
Generalisation in humans and deep neural networks

Robert Geirhos^{1-3*§}Carlos R. Medina Temme^{1*}Jonas Rauber^{2,3*}Heiko H. Schütt^{1,4,5}Matthias Bethge^{2,6,7*}Felix A. Wichmann^{1,2,6,8*}¹Neural Information Processing Group, University of Tübingen²Centre for Integrative Neuroscience, University of Tübingen³International Max Planck Research School for Intelligent Systems⁴Graduate School of Neural and Behavioural Sciences, University of Tübingen⁵Department of Psychology, University of Potsdam⁶Bernstein Center for Computational Neuroscience Tübingen⁷Max Planck Institute for Biological Cybernetics⁸Max Planck Institute for Intelligent Systems

*Joint first / joint senior authors

§To whom correspondence should be addressed: robert.geirhos@bethgelab.org

Abstract

We compare the robustness of humans and current convolutional deep neural networks (DNNs) on object recognition under twelve different types of image degradations. First, using three well known DNNs (ResNet-152, VGG-19, GoogLeNet) we find the human visual system to be more robust to nearly all of the tested image manipulations, and we observe progressively diverging classification error-patterns between humans and DNNs when the signal gets weaker. Secondly, we show that DNNs trained directly on distorted images consistently surpass human performance on the exact distortion types they were trained on, yet they display extremely poor generalisation abilities when tested on other distortion types. For example, training on salt-and-pepper noise does not imply robustness on uniform white noise and vice versa. Thus, changes in the noise distribution between training and testing constitutes a crucial challenge to deep learning vision systems that can be systematically addressed in a lifelong machine learning approach. Our new dataset consisting of 83K carefully measured human psychophysical trials provide a useful reference for lifelong robustness against image degradations set by the human visual system.

1 Introduction

1.1 Deep neural networks as models of human object recognition

The visual recognition of objects by humans in everyday life is rapid and seemingly effortless, as well as largely independent of viewpoint and object orientation [1]. The rapid and primarily foveal recognition during a single fixation has been termed *core object recognition* (see [2] for a review). We know, for example, that it is possible to reliably identify objects in the central visual field within a single fixation in less than 200 ms when viewing “standard” images [2–4]. Based on the rapidness of object recognition, core object recognition is often thought to be achieved with mainly feedforward processing although feedback connections are ubiquitous in the primate brain. Object recognition

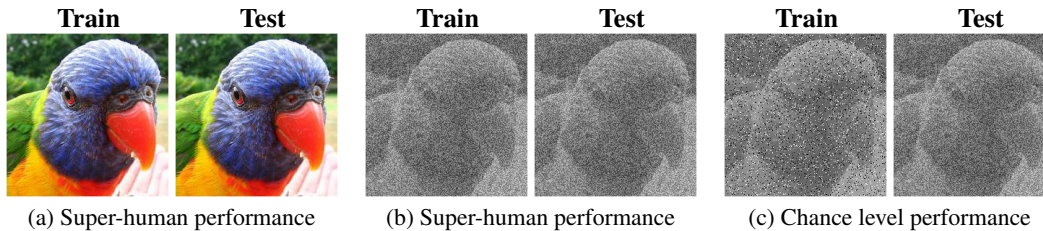


Figure 1: Classification performance of ResNet-50 trained from scratch on (potentially distorted) ImageNet images. **(a)** Classification performance when trained on standard colour images and tested on colour images is close to perfect (better than human observers). **(b)** Likewise, when trained and tested on images with additive uniform noise, performance is super-human. **(c)** Striking generalisation failure: When trained on images with salt-and-pepper noise and tested on images with uniform noise, performance is at chance level—even though both noise types do not seem much different to human observers.

in the primate brain is believed to be realised by the ventral visual pathway, a hierarchical structure consisting of areas V1-V2-V4-IT, with information from the retina reaching the cortex in V1 (e.g. [5]).

