

The contrasting shape representations that support object recognition in humans and CNNs

Gaurav Malhotra,^{1*} Marin Dujmovic¹, John Hummel,² Jeff Bowers¹

¹School of Psychological Sciences, University of Bristol, Bristol, UK

²Department of Psychology, University of Illinois Urbana-Champaign, Champaign, USA

* gaurav.malhotra@bristol.ac.uk.

1 Abstract

The success of Convolutional Neural Networks (CNNs) in classifying objects has led to a surge of interest in using these systems to understand human vision. Recent studies have argued that when CNNs are trained in the correct learning environment, they can emulate a key property of human vision – learning to classify objects based on their shape. While showing a shape-bias is indeed a desirable property for any model of human object recognition, it is unclear whether the resulting shape representations learned by these networks are human-like. We explored this question in the context of a well-known observation from psychology showing that humans encode the shape of objects in terms of relations between object features. To check whether this is also true for the representations of CNNs, we ran a series of simulations where we trained CNNs on datasets of novel shapes and tested them on a set of controlled deformations of these shapes. We found that CNNs do not show any enhanced sensitivity to deformations which alter relations between features, even when explicitly trained on such deformations. This behaviour contrasted with human participants in previous studies as well as in a new experiment. We argue that these results are a consequence of a fundamental difference between how humans and CNNs learn to recognise objects: while CNNs select features that allow them to optimally classify the proximal stimulus, humans select features that they infer to be properties of the distal stimulus. This makes human representations more generalisable to novel contexts and tasks.

Author summary

The human visual system is highly adept at recognising and reasoning about objects under a wide variety of viewing conditions. Research in vision sciences has shown that humans largely rely on shape to recognise objects, but extracting object shape from a retinal image is a challenge (indeed, it is an ill-posed problem), and until recently, models of vision have been poor at classifying naturalistic images of objects. This has changed with the development of Convolutional Neural Networks (CNNs) that perform as well as humans on certain recognition tasks. Thus it makes sense to ask if we can gain insight into how humans extract shape by studying CNNs. Here, we show that there is a fundamental difference between the two systems. While humans pay particular attention to relations between an object's parts and features, we show that the shape extracted by CNNs is akin to a template, where all diagnostic features are equally important. Thus, despite their success, these results show that CNNs are currently not good models of human object recognition and machine learning models must bridge this gap in order to capture the robustness and generalisability of human vision.

Introduction

A great deal of research into human vision is driven by the observation that we do not perceive the world as it really is. Instead, visual perception is biased. This is particularly true about our perception of objects. We prefer to group objects based on certain Gestalt principles (a bias to look for proximity, similarity, closure and continuity [9]), we prefer to view objects from certain perspectives (a bias for canonical-perspectives [37]) and we prefer to categorise objects based on certain features (a bias for shape [31, 6]).

There are two possible explanations on the origin of these biases. The first view, which we shall call the *optimisation approach*, proposes that these biases are an internalisation of the biases present in the environment relevant to a particular task. According to this view, humans prefer to view objects from a canonical perspective because these perspectives are more frequent in the visual environment and they prefer to classify objects based on shape because shape is more diagnostic during object classification. In other words, biases are a consequence of performing statistical learning with the goal of optimising behaviour on a particular task. As a person learns a task, they acquire the statistical dependencies present within the task-environment and starts mirroring these environmental biases.

The second view, which we shall call the *heuristic approach*, proposes that biases arise because of the manner in which the visual system transforms visual inputs. The retinal image of an object contains a vast amount of information in the form of luminance values at various wavelengths at each location on the retina. However, the information relevant to successful interaction with the environment, such as the identities, locations and trajectories of objects in the world (that is, information about the *distal stimulus*), is nowhere explicit in the retinal image (the *proximal stimulus*). According to the heuristic view, the visual system transforms the retinal image to create a representation of the distal stimulus. It is this transformation from proximal to distal stimulus that is the source of biases.

Of course, the simple act of transforming one representation to another should not necessarily lead to biases. But, in this case, mapping the retinal image to the distal stimulus is an ill-posed problem: there is not enough information in the proximal stimulus to unambiguously recover the properties of the distal stimulus [39, 36]. To overcome this problem, the visual system makes assumptions (i.e., employs heuristics) to determine which properties of the proximal stimulus are used to build distal representations [47, 28, 34, 40]. A striking example of such assumptions is the Kanizsa triangle [26], where the visual system encodes the multiple collinearities of edges present in the proximal image and uses these to build contours of a triangle even though these contours do not exist in the retinal image. On this view, the goal of the visual system is *not* to optimise performance on a certain task, but to arrive at a veridical representation of the (distal) cause of the retinal image. The advantage of developing these representations is that they are relevant for broad range of tasks – the same representation of an object can be used for recognition and visual reasoning [21].

Psychological experiments have provided an equivocal support for both views. For example, the optimisation view is indirectly supported by studies that show that human biases, such as the shape-bias, increases with age [46], while the heuristic view is indirectly supported by the early observations of Gestalt psychologists as well as more recent studies that show that visual perceptions are more than the sum of their sensory inputs [42, 35]. However, psychological experiments provide only an indirect method of testing these views. A more direct test can be performed by constructing a computational model that learns to identify objects and test whether this model shows similar biases to humans. For a long time, it was not possible to create such a model as most models could not match human performance on object recognition tasks. This has recently changed with the advent of Convolutional Neural Networks (CNNs).

CNNs are machine learning models that can replicate (and sometimes even exceed) human performance on some object recognition and localisation tasks [32]. Importantly, CNNs learn to recognise objects by optimising the mappings between the proximal stimulus and the object categories themselves [12]. As a consequence, the learned representations that support object recognition are specialized for image classification. There is no pressure to learn representations of objects than can perform a range of tasks, let alone learn distal representations of objects. As such, CNNs provide a concrete model to test the optimisation view. If human perceptual biases are acquired through internalising the statistics of the environment on a particular task, then training CNNs to perform classification on ecologically realistic datasets should lead to perceptual biases similar to the ones observed for humans.

A number of recent studies have tested this idea by looking at *shape-bias* in CNNs. Psychological experiments have repeatedly shown that humans categorise objects primarily based on their shape, rather than other properties such as colour, size or texture [31, 6]. One manifestation of this bias is that we can identify most objects from line drawings as quickly and accurately as we can identify them from full-color photographs [6] and we can do this even if we have no previous experience with line drawings [16]. To test whether learning to perform classification on a naturalistic images can lead to a shape-bias, Geirhos et al. [11] trained CNNs on a large collection of naturalistic images (**ImageNet**) and tested them on a cue-conflict task that combined the shape of one category (say a cat) with the texture of another category (say, an elephant). They observed that instead of showing a shape-bias, CNNs preferred to classify these conflicting images based on texture rather than shape (classifying the image as an elephant), while human participants did the opposite. In another study, Malhotra et al. [33] tested human participants and CNNs on images that simultaneously contained multiple predictive features, including global shape, but also local features such as large coloured patches or segments. While CNNs classified these images based on the most predictive features, human participants ignored these features and preferred to classify objects based on global shape, even though this was not the optimal policy in the task. These findings seem to challenge the optimisation view as they show (a) that training a network to optimise performance on a database of naturalistic images does not necessarily lead to a shape-bias, and (b) humans seem to have a shape-bias even in environments where the optimal policy is to learn based on non-shape features.

However, recently it has been argued that CNNs can also be trained to infer an object's shape given the right type of training. For example, Geirhos et al. [11] trained standard CNNs on Style-Transfer image dataset that mixes the shape of images from one class with the texture from other classes so that only shape was diagnostic of category. CNNs trained on this dataset learned to classify objects by shape. In another study, Feinman and Lake [10] found CNNs were capable of learning a shape-bias based on a small set of images, as long as the training data was carefully controlled. Similarly, Hermann et al. [15] showed that more psychologically plausible forms of data augmentation, namely the introduction of color distortion, noise, and blur to input images, make standard CNNs rely more on shape when classifying images. Indeed, the authors found that data augmentation was more effective in inducing a shape bias than modifying the learning algorithms or architectures of networks, and concluded: "Our results indicate that apparent differences in the way humans and ImageNet-trained CNNs process images may arise not primarily from differences in their internal workings, but from differences in the data that they see".

These results raise the possibility that human biases are indeed a consequence of learning the statistical properties of the environment, à la CNNs, rather than developing representations of distal objects. But studies so far have focused on judging whether or not CNNs are able to develop a shape-bias, rather than examining the type of shape

representations they acquire. If humans and CNNs indeed acquire a shape-bias through a similar process of statistical optimisation, then CNNs should not only show a shape-bias, but also develop shape representations that are similar to human shape representations. The data on this is currently controversial. Some studies have found that CNNs are sensitive to small changes in object shapes and this sensitivity correlated with humans [30], while others have found that changing some local features of objects severely impacts object recognition in CNNs but has no impact for humans [1].

A key finding about human shape representations is that humans do not give equal weight to all shape-related features. For example, it has been shown that human participants are more sensitive to distortions of shape that change relations between parts of objects than distortions that preserve these relations [2, 24]. These observations have typically been taken to support a heuristic view according to which relations present in the proximal images are used to build distal representations of objects [18]. The question we ask is whether CNNs trained to classify objects learn to encode these relational features of shape. If they do, it would suggest that CNNs and humans indeed learn similar shape representations [15] and that shape-biases in object recognition are the product of optimising performance on object classification. But if not, it would suggest that these biases are best characterized as heuristics designed to build distal representations of shape. We present two experiments below designed to tease apart the shape representations in humans and CNNs. Our results suggest that even though CNNs are able to categorise objects based on shape, they do this on the basis of qualitatively different shape representations to humans. In both experiments human behaviour is better explained by the heuristic approach than an optimisation approach.

The rest of the paper is divided into three sections. In the first section, we focus on objects that consist of multiple parts and, in the second section, on objects that consist of a single part. The deformations required to infer the shape representations of these two types of objects are different, but related. Therefore, we begin each section by describing these deformations and how these deformations are predicted to affect shape representations under the two (optimisation and heuristic) views. We then present results of experiments where humans and CNNs were trained on the same set of shapes and then presented these deformations. In the final section, we discuss how our findings pose a challenge for developing models of human vision.

Experiment 1: multi-part objects

Proximal and distal encodings of multi-part objects

What sort of deformations of the proximal stimulus should allow us to contrast the optimisation and heuristic approaches? Specific hypotheses can be derived from the structural description theory [2], which assumes we build representations of distal stimuli during the process of identifying objects. On this theory, objects are represented as collections of convex parts in specific categorical spatial relations. For example, consider two objects – a bucket and a mug – both of which consist of the same parts: a curved cylinder (the handle) and a truncated cone (the body). The encoding of objects through parts and relations between parts makes it possible to support a range of visual skills. For example, it is possible to appreciate the similarity between a mug and a bucket because they both contain the same parts (curved cylinder and truncated cone) as well as their differences (the different relations between the object parts). That is, the representational scheme supports visual reasoning. In addition, the parts themselves are coded so that they can be identified from a wide range of viewing conditions (e.g., invariance to scale, translation and viewing angle, as well as robustness to occlusion), allowing objects to be classified from novel poses and under degraded conditions.

