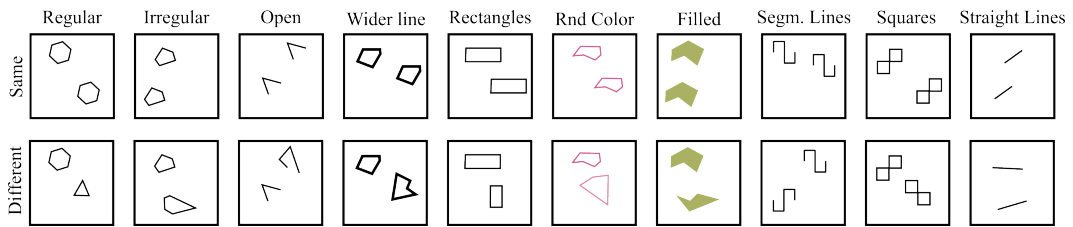


Figure 14: Illustration of the ten conditions used for the Same/Different task. Two items can be either the same or different shapes up to translation. For the ‘straight lines’ condition, the “same/different” dimension considered is the line orientation (with length kept fixed).