Until a few years ago, animate visual systems were the only ones known to be capable of broad-ranging visual object recognition. This has changed, however, with the advent of brain-inspired deep neural networks (DNNs) which, after having been trained on millions of labeled images, achieve human-level performance when classifying objects in images of natural scenes [6]. DNNs are now employed on a variety of tasks and set the new state-of-the-art, sometimes even surpassing human performance on tasks which only a few years ago were thought to be beyond an algorithmic solution for decades to come [7, 8]. Since DNNs and humans achieve similar accuracy, a number of studies have started investigating similarities and differences between DNNs and human vision [9–24]. On the one hand, the network units are an enormous simplification given the sophisticated nature and diversity of neurons in the brain [25]. On the other hand, often the strength of a model lies not in replicating the original system but rather in its ability to capture the important aspects while abstracting from details of the implementation (e.g. [26, 27]).

One of the most remarkable properties of the human visual system is its ability to generalise robustly. Humans generalise across a wide variety of changes in the input distribution, such as across different illumination conditions and weather types. For instance, human object recognition is largely unimpaired even if there are rain drops or snow flakes in front of an object. While humans are certainly exposed to a large number of such changes during their preceding lifetime (i.e., at “training time”, as we would say for DNNs), there seems to be something very generic about the way the human visual system is able to generalise that is not limited to the same distribution one was exposed to previously. Otherwise we would not be able to make sense of a scene if there was some sort of “new”, previously unseen noise. Even if one never had a shower of confetti before, one is still able to effortlessly recognise objects at a carnival parade. Naturally, such generic, robust mechanisms are not only desirable for animate visual systems but also for solving virtually any visual task that goes beyond a well-confined setting where one knows the exact test distribution already at training time. Deep learning for autonomous driving may be one prominent example: one would like to achieve robust classification performance in the presence of confetti, despite not having had any confetti exposure during training time. Thus, from a machine learning perspective, general noise robustness can be used as a highly relevant example of lifelong machine learning [28] requiring generalisation that does not rely on the standard assumption of independent, identically distributed (i.i.d.) samples at test time.

1.2 Comparing generalisation abilities

Generalisation in DNNs usually works surprisingly well: First of all, DNNs are able to learn sufficiently general features on the training distribution to achieve a high accuracy on the i.i.d. test distribution despite having sufficient capacity to completely memorise the training data [29], and

considerable effort has been devoted to understand this phenomenon (e.g. [30–32]).¹ Secondly, features learned on one task often transfer to only loosely related tasks, such as from classification to saliency prediction [33], emotion recognition [34], medical imaging [35] and a large number of other transfer learning tasks [36]. However, transfer learning still requires a substantial amount of training before it works on the new task. Here, we focus on a third setting that adopts the lifelong machine learning point of view of generalisation [37]: How well can a visual learning system cope with a new image degradation after it has learned to cope with a certain set of image distortions before. As a measure of object recognition robustness we can test the ability of a classifier or visual system to tolerate changes in the input distribution up to a certain degree, i.e., to achieve high recognition performance despite being evaluated on a test distribution that differs to some degree from the training distribution (testing under realistic, non-i.i.d. conditions). Using this approach we measure how well DNNs and human observers cope with parametric image manipulations that gradually distort the original image.

First, we assess how top-performing DNNs that are trained on ImageNet, GoogLeNet [38], VGG-19 [39] and ResNet-152 [40], compare against human observers when tested on twelve different distortions such as additive noise or phase noise (see Figure 2 for an overview)—in other words, how well do they generalise towards previously unseen distortions.² In a second set of experiments, we train networks directly on distorted images to see how well they can in general cope with noisy input, and how much training on distortions as a form of data augmentation helps in dealing with other distortions. Psychophysical investigations of human behaviour on object recognition tasks, measuring accuracies depending on image colour (greyscale vs. colour), image contrast and the amount of additive visual noise have been powerful means of exploring the human visual system, revealing much about the internal computations and mechanisms at work [42–48]. As a consequence, similar experiments might yield equally interesting insights into the functioning of DNNs, especially as a comparison to high-quality measurements of human behaviour. In particular, human data for our experiments were obtained using a controlled lab environment (instead of e.g. Amazon Mechanical Turk without sufficient control about presentation times, display calibration, viewing angles, and sustained attention of participants). Our carefully measured behavioural datasets—twelve experiments encompassing a total number of 82,880 psychophysical trials—as well as materials and code are available online at <https://github.com/rgeirhos/generalisation-humans-DNNs>.