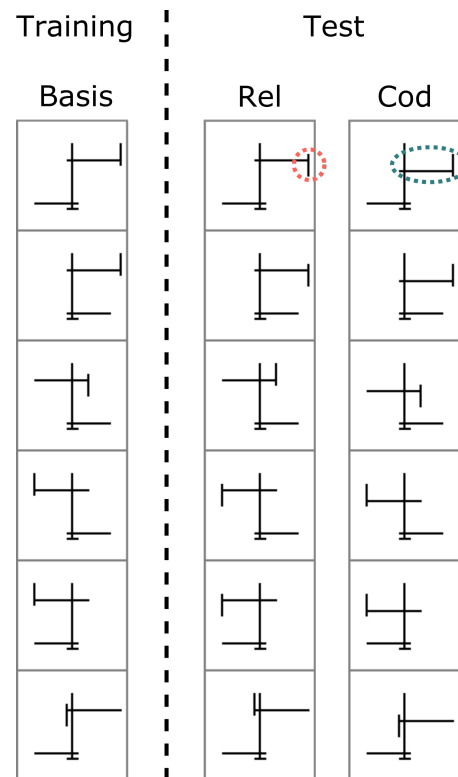


Fig 1. Stimuli used by Hummel and Stankiewicz [24]. The first column shows a set of six (*Basis*) shapes that participants were trained to recognise. Participants were then tested on shapes in the second and third columns, which were generated by deforming the *Basis* shape in the corresponding row. In the second column (*Rel* deformation) a shape is generated by changing one categorical relation (highlighted in red circle). In the third column (*Cod* deformation) all categorical relations are preserved but coordinates of some elements are shifted (highlighted in blue ellipse).

Note that the reliance on categorical relations to build up distal representations of multi-part objects is a built-in assumption of the model (one of the model's heuristics), and it leads to the first hypothesis we test, namely that image deformations that change a categorical relation between an object's parts should have a larger impact on the object's representation than metrically-equivalent deformations that leave the categorical relations intact (as might be produced by viewing a given object from different angles). By contrast, any model that relies only on the properties of the proximal stimulus might be expected to treat all metrically-equivalent deformations as equivalent. Such a model may learn that some distortions are more important than others in the context of specific objects, but it is unclear why they would show a general tendency to treat categorical deformations as different than metric ones since there is no heuristic that assumes that categorical relations between parts is central feature of object shape representations. (Indeed, there is no explicit encoding of parts at all.) Instead, all deformations are simply changes in the locations of features in the image.

Hummel and Stankiewicz [24] explored this question in the context of comparing structural description and view based models of human vision. They created a collection of shapes modeled on Tarr and Pinker's (1989) simple "objects". Each object consisted of a collection of lines connected at right angles (Figure 1). Hummel and Stankiewicz then created two deformations of each of these *Basis* object. One deformation, the

relational deformation (Rel), was identical to the Basis object from which it was created except that one line was moved so that its “above/below” relation to the line to which it was connected changed (from above to below or vice-versa). This deformation differed from the Basis object in the coordinates of one part and in the categorical relation of one part to another. The other deformation, the *coordinates* deformation (Cod), moved two lines in the Basis object in a way that preserved the categorical spatial relations between all the lines composing the object, but changed the coordinates of two lines. Note that both variants deformed the objects (proximal stimulus) but the relational variant changes categorical relations between parts of the object.

Across five experiments participants first learned to classify a set of base objects and then in a test phase were asked to identify the base objects and reject the Rel and Cod stimuli. The experiments differed in the specific set of images used, the specific tasks, the duration of the stimuli, but across all experiments, participants found it easy to reject the Rel foils and difficult to reject the Cod foils. The effects were not subtle. In Experiment 1 (that used the stimuli from Figure 1) participants mistook the Rel and Cod images as the base approximately 10% and 90%, respectively, with similar findings observed across experiments. Hummel and Stankiewicz took these findings to support the claim that humans encode objects in terms of the categorical relations between their parts, consistent with the predictions of the structural description theories that propose a heuristic approach to human shape representation [18].

Testing CNNs on the Hummel and Stankiewicz [24] stimuli

The findings of Hummel and Stankiewicz [24] provide a critical test for any model that claims to be a theory of human vision: if shape-bias in humans is a consequence of optimising over an object recognition task, then it should also lead to shape representations that are more sensitive to relational deformations than coordinate deformations. To test this hypothesis, we replicated the experimental setup of Hummel and Stankiewicz, replacing human participants with CNNs. We trained the network on the Basis shapes shown in Figure 1. Once the network had learnt to categorise these Basis shapes, we tested how it categorised the Basis stimuli as well as Rel (relational) and Cod (coordinate) deformations of each shape (see Methods for details). The results are shown in Figure 2 following three different training conditions. The panel on the left shows the accuracy when no data augmentation was used during training, consistent with the training conditions in [24]. That is, the network learned six different Basis

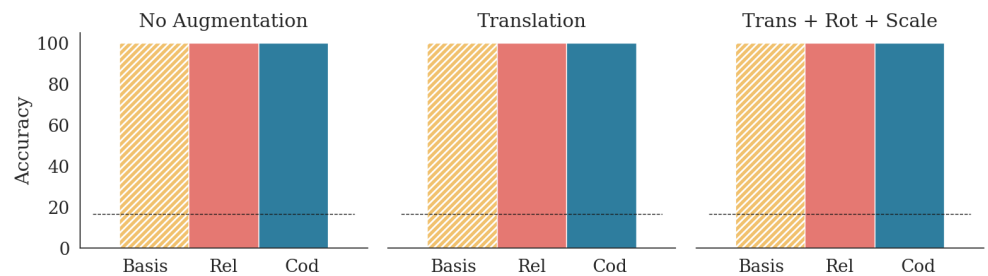


Fig 2. Accuracy on the Basis shapes as well as the two types of deformations (Rel and Cod) for a VGG-16 network. In each case, the model was trained on the set of Basis shapes shown in Figure 1 above, and the training set consisted of (left) exactly these shapes in fixed position, scale and rotation, (middle) Basis shapes were translated to different positions on the canvas but presented at a fixed scale and rotation, and (right) Basis shapes were translated, rotated and scaled.

stimuli (one for each category) with no variation of the Basis images across training trials. The network learned the task perfectly (accuracy for Basis shapes is 100%). But unlike humans, it also generalised perfectly to the Rel and Cod deformations (accuracy remained 100% under the two conditions). The two panels on the right show the results following two different data augmentation training conditions, consistent with standard machine learning approaches. The same results were obtained.

These results are qualitatively different from the observations of Hummel and Stankiewicz [24]: while human participants also generalised to Cod (coordinate) deformations, their performance for Rel (relational) deformations was very different (mistaking the Rel deformations as Basis shapes less than 10 percent of the time).

One possible explanation for the perfect performance on the Rel and Cod stimuli across the three training conditions is that the CNNs were forced to make a response, and the softmax classification function obscured the fact that the model was treating the Rel and Cod stimuli differently in the hidden layers of the network. In other words, it is possible that the Basis-Rel similarity was less than the Basis-Cod similarity (like humans) but these dissimilarities did not manifest in the outputs of the model.

We tested this hypothesis by examining the internal representations of the trained network. Figure 3 shows the average similarity between internal representations for a Basis image and its Rel deformation (Ba-Rel), as well as its Cod deformation (Ba-Cod). (see Methods for how this distance was computed). The internal representations are computed at all convolutional and fully connected layers within the network. We compared these similarities to two baselines: the similarity between two Basis images

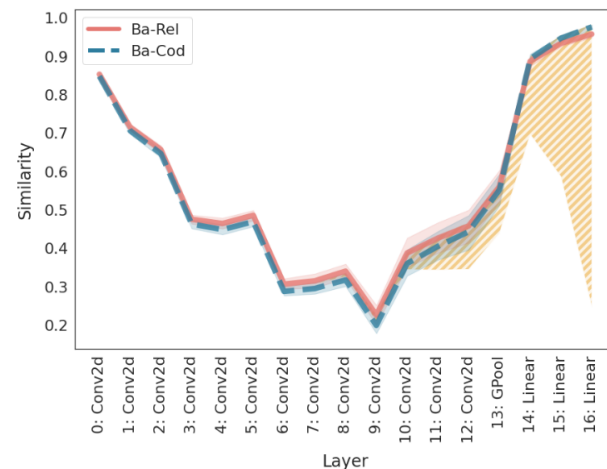


Fig 3. Cosine similarity between internal representations at convolutional and fully connected layers of a trained VGG-16 network. The solid (red) line plots the cosine similarity between the internal representations of a Basis shape and its Rel deformation, while the dashed (blue) line plots the cosine similarity between the internal representations of a Basis shape and its Cod deformation. The hatched (yellow) area shows the upper and lower bounds on cosine similarity (obtained by computing the cosine similarity of images from the same and different categories, respectively). Shaded regions around each line show 95% confidence interval. Based on the results of Hummel and Stankiewicz [24], we would expect the solid (red) line (Ba-Rel) to be closer to the lower, rather than upper bound. Instead we observe that it stays at the upper bound throughout the network and is statistically indistinguishable from the dashed (blue) line, showing that there is no significant difference between the Basis shape and either (relational or coordinate) deformation at any layer of the network.

that belong to the same category and the similarity between two Basis images that belong to different categories. These two baselines provide the upper and lower bounds on similarities respectively. The results show that the network starts with representing all types of images in a similar manner (there is no statistical difference between similarities within or between categories in the early layers) but representations begin to separate in the deeper convolution and fully connected layers (the hatched (yellow) region increases in size as we move left to right because images from different categories have lower similarity than images from the same category). Crucially, there is no statistical difference between Ba-Rel and Ba-Cod, at any layer, both of which are at the upper bound of similarity between representations. That is, the similarity between the representation of a Basis image and its Rel deformation is no less than that of a Basis image and its Cod deformation or even the similarity between two different Basis images from the same category. These results suggest that the lack of difference in the network's performance on Basis, Rel and Cod stimuli extends to its internal representations – that is, we did not find any evidence that suggests that the CNN represents a relational change to an image in any privileged manner compared to a coordinate change.

Although these findings suggest that CNNs represent shape in a dissimilar manner to humans, two alternative explanations should be considered. First, we have assumed that pre-training CNNs on a large dataset of naturalistic images, such as **ImageNet**, should be enough to make them acquire similar shape representations as humans. This is indeed a standard assumption when comparing these models with human data [see, for example, 7]. However, it is possible that pre-training the model in this manner is not enough to teach the models about the importance of relations. Perhaps if CNNs were explicitly trained that relations matter, they would do a better job in characterizing human shape representations. Second, performance on the two deformations was perfect in the simulations above, and perhaps ceiling effects obscured more subtle differences between Rel and Cod stimuli consistent with human vision. We consider these two possibilities in turn.

Teaching relational representations of multi-part objects

In the above experiments, the training environment did not contain an instance of a relational (Rel) or coordinate (Cod) deformation. What if the network was trained to recognise that relational changes are important? In the next set of simulations, we created a training environment with a “relational bias”. We show next that when we do this, the network can learn specific changes to relations but it does not generalise this knowledge to novel (but highly similar) relational changes.

Consider the three augmented training sets shown in Figure 4. In each set the network is trained on the six Basis shapes (and their translation, rotation and scale transformations) just like in the experiments above. In addition, it is also trained on five new shapes. These five shapes are the Rel deformations of the first five Basis shapes. In other words, the training set assigns different categories to a shape and its Rel deformation for five out of six figures. After the network has been trained on these eleven (5 + 5 + 1) shapes, it is tested on the Rel and Cod deformations of the final (unpaired) Basis shape.

The difference between the three datasets lies in the degree of novelty of test images. In all three datasets in Figure 4, the same relation (dashed red circle) is changed between the unpaired Basis shape and its Rel deformation. However, in the first set, there were four other categories (two pairs, highlighted in red rectangles) in the training set where a similar change in relation occurred – that is, for all highlighted categories, there existed another category where the short red segment at the left end of the top bar flipped from “above” to “below” or vice-versa. In the second training set (Set 2 in

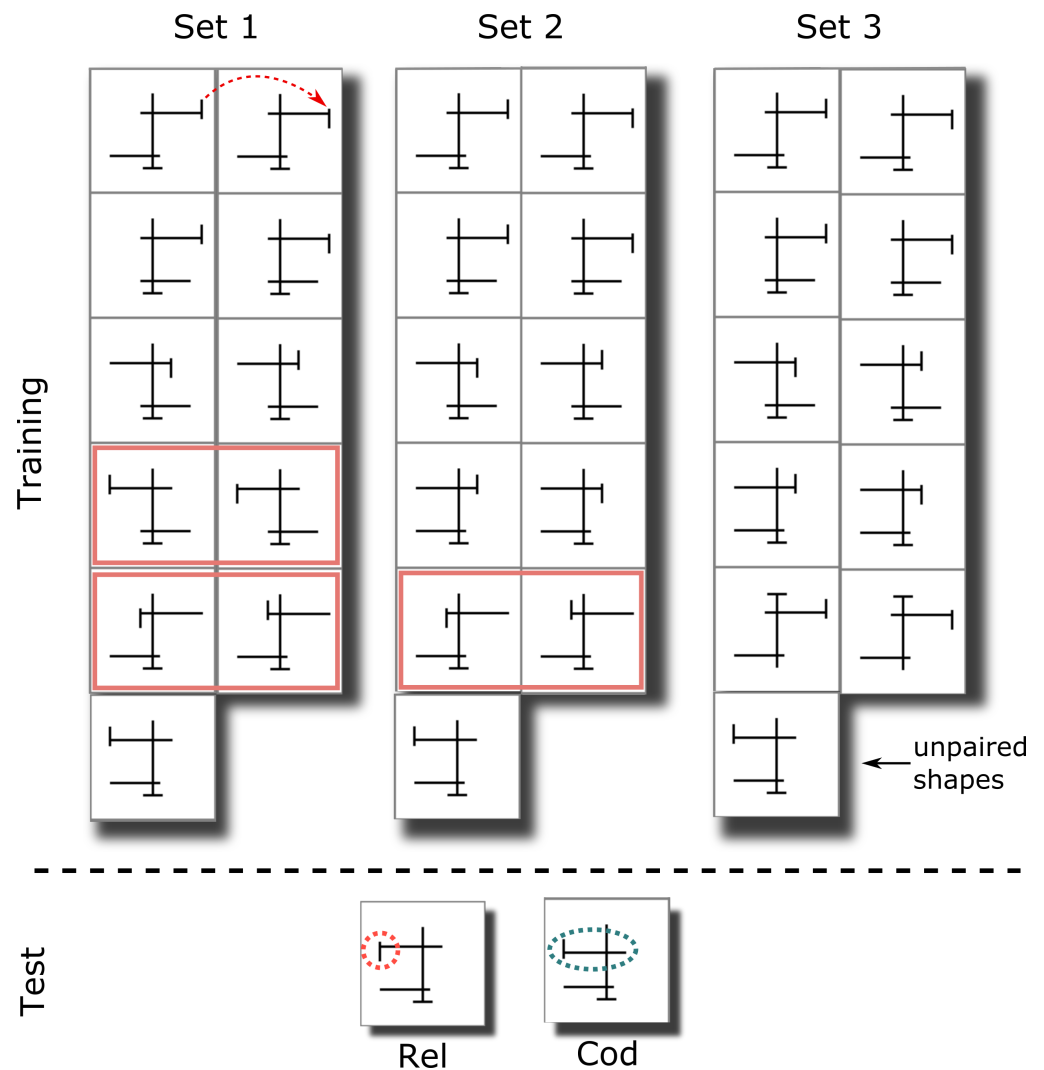


Fig 4. Three training sets that try to teach the network to recognise relational changes. In each set, the first column shows a set of six unique Basis shapes, while the second column shows Rel deformations of the first five shapes (see red arrow). At the bottom are the two test shapes. These test shapes are identical to the eleventh (unpaired) training shape, except for one relational (dashed red circle) or coordinate (dashed blue ellipse) deformation. In Set 1 and Set 2, the difference between the untrained shape and the tested Rel deformation overlaps with some pairs of shapes in the training set (highlighted in solid red rectangles), while in Set 3 there is no overlap between the tested deformation and any trained deformation.

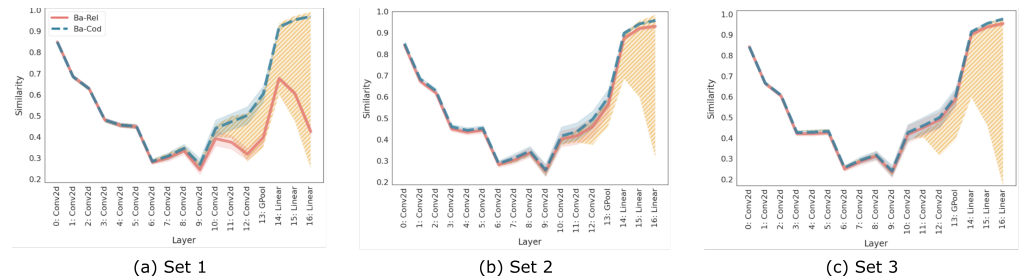


Fig 5. Cosine similarity between Basis image and two types of deformations for a VGG-16 network trained on augmented datasets. Like Figure 3, the hatched (yellow) region shows the upper and lower bound on similarity, the solid (red) line shows similarity between Basis and the relational deformation while the dashed (blue) lines shows similarity between Basis and the coordinate deformation.

Figure 4) there were two categories in the training set where the tested relation changed. However, in this case, this relational change occurred in a different location (closer to central vertical line). In the third training set, the tested relational change was truly novel – that is, none of the trained categories differed in the tested relational change. Of course, the training set still contained categories that differed in a relational change but none of the trained differences overlapped with the tested difference.

Figure 5 shows the cosine similarity in internal representations for CNNs trained on these three modified data sets (as in previous simulations, all networks were pre-trained on ImageNet and re-tuned to each training set). In panel (a) the cosine similarity Ba-Rel is lower than Ba-Cod in deeper layers of the CNN. That is, the network treats the relational deformation as *less* similar to Basis figures than the coordinate deformations. This looks much more like the behaviour of human participants in Hummel and Stankiewicz [24]. But note this training set contains two categories with exactly the same relational change that distinguishes the tested Rel deformation from the corresponding Basis figure. A stronger test is provided in the second case. Here, we observed that the this effect is significantly reduced (panel (b)) – the cosine similarity Ba-Rel is slightly lower than Ba-Cod but by a much smaller degree and the difference only exists in the fully connected layers (also compare results in Figure 15 in Appendix for AlexNet, where this effect is slightly more pronounced but qualitatively similar). The strongest test for whether the network learns relational representations is provided by the third case, where none of the categories in the training set change the exact relation that distinguishes the Rel deformation from the Basis image. Here, we observed (Figure 5(c)) that the effect disappears completely – the cosine similarity Ba-Rel was indistinguishable from Ba-Cod and both similarities were at the upper bound. The network failed to learn that novel relational changes are more important for classification than coordinate changes even when the learning environment contained a “relational bias” – i.e., changing relations led to a change in an image’s category mapping.

Experiment 2: single-part objects

Deformations for testing single-part objects

As detailed above, structural description theories claim that the categorical relations between object parts are encoded in order to build distal representations of multi-part objects. But of course, in order to build distal representations of complex objects, it is also necessary to build distal representations of the parts themselves. This raises the question of what sorts of deformations of the proximal stimulus should allow us to

contrast optimisation and heuristic approaches for identifying the component parts of complex objects or single-part objects? According to the structural description theory [2], certain shape properties of the proximal image are taken by the visual system as strong evidence that individual parts have those properties. For example, if there is a straight or parallel line in the image, the visual system infers that the part contains a straight edge or parallel edges. If the proximal stimulus is symmetrical, it is assumed that the part is symmetrical [see, for example, 41]. These (and other) shape features used to build a distal representation of the object part are called nonaccidental because they would only rarely be produced by accidental alignments of viewpoint. The visual system ignores the possibility that a given nonaccidental feature in the proximal stimulus (e.g., a straight line) is the product of an accidental alignment of eye and distal stimulus (e.g., a curved edge). That is, the human visual system uses nonaccidental proximal features as a heuristic to infer distal representations of object parts.

Critical for our purpose, many of the nonaccidental features described by Biederman [2] are relational features, and indeed, many of the features are associated with Gestalt rules of perceptual organization, such as good continuation, symmetry, and Pragnanz (simplicity). Accordingly, any deformations of the proximal stimulus that alter these nonaccidental features (such as disrupting symmetry) should have a larger impact on classifications than deformations that do not. By contrast, it is not clear that CNNs optimized to classify objects will encode symmetry or other relational features used to build distal representations. Accordingly, CNNs may be insensitive to deformations of symmetry or other relations present in the proximal stimulus.

With this in mind, we created set of seven symmetrical pentagons (Figure 6(a)), and made deformations of these polygons by altering the locations of the vertices composing the polygons in a way that precisely controlled the metric change in the vertices' locations (in the retinal image). Like Experiment 1, we created two types of deformations: (a) a coordinate deformation that parametrically varied the degree to which a polygon rotated in the visual image, vs. (b) a relational change that had an equivalent impact as the corresponding rotation, but instead introduced a shear that changed relative location of the polygon's vertices. Although the specific manipulation is different from that used in Experiment 1, the general logic is the same: one of the deformations preserves the relations between object features while the other changes them. To a model that looks only at proximal stimulus, both deformations lead to an equivalent pixel-by-pixel change, while to a model that infers properties such as symmetry and solidity of the distal stimulus, the coordinate deformation preserves these properties while the relational deformation changes them.

Figure 6(b) shows some examples of test images for one of the trained shapes. These test shapes are organised based on the degree and type of deformation. The degree of relational deformation (shear) of a test image increases as we move from left to right, while the degree of coordinate deformation (rotation) increases as we move from top to bottom. We can also construct test shapes that are a combination of these relational and coordinate deformations. Every shape in Figure 6(b) is a combination of a rotation and a shear of the basis shape in the top-left corner. We have organised these test shapes based on their distance to the basis figure: all shapes along each diagonal have the same cosine distance to the basis shape.¹ and diagonals farther from the basis shape are at a larger distance. Thus, this method gives us a set of test shapes organised according to increasing relational and coordinate changes and matched based on the distance to the basis shape. We could now ask how accuracy degrades on this landscape of test shapes. If the visual system encodes shape as a set of diagnostic features of the retinal image, accuracy should fall as one moves across (perpendicular to) the diagonals

¹We obtained qualitatively similar results when deformations were organised based on their Euclidean distance to the Basis shape.