2 Methods

We here report the core elements of employed paradigm, procedure, image manipulations, observers and DNNs; this is aimed at giving the reader just enough information to understand experiments and results. For in-depth explanations we kindly refer to the comprehensive supplementary material, which seeks to provide exhaustive and reproducible experimental details.

2.1 Paradigm, procedure & 16-class-ImageNet

For this study, we developed an experimental paradigm aimed at comparing human observers and DNNs as fair as possible by using a forced-choice image categorisation task.³ Achieving a fair psychophysical comparison comes with a number of challenges: First of all, many high-performing DNNs are trained on the ILSRVR 2012 database [50] with 1,000 fine-grained categories (e.g., over a hundred different dog breeds). If humans are asked to name objects, however, they most naturally categorise them into so-called entry-level categories (e.g. dog rather than German shepherd). We thus developed a mapping from 16 entry-level categories such as dog, car or chair to their corresponding ImageNet categories using the WordNet hierarchy [51]. We term this dataset “16-class-ImageNet” since it groups a subset of ImageNet classes into 16 entry-level categories (airplane, bicycle, boat, car, chair, dog, keyboard, oven, bear, bird, bottle, cat, clock, elephant, knife, truck). In every experiment, then, an image was presented on a computer screen and observers had to choose the correct category by clicking on one of these 16 categories. For pre-trained DNNs, the sum of all softmax values mapping to

¹Still, DNNs usually need orders of magnitude more training data in comparison to humans, as explored by the literature on one-shot or few-shot learning (see e.g. [23] for an overview).

²We have reported a subset of these experiments on arXiv in an earlier version of this paper [41].

³This is the same paradigm as reported in [49].

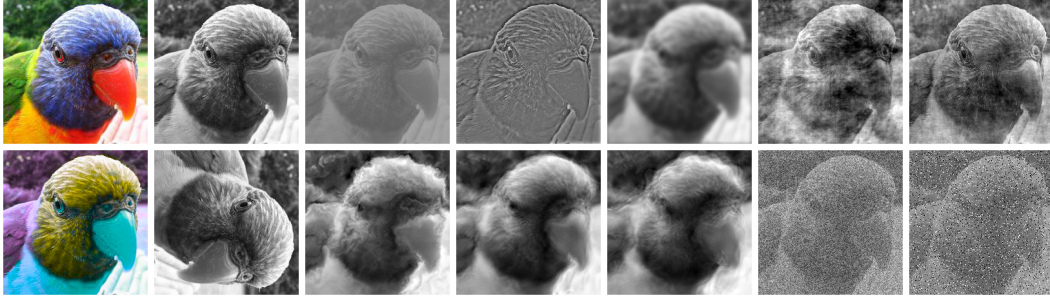


Figure 2: Example stimulus image of class bird across all distortion types. From left to right, image manipulations are: colour (undistorted), greyscale, low contrast, high-pass, low-pass (blurring), phase noise, power equalisation. Bottom row: opponent colour, rotation, Eidolon I, II and III, additive uniform noise, salt-and-pepper noise. Example stimulus images across all used distortion levels are available in the supplementary material.

a certain entry-level category was computed. The entry-level category with the highest sum was then taken as the network’s decision. A second challenge is the fact that standard DNNs only use feedforward computations at inference time, while recurrent connections are ubiquitous in the human brain [52, 53].⁴ In order to prevent this discrepancy from playing a major confounding role in our experimental comparison, presentation time for human observers was limited to 200 ms. An image was immediately followed by a 200 ms presentation of a noise mask with $1/f$ spectrum, known to minimise, as much as psychophysically possible, feedback influence in the brain.

2.2 Observers & pre-trained deep neural networks

Data from human observers were compared against classification performance of three pre-trained DNNs: VGG-19 [39], GoogLeNet [38] and ResNet-152 [40]. For each of the twelve experiments that were conducted, either five or six observers participated (with the exception of the colour experiment, for which only three observers participated since similar experiments had already been performed by a number of studies [48, 55, 56]). Observers reported normal or corrected-to-normal vision.