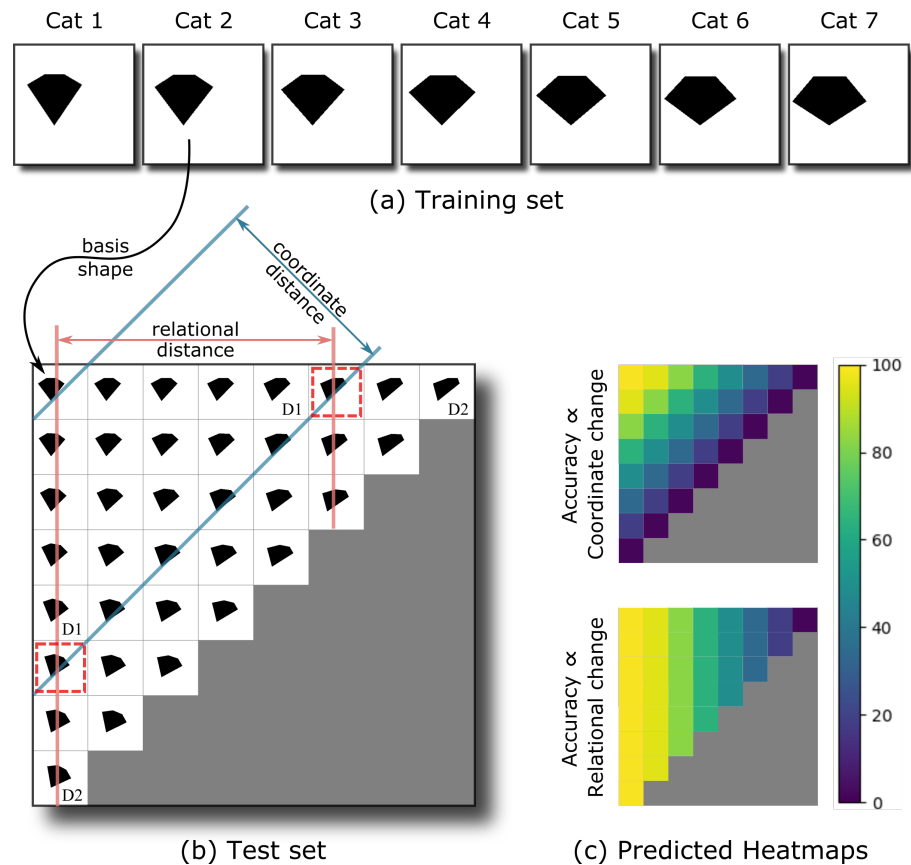


Fig 6. Stimuli used to test shape representations in single-part objects. (a) The shapes in the Basis set used for training. Each shape is presented at various translations and scales. (b) The test set for one of the categories (Cat 2) is obtained by deforming the Basis shape (in the top-left corner) through a combination of rotation and shear operations. Here we have organised these deformations in a matrix based on their coordinate distance (measured as cosine distance) and relational distance (measured as change in relative location of vertices) from the basis shape. All deformations on a diagonal of this matrix are at the same coordinate distance from the Basis shape and all deformations in a column are at the same relational distance from the Basis shape. Highlighted (red) squares show stimuli for computing cosine distance in Figure 8 below. Deformations marked D1 and D2 are used for testing human participants. (c) The predicted accuracy on the test set presented as heat-maps, assuming that accuracy is a function of coordinate distance (top), or relational distance (bottom).

on the landscape. On the other hand, if the visual system encodes shape as a property of the distal stimulus, then changing internal relations should lead to a larger change in classification accuracy than an equivalent coordinate change – that is, the accuracy should fall sharply as one moves left to right along each diagonal. Figure 6(c) shows predicted accuracy on this landscape for the two types of shape representations.

Performance of CNNs

The performance of the network on each test category is shown in Figure 7. Accuracy was highest on the top-left corner for each category (i.e., for the Basis shape) and reduced as the degree of relational and coordinate change was increased. Thus, unlike

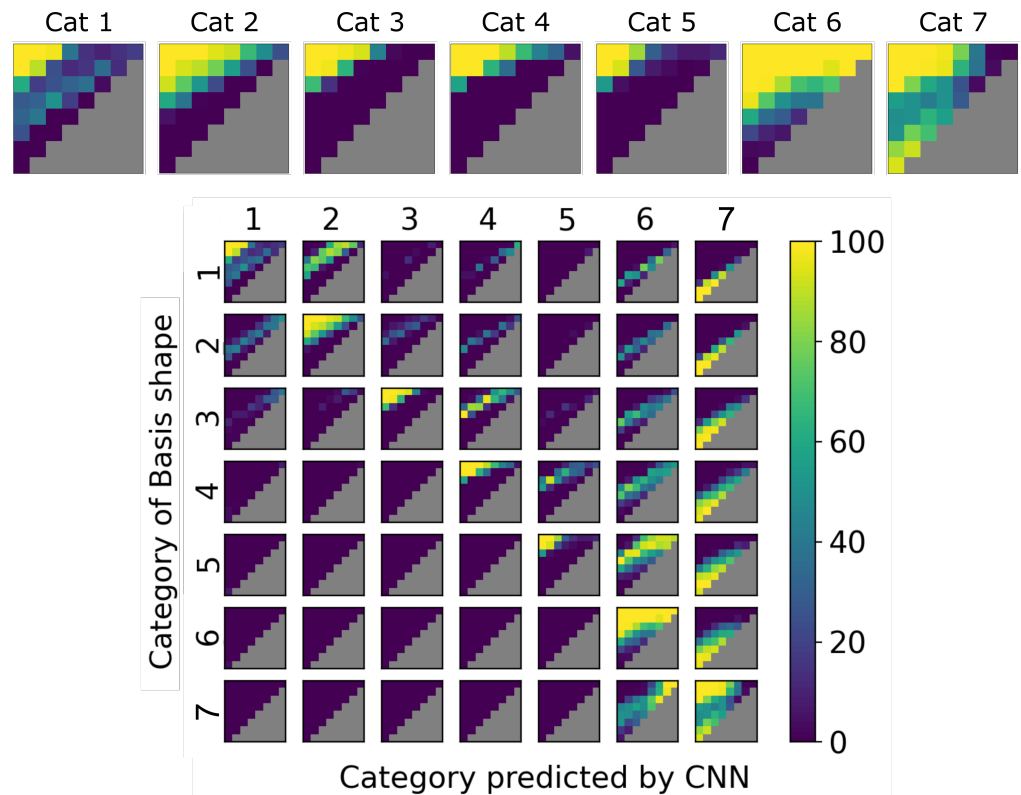


Fig 7. Accuracy on the landscape of relational and coordinate deformations for a VGG-16 trained on a set of seven polygons shown in Figure 6(a). (Top) For each category, every point on this landscape shows the percent of shapes (with a relational and coordinate deformation given by the position on the landscape) accurately classified as the given category by the trained network. (Bottom) The confusion matrix for all deformations. For any heat map, the category label along the row shows the ground truth – i.e., all test shapes used to obtain the heat map were obtained by distorting the basis shape from that category. The category label along the column shows output class label assigned by the network. Therefore, in each row, the diagonal heat map shows the correct classifications, while the off-diagonal heat maps show how each deformation was misclassified.

Experiment 1, where we were able to observe only ceiling performance for both deformations, the design of Experiment 2 allows us to compare how performance degrades for the two types of deformations. We observed that for most categories, accuracy decreased as a function of distance to the Basis shape (perpendicular to the diagonals), rather than relational change (along the diagonals, left to right). In fact, for some categories accuracy *improved* as one moved from left to right along the diagonals.

We also observed a variability across categories in how the performance reduced across deformations. For example, in Figure 7(top), Category 5 shows much larger decrease in performance with increase in rotation of test images than Category 7. To understand why this was the case, it is useful to look at the errors made by the network. Figure 7(bottom) shows the confusion matrix, where each heat map now shows the number of times an output class was chosen for a given input. This confusion matrix shows that the network was prone to mis-classify large rotations from any category as belonging to Category 7. These Type I errors create a bias in the accuracy results for Category 7 in Figure 7(top). The high accuracy for large rotations for Category 7

category are, in fact, misleading as the network classifies large rotations for any category as Category 7. This confusion matrix also shows that the network showed a “rightward” bias – there are more Type I errors in the upper triangle of the matrix than the lower triangle. In other words, the network was more likely to mis-classify images from each category as the category above rather than the category below.

These results suggest that the network does not represent the shapes in this task in a relational manner. If it did, its performance on relational changes should have been a lot worse than its performance on relation-preserving rotations. But, like in the case of the stimuli from Hummel and Stankiewicz [24], accuracy only provides an indirect measure of internal representations. It is possible that even though the network correctly classifies many images with large relational change (images on the right hand side in Figure 6(b)), the internal representation of these images is quite different from the Basis image. We tested this hypothesis by examining the cosine distance between internal representations for the Basis image and two test images that were equidistant from it. An example of these images is highlighted (dashed red squares) in Figure 6(b).

These cosine similarities for an example VGG-16 network are plotted in Figure 8 (we obtained qualitatively similar results for AlexNet – see Figure 17 in Appendix). At all internal layers, we observed that the average similarity between Basis and relational (shear) deformation was higher than the average similarity between the Basis image and its coordinate (rotation) deformation (compare solid and dashed lines in Figure 8). In other words, relational deformation of an image was closer to the Basis image than its coordinate deformation. This is the opposite of what one would expect if the network represented the stimuli in a relational manner.

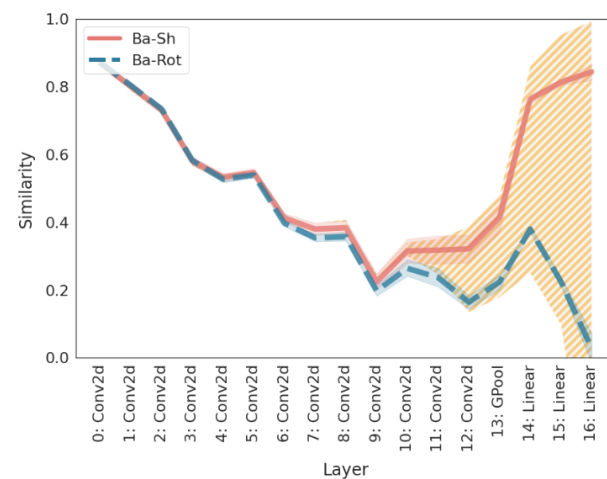


Fig 8. Cosine similarity in internal representations of VGG-16 in Experiment 2. The solid (red) and dashed (blue) lines show the average cosine similarity between Basis images and relational (shear) and coordinate (rotation) deformations, respectively. The hatched (yellow) region shows the bounds on this similarity, with the upper bound determined by the average similarity between Basis images from the same category and lower bound determined by the average similarity between Basis images of different categories. If relational (shear) deformation has a larger affect on internal representations than a coordinate (rotation) deformation, one would expect the solid (red) line to be below the dashed (blue) line.

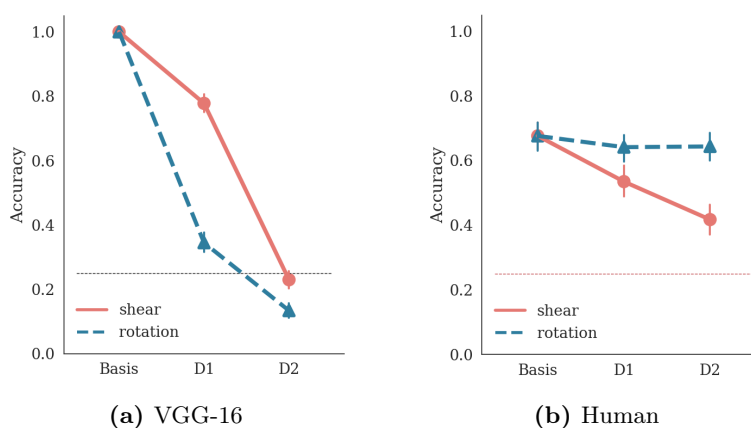


Fig 9. Comparison of (a) network performance and (b) human participants in classifying sheared (dashed line) and rotated (solid line) images. Each panel shows performance under three conditions: basis image, deformation D1 and deformation D2. For the shear deformation, D1 and D2 consists of images in the top row in the fourth and eighth column in Figure 6. For the rotation deformation, D1 and D2 consist of images in the first column and fourth and eighth rows. Error bars show 95% confidence interval and dashed red line shows chance performance.