2.3 Image manipulations

A total of twelve experiments were performed in a well-controlled psychophysical lab setting. In every experiment, a (possibly parametric) distortion was applied to a large number of images, such that the signal strength ranged from ‘no distortion / full signal’ to ‘distorted / weak(er) signal’. We then measured how classification accuracy changed as a function of signal strength. Three of the employed image manipulations were dichotomous (colour vs. greyscale, true vs. opponent colour, original vs. equalised power spectrum); one manipulation had four different levels (0, 90, 180 and 270 degrees of rotation); one had seven levels (0, 30, ..., 180 degrees of phase noise) and the other distortions had eight different levels. Those manipulations were: uniform noise, controlled by the ‘width’ parameter indicating the bounds of pixel-wise additive uniform noise; low-pass filtering and high-pass filtering (with different standard deviations of a Gaussian filter); contrast reduction (contrast levels from 100% to 1%) as well as three different manipulations from the eidolon toolbox [57]). The three eidolon experiments correspond to different versions of a parametric image manipulation, with the ‘reach’ parameter controlling the strength of the distortion. Additionally, for experiments with training on distortions, we also evaluated performance on stimuli with salt-and-pepper noise (controlled by parameter p indicating probability of setting a pixel to either black or white; $p \in [0, 10, 20, 35, 50, 65, 80, 95]\%$). More information about the different image manipulations is provided in the supplementary material (Section Image preprocessing and distortions), where we also show example images across all manipulations and distortion levels (Figures 10, 11, 12, 13, 14). For a brief overview, Figure 2 depicts one exemplary manipulation per distortion. Overall, the manipulations we used were chosen to reflect a large variety of possible distortions.

⁴But see e.g. [54] for a critical assessment of this argument.