Performance of human participants

The optimisation view (CNNs) and heuristic view (structural description theory) make contrasting predictions of how performance should degrade when a learned shape is deformed through rotation and shear transformations. In our next experiment, we examined which of these predictions holds for human participants.

We trained 23 participants on the same categorisation task used to train the CNN above. Participants saw the polygons shown in Figure 6(a) and had to learn to categorise them. Once they had learned this task, they were tested on four deformations of each Basis shape – two shears and two rotations. These deformations are marked as D1 and D2 in Figure 6(b) (see Methods for details).

The average accuracy of classification on each of these deformations is shown in Figure 9. In Figure 9(a), we can see that the CNN is more sensitive to rotation than to shear. While performance decreases for both deformations, it decreases more rapidly for rotations. Human participants showed the opposite pattern (Figure 9(b)). There was no significant difference in performance between the basis image and the two rotation deformations (both $t(22) < 3.48$, $p > .28$), while performance decreased significantly for each of shear deformations (both $t(22) > 14.10$, $p < .001$, $d_z > .83$). The largest shear resulted in largest decrease in performance ($M_{difference} = 25.87\%$). Thus, the behaviour of participants was in line with the prediction of structural description theories, where shape is encoded based on relations between internal parts, and in the opposite direction to the performance of the CNN.

Teaching relational representations for single-part objects

One response to this difference between CNNs and humans is that it arises due to the difference in experience of the two systems. Humans experience objects in a variety of rotations and consequently represent a novel object in a rotation invariant manner. CNNs, on the other hand, have not been explicitly trained on objects in different orientations (although ImageNet includes objects in various poses). It could therefore be argued that the CNN does not learn relational representations in the polygons task

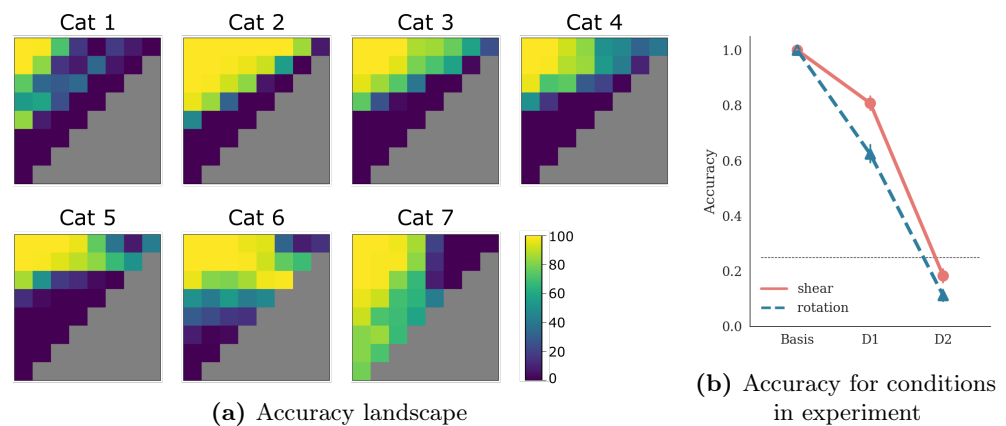


Fig 10. Performance of VGG-16 trained on an augmented dataset where the Basis shapes are not only translated and scaled, but also rotated randomly in the range $[-45^\circ, 0^\circ]$. The network is then tested in the same way as above, where the Basis shapes are deformed by shear or rotation (in the range $[0^\circ, +45^\circ]$) transform. (a) Accuracy (plotted) as percent correct on the landscape of deformations, and (b) shows accuracy for deformations used in experiment with human participants (compare with Figure 9(b) above).

because the training set did not provide an incentive for learning such a representation. Indeed, the optimisation view argues that a bias must be present in the training environment for the visual system to internalise it.

To give the network a better chance of learning to classify based on internal relations, we conducted two further simulations. In the first simulation, we trained the network on rotational deformations of Basis shapes, in addition to translation and scale deformations. To test whether the network generalises based on relational representations, we tested the network on the grid of deformations shown in Figure 6(b) and excluded the rotations used in that grid from the training set. Specifically, the Basis shapes were presented at random rotations in the training set in the range $[-45^\circ, 0]$ and tested on rotations in the range $[0, +45^\circ]$. The network performance on this test grid for each category is shown in Figure 10.

We observed that, despite being trained on this augmented dataset, results remained qualitatively similar. For most categories performance degraded equally or more with a change in rotation than with an equivalent change in shear. That is, the network was *better* at generalising to large relational deformations (shears) than large relation-preserving deformations (rotation). The pattern was different for the final category, where the network showed good performance on large rotations. But examining the confusion matrix again revealed that the high accuracy at large rotations for these two categories was misleading as it was accompanied with large Type I errors: large rotations for shapes of any category were mis-classified as belonging to the final category. Overall, we did not find any evidence for the network learning shapes based on their internal relations.

In the second set of simulations, we selected six (out of seven) categories and trained the network on random translations, scales and *all* rotations ($[0, 360^\circ]$) for these categories. For the seventh category (Cat 3), images were still randomly translated and scaled, but always presented in the upright orientation. We then tested how the network generalised to the two types of deformations for this critical category.

The results of this simulation are shown in Figure 11. Figure 11(a) shows the heat-map of accuracy on the test grid for the left-out category. This heat map showed

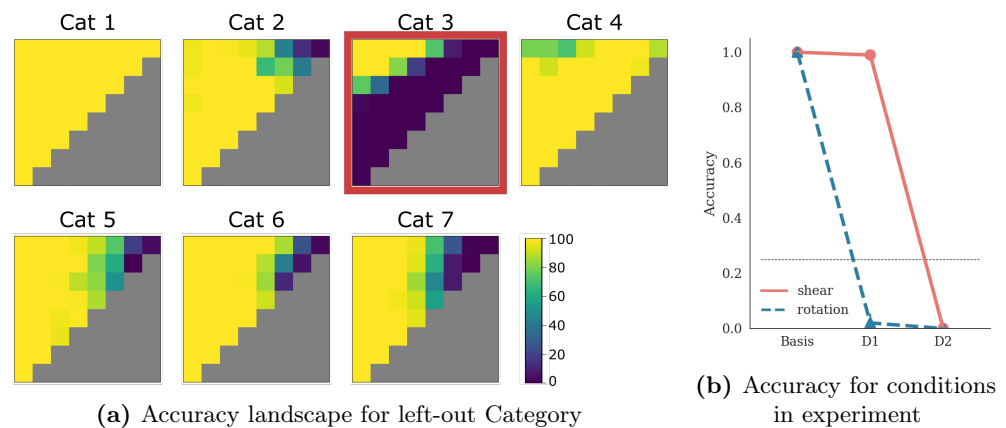


Fig 11. Performance of VGG-16 trained on an augmented dataset, where the Basis shapes are not only translated and scaled for all categories but also rotated at all angle ($[0^\circ, 360^\circ]$) for six out of seven categories. We test whether this augmented training allow the network to generalise better on the relation-preserving deformation (rotation) than the relation-changing deformation (shear) for the left-out category (Cat 3, highlighted). (b) Accuracy for the set of deformations tested with human participants (compare with Figure 9(b) above).

that the network continued showing the pattern observed above – it’s performance decreases across (perpendicular to) the diagonals, but increases as one moves from left-to-right along these diagonals. Figure 11(b) shows the performance on the same conditions as the human experiment (see Figure 9). Again, we see that the performance drops less rapidly across the two shear deformations (dashed line) than the two rotation deformations (solid line). This figure makes it clear that training other orientations on all rotations does not help the network generalise better to novel orientations for the left-out category. In fact, the performance drops more quickly than when none of the categories were rotated in the training set (compare with Figure 9(a)). This is because the network starts classifying novel orientations of the left-out shape as the shapes that it had seen being rotated in the training set.

It may be tempting to think that the differences between humans and CNNs can be reconciled by training CNNs that learn rotation-invariant shapes. However, consider how a CNN achieves rotation-invariance. Figure 12, taken from Goodfellow et al. [12, chap. 9], illustrates how a network consisting of convolution and pooling layers may learn to recognise digits in different orientations. As a result of training on digits (here, the digit 5) oriented in three different directions, the convolution layer develops three different filters, one for each orientation. A downstream pooling unit then amalgamates this knowledge and fires when any one of the convolution filters is activated. Therefore, this pooling unit can be considered as representing the rotation-invariant digit 5. During testing, when the network is presented the digit 5 in any orientation, the corresponding convolution filter gets activated, resulting in a large response in the pooling unit and the network successfully recognises the digit 5, irrespective of it’s orientation.

In contrast, a relational account of shape representation does not rely on developing filters for each orientation of a shape. Indeed, it is not even necessary to observe a shape in all orientations to get, at least some degree of, rotation invariance. All that is needed is to be able to recognise the internal parts of an object and check whether they are in the same relation as the learned shape. Accordingly, many psychological studies have shown that invariance, such as rotation invariance, precedes recognition [3, 4, 5, 20].

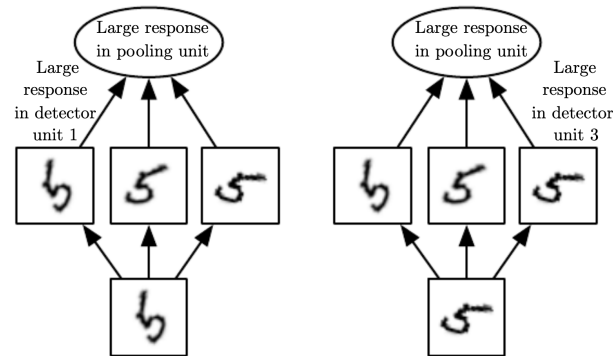


Fig 12. An example of how the operations of convolution and pooling can help a CNN achieve rotation invariance. Taken from Goodfellow et al. [12, chap. 9]

Discussion

In two sets of experiments we have shown that CNNs and humans represent shape in qualitatively different ways. In Experiment 1 we compared how CNNs and humans encode multi-part objects following deformations in the categorical relation between parts (Relational Variants) and deformations that maintained the categorical relations between parts (Coordinate Variants). Whereas humans are highly sensitive to deformations in the categorical relations between parts [24, 18], we found that CNNs are entirely insensitive to these deformations, with performance only a function of the cosine distance between images in pixel space. Furthermore, we could not train CNNs to classify objects on the basis of the relations between parts. In Experiment 2, we compared CNNs and humans in the classification of single part objects when they were deformed by sheering (Relational) or by rotating (Coordinate) manipulations. Again, we found that humans are highly sensitive to relational deformations, whereas CNNs are only sensitive to coordinate manipulations, and once again, CNNs could not learn to be sensitive to relational manipulations.

These findings speak to a current debate concerning the representations that support object recognition in CNNs and humans. On the one hand, some studies have shown that CNNs often classify objects on the basis of texture and other non-shape features [11, 1] whereas humans primary rely on shape [6, 41], suggesting fundamental differences between the two systems. On the other hand, recent studies have shown that changes in the training environment, objective functions, and minor changes to the architectures of CNNs can induce a shape-bias for object recognition in CNNs [10, 11, 15, 14]. Although these later CNNs still do not rely on shape to the same extent as humans, the findings suggests that CNNs may be able account for the the human shape-bias results, and more generally, suggests the goal of optimizing classification performance in CNNs is a promising approach for developing models of human object recognition.