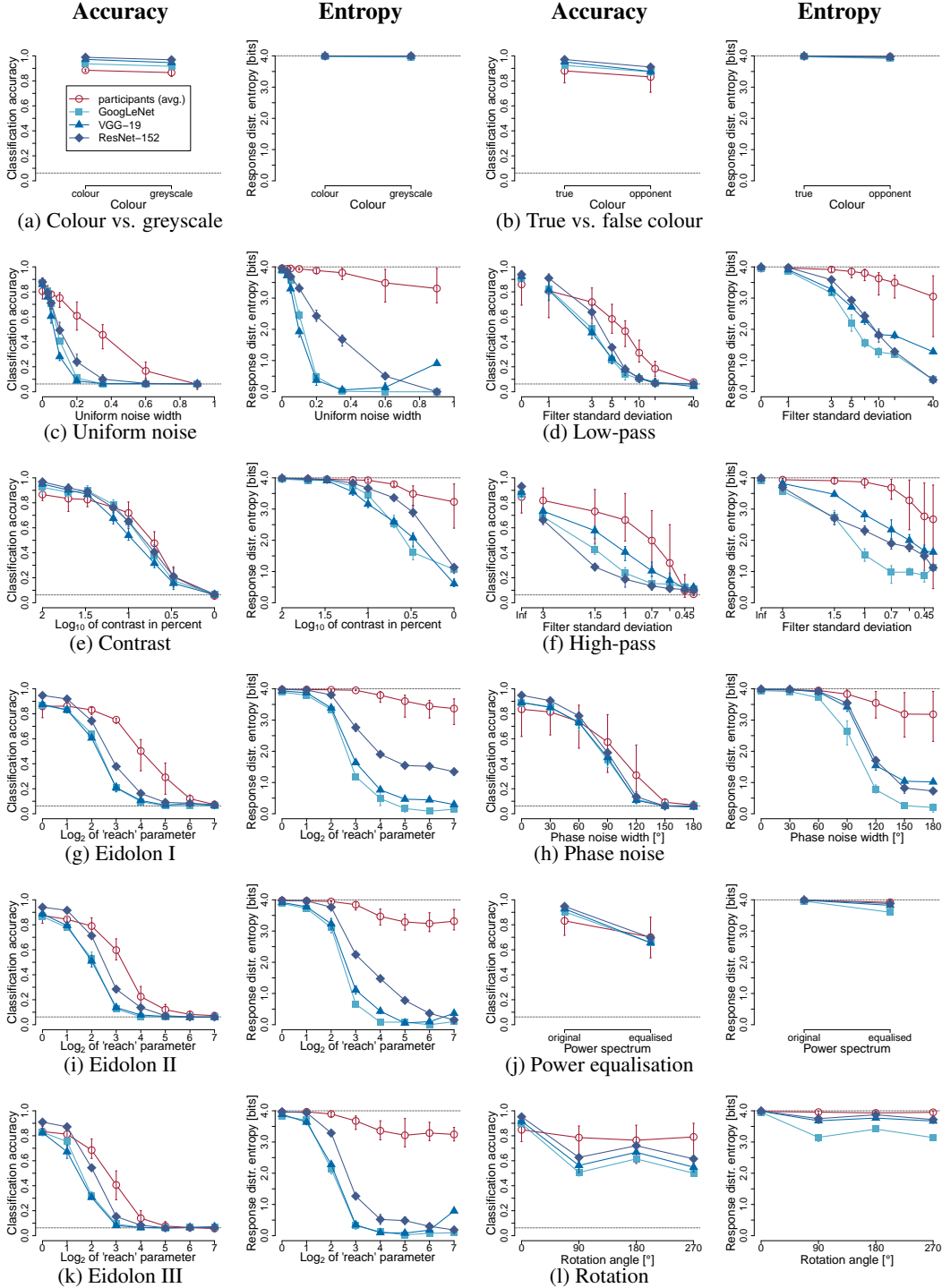


Figure 3: Classification accuracy and response distribution entropy for **GoogLeNet**, **VGG-19** and **ResNet-152** as well as for **human observers**. ‘Entropy’ indicates the Shannon entropy of the response/decision distribution (16 classes). It here is a measure of bias towards certain categories: using a test dataset that is balanced with respect to the number of images per category, responding equally frequently with all 16 categories elicits the maximum possible entropy of four bits. If a network or observer responds prefers some categories over others, entropy decreases (down to zero bits in the extreme case of responding with one particular category all the time, irrespective of the ground truth category). Human ‘error bars’ indicate the full range of results across participants. Image manipulations are explained in Section 2.3 and visualised in Figures 10, 11, 12, 13 and 14.

2.4 Training on distortions

Beyond evaluating standard pre-trained DNNs on distortions (results reported in Figure 3), we also trained networks directly on distortions (Figure 4). These networks were trained on 16-class-ImageNet, a subset of the standard ImageNet dataset as described in Section 2.1. This reduced the size of the unperturbed training set to approximately one fifth. To correct for the highly imbalanced number of samples per class, we weighted each sample in the loss function with a weight proportional to one over the number of samples of the corresponding class. All networks trained in these experiments had a ResNet-like architecture that differed from a standard ResNet-50 only in the number of output neurons that we reduced from 1000 to 16 to match the 16 entry-level classes of the dataset. Weights were initialised with a truncated normal distribution with zero mean and a standard deviation of $\frac{1}{\sqrt{n}}$ where n is the number of output neurons in a layer. While training from scratch, we performed on-the-fly data augmentation using different combinations of the image manipulations. When training a network on multiple types of image manipulations (models B1 to B9 as well as C1 and C2 of Figure 4), the type of manipulation (including *unperturbed*, i.e. standard colour images if applicable) was drawn uniformly and we only applied one manipulation at a time (i.e., the network never saw a single image perturbed with multiple image manipulations simultaneously, except that some image manipulations did include other manipulations per construction: uniform noise, for example, was always added after conversion to greyscale and contrast reduction to 30%). For a given image manipulation, the amount of perturbation was drawn uniformly from the levels used during test time (cf. Figure 3). The remaining aspects of the training followed standard training procedures for training a ResNet on ImageNet: we used SGD with a momentum of 0.997, a batch size of 64, and an initial learning rate of 0.025. The learning rate was multiplied with 0.1 after 30, 60, 80 and 90 epochs (when training for 100 epochs) or 60, 120, 160 and 180 epochs (when training for 200 epochs). Training was done using TensorFlow 1.6.0 [58]. In the training experiments, all manipulations with more than two levels were included except for the eidolon stimuli, since the generation of those stimuli is computationally too slow for ImageNet training. For comparison purposes, we additionally included colour vs. greyscale as well as salt-and-pepper noise (for which there is no human data, but informal comparisons between uniform noise and salt-and-pepper noise strongly suggest that human performance will be similar, see Figure 1c).