However, our results show that CNNs that learn to classify objects on the basis of shape (our stimuli did not contain any diagnostic texture or colour) learn the wrong sort of shape representation. These findings add to other studies that also highlight the different types of shape representation used by CNNs and the human visual system. For example, Puebla and Bowers [43] have found that CNNs fail to support a simple relational judgement with shapes, namely, whether two shapes are the same or different. Again, this highlights how CNNs trained to process shape ignore relational information. In addition, Baker et al. [1] have shown that CNNs that classify objects based on shape focus on local features and ignore how local features relate to one another in order to encode the global structure of objects. They attribute this failure to a range of

processes that are present in humans but appear to be absent in CNNs, including figure-ground segregation, completing objects behind occluders, encoding boarder ownership, and inferring 3D properties about the object. Consistent with this hypothesis, Jacob et al. [25] have recently highlighted a number of these failures in CNNs, including a failure to represent 3D structure, occlusion, and parts of objects.

One interesting study that provides some evidence to suggest that standard CNNs have similar shape representations to humans was reported by Kubilius et al. [30]. In one of their experiments (Experiment 3), they compared the similarity of representations in various CNNs in response to a change in metric and non-accidental features of single-part objects. For instance, they compared a base object that looked like a slightly curved brick to two objects: one object that was obtained by deforming the base object into a straight brick (a non-accidental change) and a second object that was obtained by deforming the base object into a greatly curved brick (a metric change). Kubilius et al. reported that, like humans, CNNs were more sensitive to non-accidental changes. However, it is unclear whether CNNs were more sensitive to one of their manipulations because of the non-accidental change or because of other confounds accompanying these manipulations. For example, when Kubilius et al. modified some of the base shapes to non-accidental deformations, it was accompanied by a change in shading (luminance) and local features. Recent research [1, 11] has shown that, unlike humans, CNNs are in fact highly sensitive to change in local and textural features and it is unclear whether it is these types changes that are driving the effects observed by Kubilius et al. [30]. More work is required to reconcile their findings with our own.

More generally, our findings raise the question as to whether optimizing CNNs on classification tasks is even the right approach to developing models human object recognition. It is striking how well our findings are well predicted by a classic structural description theory of object recognition that builds a distal representation of objects using heuristics [e.g., 2]. As detailed above, on this theory, the visual system encodes specific features of the proximal stimulus that are best suited for making inferences about the distal object. This includes explicitly coding the relations between parts in order to supporting visual reasoning about objects (e.g., appreciating the similarity and differences of buckets and mugs as discussed above), and encoding parts in terms of non-accidental features that often include relations between features, such as symmetry, in order to infer their 3D distal shape from variable proximal 2D images. Just as predicted, humans are selectively sensitive to these deformations (changes in the relations between parts in Figure 1 and changes in symmetry in Figure 6), whereas CNNs treated these deformations no differently than others.

There is reason to believe that building structural descriptions of object shape based on explicit representations of the spatial relations among an object's parts will prove especially challenging for current CNN models of object recognition. Explicitly relational representations require two representational degrees of freedom, one to specify the parts and relations involved in the representation (e.g., there is a brick shape and a cone shape, and one of them is above the other) and a second to specify their bindings [e.g., to distinguish whether the brick is above the cone or vice-versa; 8, 19, 21, 22, 23]. The units in a CNN have only a single degree of freedom, namely activation, with which to express information, rendering them formally too weak to represent relations, and therefore structural descriptions, explicitly. Other Deep Learning architectures such as Capsule Networks [44], Transformers [49], LSTMs [17] or Neural Turing machines [13] may provide the representational power necessary to represent structural descriptions, but to date this has yet to be demonstrated.

Methods

Generating training and test sets

Training and test sets for Hummel and Stankiewicz [24] We constructed six basis shapes that were identical to the shapes used by Hummel and Stankiewicz [24] in their Experiments 1–3. Each image was sized 196x196 pixels and consisted of five black line segments on a white background organised into different shapes. All images had one short (horizontal) segment at the bottom and one long (vertical) segment in the middle. This left three segments, two long, which were always horizontal, and one short, which was always vertical. The two horizontal segments could be either left-of or right-of the central vertical segment. Additionally, the short vertical segment could be attached to the left-of or the right-of the upper horizontal segment. This means that there were a total of 8 (2x2x2) possible Basis shapes. We selected six out of these to match the six shapes used by Hummel and Stankiewicz [24]. Following Hummel and Stankiewicz [24], we constructed Rel (relational) deformations (called V1 variants by Hummel and Stankiewicz [24]) of each Basis shape by shifting the location of the top vertical segment, so that it's categorical relation to the upper horizontal segment changed from “above” to “below”. Similarly, we constructed Cod (coordinate) deformations (called V2 variants by Hummel and Stankiewicz [24]) by shifting the location of *both* the top horizontal line and the short vertical segments together, so that the categorical relations between all the segments remained the same but the pixel distance (e.g. cosine distance) was at least as large as the pixel distance for the corresponding Rel deformation. Each training set contained 5000 images in each category. When no augmentation was used, all the images in each category were identical – i.e. the shape appeared at the identical location on the canvas for all images in a category. In addition, we constructed two augmented datasets, one in which the Basis image was translated to a random locations (in the range $[-50, +50]$ pixels) on the canvas and another in which it was additionally randomly scaled ($[\frac{1}{2}, 1]$) or rotated ($[-20^\circ, +20^\circ]$).

As described above, we generated three additional datasets for teaching CNNs to recognise relational deformations on Hummel and Stankiewicz's stimuli (see Figure 4). Each of these training sets contained five pairs of Basis shape and one unpaired shape. Each pair consisted of a shape and it's Rel deformation. The test set consisted of Rel and Cod deformations of the unpaired shape. The difference in the three datasets was the amount of overlap between the trained Rel deformations and tested deformation. In the first dataset, there were two pairs of Basis shapes where a Rel deformation was constructed by changing the same categorical relation as the one that differed between the unpaired shape and the tested Rel deformation. In the second set, there was only one such pair. And in the third set, none of the trained shapes differed in the tested deformation. Each training set again consisted of 5000 images, where each image was constructed by translating, scaling and rotating the Basis shape for that category. The test set consisted of 1000 images where each image was constructed by randomly translating, scaling and rotating the Rel and Cod deformations of the unpaired Basis shape.

Training and test sets for polygons task The training set for Experiment 2 consisted of seven symmetric filled pentagons, presented on a white canvas. Each category contained 5000 training images. The training set presented these polygons at different translations and scales, so it was not possible to classify them based on the position of a local feature or the area of the polygon. The difference between Basis shapes for two categories was the angles between the edges. The test set consisted of a grid of shapes that were obtained by deforming the Basis shape of the corresponding category. We used two deformations: rotation, which preserved the internal angles

between edges, and shear, which changed internal angles. To shear a shape, its vertices were horizontally moved by a distance that depended on the vertical distance to the apex. For a vertex with coordinates (x_{old}, y_{old}) , we obtained a new set of vertices, $(x_{new}, y_{new}) = (x_{old} + \lambda(\Delta y)^2, y_{old})$, where λ was the degree of shear and Δy was the distance between y_{old} and y_{apex} , the y-coordinate of the vertex at the apex. Images could also be combination of rotations and shears. To do this, the Basis image was first sheared, then rotated. We measured the distance of a deformed image and test images were organised as shown in Figure 6, where images in each column had the same degree of shear and images along each diagonal had the same (cosine or euclidean) distance to the Basis image. We then obtained twenty exemplars of each deformed image on the grid by randomly translating and scaling the image.

CNN Simulations

We evaluated two deep convolutional neural networks, VGG-16 [45] and AlexNet [29] on the image classification tasks described in the Results section. We obtained qualitatively similar results for both architectures. Therefore, we focus on the results of VGG-16 in the main text and describe the results of AlexNet in Appendix 1.1. Since human participants had a lifetime experience of classifying naturalistic objects prior to the experiment, we used network implementations that had been pre-trained on a set of naturalistic images (ImageNet). In each experiment, the pre-trained network was fine-tuned to classify the 5000 images per category. Each of these images were obtained from the corresponding Basis image in the manner described above. This fine-tuning was performed in the standard manner [50] by replacing the last layer of the classifier to reflect the number of target classes in each dataset. The models learnt to minimise the cross-entropy error by using the Adam optimiser [27] with a learning rate of 10^{-5} and a weight-decay of 10^{-3} . In all simulations, learning continued till the loss function had converged. In most cases, the networks achieved nearly perfect classification on the training set. For the first set of experiments (Hummel and Stankiewicz [24] stimuli), the test set consisted of 1000 Basis, Rel and Cod deformations of each category. In the second set of experiments (polygons stimuli) the test set consisted of 100 exemplars (translation and scale variants) of each test image on the grid (see Figure 6). All simulations were performed using the Pytorch framework [38] and we used torchvision implementation of all models.

To test the similarity of internal representations (Figures 3, 5 and 8), we obtained the embedding of an image at each convolution and fully connected layer of the CNN. For the first set of simulations (Hummel and Stankiewicz [24] stimuli), we selected one (of the six) category and randomly chose 100 pairs of images from the Basis and Rel test set. We then computed the cosine similarity between embeddings of each pair. This gives the estimated average distance in the Ba-Rel condition (solid red line in Figure 3). Similarly the cosine similarity between 100 pairs of Basis and Cod test images gives the Ba-Cod distance (dashed blue line). These distances are compared against two baseline conditions. The upper limit of similarity is given by the similarity of 100 pairs of Basis images from the same category (upper bound of the hatched yellow area in Figure 3). The lower limit is given by the similarity of 100 pairs of Basis images from different category (in each pair, one of the images was from one category and the other from one of the other six categories). The similarity of internal representations for the polygons stimuli is obtained in a similar manner. The similarity Ba-Sh (solid red line in Figure 8) is estimated by measuring the average cosine similarity between embeddings of 100 pairs images from the Basis and sheared sets of the same category. Similarly, Ba-Rot is estimated by measuring the average cosine similarity between embeddings of 100 pairs of images from Basis and rotated sets of the same category.

Behavioral experiment

717

Participants

718

Participants ($N = 37$, $M_{age} = 33$, 70% female) with normal or corrected-to-normal vision were recruited via Prolific and the experiment was conducted on the Pavlovia platform. They were reimbursed a fixed 2 GBP and participants who proceeded to the testing phase ($N = 23$) had a chance to earn a bonus of up to another 2 GBP depending on performance during testing. The average payment was 8 GBP/hour. An written ethics approval for the study was obtained for the study from the University of Bristol Ethics board.

719

720

721

722

723

724

725

Stimuli

726

Four categories were chosen from the total data set for the behavioral study. These are Cat 1, Cat 3, Cat 5 and Cat 7 from Figure 6a. For the test data, we selected two deformations of each type that were matched according to the cosine distance from the basis (trained) image. For the relational deformation, these were the fifth (Deformation D1) and final (Deformation D2) shear in the top row of Figure 6b. For the coordinate deformation, these were the fifth (D1) and final (D2) rotations in the left most column of Figure 6b. This made up the 5 conditions in the experiment: Basis, D1 (Shear), D2 (Shear), D1 (Rotation) and D2 (Rotation). The original stimuli were 224x224 pixels but were re-scaled for each participant to 50% of the vertical resolution to account for the variability in screen size and resolution when running the study online.

727

728

729

730

731

732

733

734

735

736

Procedure

737

Participants completed a supervised training phase in which they learned to categorize basis versions of the four categories. Each training block consisted of 40 stimuli for a total of 200 training trials (50 per category). Feedback on overall accuracy was given at the end of each block. Participants completed up to a maximum of 5 training blocks, or until they reached 85% categorization accuracy in a block. Participants who managed to reach 85% accuracy continued to the test block. The order of trials was randomised for each participant. Each trial started with a fixation cross (750 ms), then the stimulus was presented (500 ms) followed by four response buttons corresponding to the four categories (until response). After participants responded, feedback was given - CORRECT (1 s) if the response was correct, and INCORRECT with additional information about what the correct response should have been (1.5 s) if the response was incorrect.

738

739

740

741

742

743

744

745

746

747

748

749

The training phase was followed by a test phase consisting of five test blocks. Each block consisted of 20 trials for a total of 100 test trials (25 per condition). Like the training phase, the order of test trials was randomised for each participant. The procedure for each test trial was the same as in the training phase apart from the fact that participants were not given any feedback during testing.

750

751

752

753

754

Analysis

755

Four planned comparisons (t-tests) were conducted in order to test whether accuracy rates in each of the shear and rotation conditions differed from accuracy in the basis condition.

756

757

758

Code and Data

759

All code for generating the datasets, simulating the model as well as participant data from Experiment 2 can be downloaded from: <https://github.com/gammagit/distal>

760

761

References

- [1] Nicholas Baker, Hongjing Lu, Gennady Erlikhman, and Philip J Kellman. Deep convolutional networks do not classify based on global object shape. *PLoS computational biology*, 14(12):e1006613, 2018. 762-765
- [2] Irving Biederman. Recognition-by-components: a theory of human image understanding. *Psychological review*, 94(2):115, 1987. 766-767
- [3] Irving Biederman and Eric E Cooper. Evidence for complete translational and reflectional invariance in visual object priming. *Perception*, 20(5):585–593, 1991. 768-769
- [4] Irving Biederman and Eric E Cooper. Size invariance in visual object priming. *Journal of Experimental Psychology: Human Perception and Performance*, 18(1):121, 1992. 770-772
- [5] Irving Biederman and Peter C Gerhardstein. dependent mechanisms in visual object recognition: Reply to tarr and bülhoff (1995). 1995. 773-774
- [6] Irving Biederman and Ginny Ju. Surface versus edge-based determinants of visual recognition. *Cognitive psychology*, 20(1):38–64, 1988. 775-776
- [7] Charles F Cadieu, Ha Hong, Daniel LK Yamins, Nicolas Pinto, Diego Ardila, Ethan A Solomon, Najib J Majaaj, and James J DiCarlo. Deep neural networks rival the representation of primate IT cortex for core visual object recognition. *PLoS Comput Biol*, 10(12):e1003963, 2014. 777-780
- [8] Leonidas AA Doumas, Guillermo Puebla, Andrea E Martin, and John E Hummel. A theory of relation learning and cross-domain generalization. *Psychological Review*, in press. 781-783
- [9] Willis D Ellis. *A source book of Gestalt psychology*. Routledge, 2013. 784
- [10] Reuben Feinman and Brenden M Lake. Learning inductive biases with simple neural networks. *arXiv preprint arXiv:1802.02745*, 2018. 785-786
- [11] Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A Wichmann, and Wieland Brendel. Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness. *arXiv preprint arXiv:1811.12231*, 2018. 787-790
- [12] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT press, 2016. 791-792
- [13] Alex Graves, Greg Wayne, and Ivo Danihelka. Neural turing machines. *arXiv preprint arXiv:1410.5401*, 2014. 793-794
- [14] Katherine Hermann and Andrew Lampinen. What shapes feature representations? exploring datasets, architectures, and training. *Advances in Neural Information Processing Systems*, 33, 2020. 795-797
- [15] Katherine Hermann, Ting Chen, and Simon Kornblith. The origins and prevalence of texture bias in convolutional neural networks. *Advances in Neural Information Processing Systems*, 33, 2020. 798-800
- [16] Julian Hochberg and Virginia Brooks. Pictorial recognition as an unlearned ability: A study of one child’s performance. *The American Journal of Psychology*, 75(4):624–628, 1962. 801-803

- [17] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997. 804
805
- [18] John E Hummel. Reference frames and relations in computational models of object recognition. *Current Directions in Psychological Science*, 3(4):111–116, 1994. 806
807
- [19] John E Hummel. Getting symbols out of a neural architecture. *Connection Science*, 23(2):109–118, 2011. 808
809
- [20] John E Hummel. Object recognition. *Oxford handbook of cognitive psychology*, pages 32–46, 2013. 810
811
- [21] John E Hummel and Irving Biederman. Dynamic binding in a neural network for shape recognition. *Psychological review*, 99(3):480, 1992. 812
813
- [22] John E Hummel and Keith J Holyoak. Distributed representations of structure: A theory of analogical access and mapping. *Psychological review*, 104(3):427, 1997. 814
815
- [23] John E Hummel and Keith J Holyoak. A symbolic-connectionist theory of relational inference and generalization. *Psychological review*, 110(2):220, 2003. 816
817
- [24] John E Hummel and Brian J Stankiewicz. Categorical relations in shape perception. *Spatial vision*, 10(3):201–236, 1996. 818
819
- [25] Georgin Jacob, RT Pramod, Harish Katti, and SP Arun. Qualitative similarities and differences in visual object representations between brains and deep networks. *Nature communications*, 12(1):1–14, 2021. 820
821
822
- [26] Gaetano Kanizsa. *Organization in vision: Essays on Gestalt perception*. Praeger Publishers, 1979. 823
824
- [27] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 825
826
- [28] David C Knill. Perception of surface contours and surface shape: from computation to psychophysics. *JOSA A*, 9(9):1449–1464, 1992. 827
828
- [29] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25:1097–1105, 2012. 829
830
831
- [30] Jonas Kubilius, Stefania Bracci, and Hans P Op de Beeck. Deep neural networks as a computational model for human shape sensitivity. *PLoS computational biology*, 12(4):e1004896, 2016. 832
833
834
- [31] Barbara Landau, Linda B Smith, and Susan S Jones. The importance of shape in early lexical learning. *Cognitive development*, 3(3):299–321, 1988. 835
836
- [32] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436–444, 2015. 837
838
- [33] Gaurav Malhotra, Marin Dujmovic, and Jeffrey S Bowers. Feature blindness: a challenge for understanding and modelling visual object recognition. *bioRxiv*, 2021. 839
840
- [34] Pascal Mamassian and Michael S Landy. Observer biases in the 3d interpretation of line drawings. *Vision research*, 38(18):2817–2832, 1998. 841
842
- [35] Ken Nakayama and Shinsuke Shimojo. Experiencing and perceiving visual surfaces. *Science*, 257(5075):1357–1363, 1992. 843
844

- [36] Ken Nakayama, Zijiang J He, and Shinsuke Shimojo. Visual surface representation: A critical link between lower-level and higher-level vision. *Visual cognition: An invitation to cognitive science*, 2:1–70, 1995. 845–847
- [37] Stephen Palmer. Canonical perspective and the perception of objects. *Attention and performance*, pages 135–151, 1981. 848–849
- [38] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017. 850–852
- [39] Zygmunt Pizlo. Perception viewed as an inverse problem. *Vision research*, 41(24):3145–3161, 2001. 853–854
- [40] Zygmunt Pizlo and Adam K Stevenson. Shape constancy from novel views. *Perception & Psychophysics*, 61(7):1299–1307, 1999. 855–856
- [41] Zygmunt Pizlo, Tadamasa Sawada, Yunfeng Li, Walter G Kropatsch, and Robert M Steinman. New approach to the perception of 3d shape based on veridicality, complexity, symmetry and volume. *Vision research*, 50(1):1–11, 2010. 857–859
- [42] James R Pomerantz, Lawrence C Sager, and Robert J Stoeber. Perception of wholes and of their component parts: some configural superiority effects. *Journal of Experimental Psychology: Human Perception and Performance*, 3(3):422, 1977. 860–862
- [43] Guillermo Puebla and Jeffrey Bowers. Can deep convolutional neural networks support relational reasoning in the same-different task? *bioRxiv*, 2021. 863–864
- [44] Sara Sabour, Nicholas Frosst, and Geoffrey E Hinton. Dynamic routing between capsules. *arXiv preprint arXiv:1710.09829*, 2017. 865–866
- [45] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 867–868
- [46] Linda B Smith, Susan S Jones, Barbara Landau, Lisa Gershkoff-Stowe, and Larissa Samuelson. Object name learning provides on-the-job training for attention. *Psychological science*, 13(1):13–19, 2002. 869–871
- [47] Kent A Stevens. The visual interpretation of surface contours. *Artificial Intelligence*, 17(1-3):47–73, 1981. 872–873
- [48] Michael J Tarr and Steven Pinker. Mental rotation and orientation-dependence in shape recognition. *Cognitive psychology*, 21(2):233–282, 1989. 874–875
- [49] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *arXiv preprint arXiv:1706.03762*, 2017. 876–878
- [50] Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How transferable are features in deep neural networks? *arXiv preprint arXiv:1411.1792*, 2014. 879–880

Acknowledgments 881

This research was supported by the European Research Council Grant Generalization in Mind and Machine, ID number 741134. 882–883

Appendix

884

1.1 Results with AlexNet

885

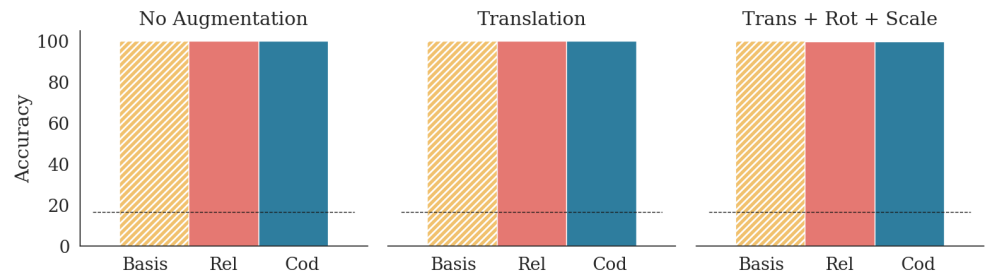


Fig 13. Performance of AlexNet in the test set for Experiment 1. Each panel shows accuracy on the Basis shapes as well as the two types of deformations: relational (Rel) which changes a categorical relation and coordinate (Cod), which preserves all categorical relations. In each case, the model was trained on the set of Basis shapes shown in Figure 1 in the main text, and the training set consisted of (a) exactly these shapes in fixed position, scale and rotation, (b) Basis shapes were translated to different positions on the canvas but presented at a fixed scale and rotation, and (c) Basis shapes were translated, rotated and scaled. Compare with performance of VGG-16 in Figure 2.

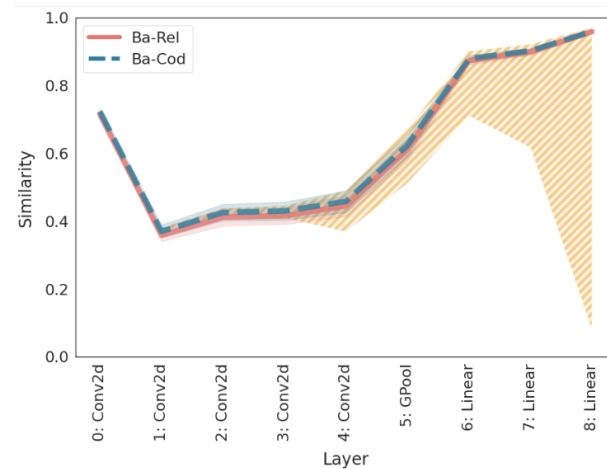


Fig 14. Cosine similarity between the internal representations of Basis images and the internal representation of Rel (relational) and Cod (coordinate) deformations of the Basis image. Like the results for VGG-16 (compare with Figure 3 in the main text), the similarity between Basis images and both types of deformations is at the upper bound throughout the network, showing that the network does not distinguish the trained (Basis) image from its Rel and Cod deformations.

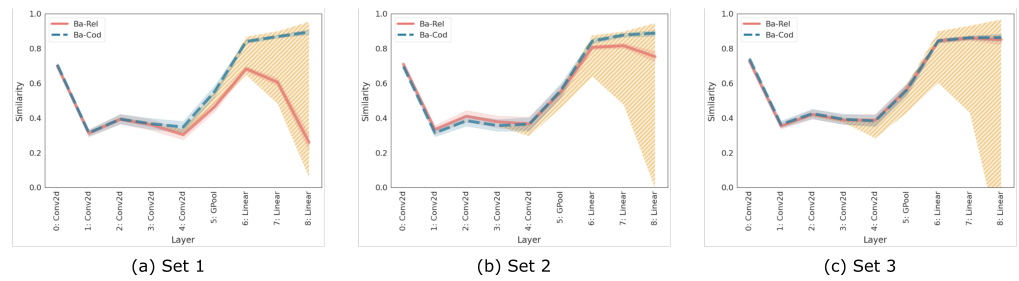


Fig 15. Cosine similarity between Basis image and two types of deformations for AlexNet trained on (a) Set 1, (b) Set 2, and (c) Set 3 in Figure 4. Like the results for VGG-16 (compare with Figure 5), we see that the network learns to distinguish the Rel deformation from the Basis image for Set 1, when it has seen the specific deformation in the training set. But this sensitivity to Rel deformation diminishes in Set 2, when only one pair of trained shapes have a similar deformation and completely lost for Set 3 when the network has been trained on the Rel deformations, but the specific deformation tested is novel.

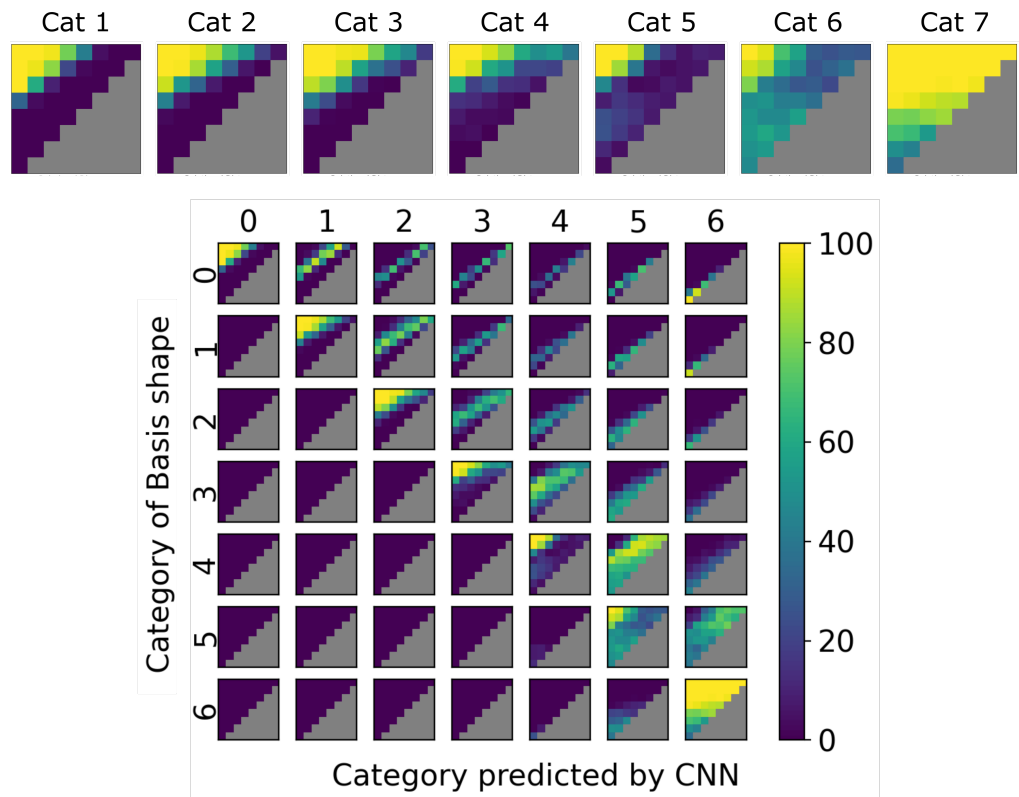


Fig 16. Each heatmap shows accuracy for on test items for a particular category for AlexNet pre-trained on ImageNet and fine-tuned on the dataset in Figure 6. Each cell in the heatmap corresponds to a deformation that is a combination of relational (shear) and coordinate (rotation) transformations of the trained Basis shapes (see Figure 6(a)). The grid at the bottom shows the “confusion matrix” – each heatmap in the grid shows the proportion of responses predicted as the category along the column for a deformation with basis shape taken from the category along the row. Like the results for VGG-16 (compare with Figure 7), we see that accuracy decreases as a function of coordinate distance from the basis shape, rather than the relational distance.

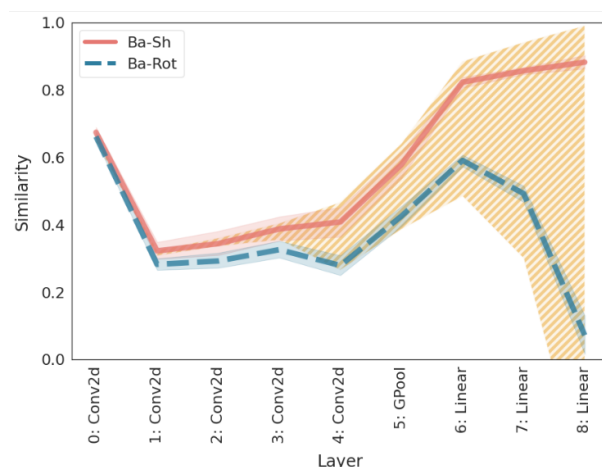


Fig 17. Cosine similarity between internal representations for the Basis shapes and two deformations of the basis shape (dashed red squares in Figure 6(b)) from the polygons dataset at each convolution and fully connected layer of **AlexNet**. Solid (red) line shows the average similarity between representations for a basis shape and its relational (shear) deformation, while dashed (blue) line shows the average similarity between a basis shape and its coordinate (rotation) transformation. The hatched area shows the bounds on similarity, with the upper bound determined by the average similarity between two basis shapes from the same category and lower bounds determined by the average similarity between two basis shapes of different categories. Like the results for **VGG-16** (compare with Figure 8), we observed that the network treated the relational (shear) deformation as being more similar to the basis shape than the coordinate (rotation) deformation. This was the opposite behaviour to the human participants (see Figure 18(b)).

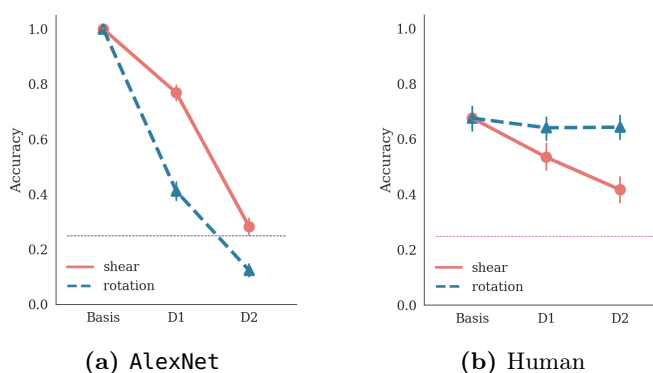


Fig 18. Comparison of (a) network performance and (b) human participants in classifying sheared (dashed line) and rotated (solid line) images. Each panel shows performance under three conditions: basis image, deformation D1 and deformation D2. For the shear deformation, D1 and D2 consists of images in the top row in the fourth and eighth column in Figure 6. For the rotation deformation, D1 and D2 consist of images in the first column and fourth and eighth rows. Error bars show 95% confidence interval and dashed red line shows chance performance. Note that the results in (b) are reproduced here for convenience but are the results of the same experiment reported in Figure 9(b) above.

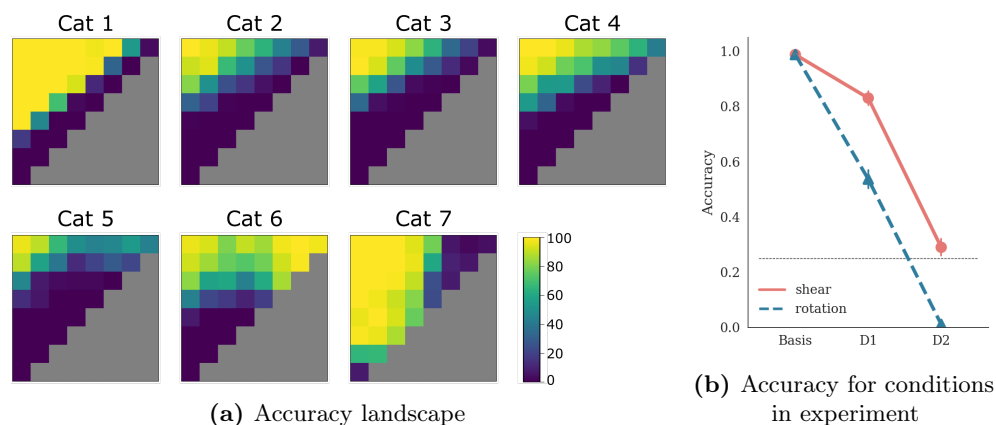


Fig 19. Performance of AlexNet trained on an augmented dataset where the Basis shapes are not only translated and scaled but also randomly rotated in the range $[-45^\circ, 0^\circ]$. The network is then tested on on shear and rotation deformations in the range $[0^\circ, +45^\circ]$. Like the results for VGG-16 (compare with Figure 10), we observed that even when the network was trained on some rotations, it's performance on untrained rotations (a coordinate transformation) was still worse than shears (a relational transformation). (b) shows accuracy for deformation level D1 and D2 used for testing human participants (compare with human performance in Figure 18(b) above).

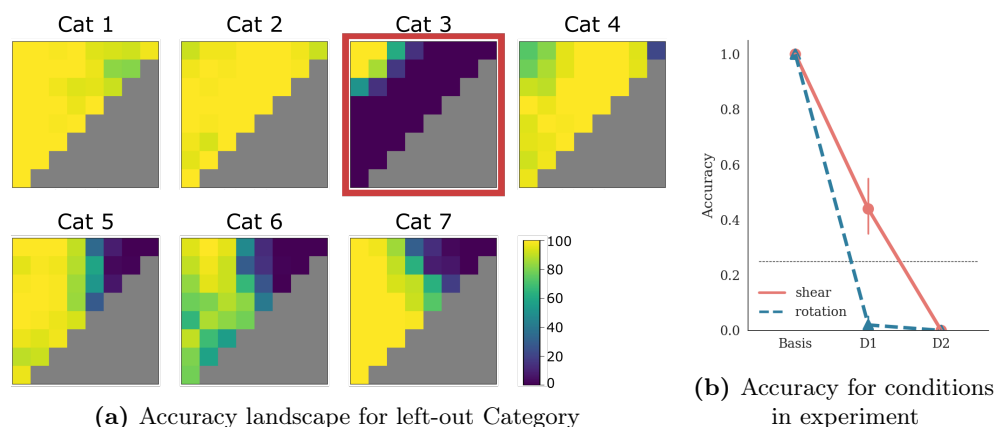


Fig 20. Performance of AlexNet trained on an augmented dataset where the basis shapes are not only randomly translated and scaled but also rotated. For six out of seven categories, the network is trained on *all* rotations ($[0, 360^\circ]$). We then tested the network on the left-out category (Cat 3, highlighted with red square in (a)) on untrained rotations and shears. However, we observed that despite being trained in this manner, the accuracy degraded as a function of the coordinate deformation, rather than the relation deformation. (b) shows the performance of this network for deformations D1 and D2 used to test human participants (compare with results in 18(b) above).