3 Generalisation of humans and pre-trained DNNs towards distortions

In order to assess generalisation performance when the signal gets weaker, we tested twelve different ways of degrading images. These images at various levels of signal strength were then shown to both human observers in a lab and to pre-trained DNNs (ResNet-152, GoogLeNet and VGG-19) for classification. The results of this comparison are visualised in Figure 3. While human and DNN performance was similar for comparatively minor colour-related distortions such as conversion to greyscale or opponent colours, we find human observers to be more robust for all of the other distortions: by a small margin for low contrast, power equalisation and phase noise images and by a larger margin for uniform noise, low-pass, high-pass, rotation and all three eidolon experiments. Furthermore, there are strong differences in the error patterns as measured by the response distribution entropy (indicating biases towards certain categories). Human participants' responses were distributed more or less equally amongst the 16 classes, whereas all three DNNs show increasing biases towards certain categories when the signal gets weaker. These biases are not completely explained by the prior class probabilities, and deviate from distortion to distortion. For instance, ResNet-152 almost solely predicts class `bottle` for images with strong uniform noise (irrespective of the ground truth category),⁵ and classes `dog` or `bird` for images distorted by phase noise. One might think of simple tricks to reduce the discrepancy between the response distribution entropy of DNNs and humans. One possible way would be increasing the softmax temperature parameter and assuming that model decisions are sampled from the softmax distribution rather than taking the argmax. However, increasing the response DNN distribution entropy in this way dramatically decreases classification accuracy and thus comes with a trade-off (cf. Figure 8 in the supplementary material).

These results are in line with previous findings reporting human-like processing of chromatic information in DNNs [19] but strong decreases in DNN recognition accuracy for image degradations like

⁵A category-level analysis of decision biases for the uniform noise experiment is provided in the supplementary material, Figure 9.

Evaluation condition	human observers	Model																			
	A1	A2	A3	A4	A5	A6	A7	A8	A9	B1	B2	B3	B4	B5	B6	B7	B8	B9	C1	C2	
colour	88.5	96.7	90.6	50.0	83.1	86.1	84.2	90.8	10.4	8.1	97.9	95.4	72.3	93.0	91.1	92.4	94.9	10.2	11.2	95.5	95.9
greyscale	86.6	87.8	95.6	<u>94.1</u>	86.2	93.2	87.8	90.5	10.3	9.8	94.0	96.8	96.2	93.3	95.7	94.3	90.9	11.4	12.8	94.8	95.1
contrast (5%)	47.6	13.1	14.2	89.4	19.6	39.8	17.1	10.2	28.6	29.0	46.3	51.7	95.1	50.5	79.1	59.4	45.2	34.6	37.9	90.9	88.2
low-pass (std=7)	48.5	18.9	16.1	16.4	78.4	11.9	16.0	9.8	6.9	6.6	16.0	18.6	14.4	87.2	20.5	13.8	13.5	7.1	9.3	74.7	74.9
high-pass (std=0.7)	49.8	21.1	24.7	29.9	11.7	92.6	27.7	8.3	10.4	20.6	25.1	22.8	29.2	25.0	94.3	27.5	28.3	18.9	19.8	91.4	90.7
phase noise (90°)	57.4	23.3	28.3	31.2	27.0	46.6	81.4	24.4	7.4	8.9	30.8	31.4	30.6	31.4	43.4	87.4	24.1	7.8	7.6	82.9	82.6
rotation (90°)	78.5	36.5	43.3	39.9	31.8	40.4	37.7	89.0	8.5	8.0	38.5	41.9	40.3	35.2	40.1	40.5	89.0	8.3	8.8	80.1	80.5
salt-and-pepper noise (0.2)	NA	6.1	6.4	5.8	7.9	6.2	6.2	6.4	79.4	6.2	6.2	6.1	6.3	5.4	5.8	5.7	6.2	89.6	6.2	78.6	13.6
uniform noise (0.35)	45.6	6.2	7.3	6.9	9.0	7.3	6.2	6.0	10.2	80.3	84.6	83.3	85.0	84.6	83.7	82.5	83.8	85.4	89.8	11.0	71.5

= manipulation included in training data

Figure 4: Classification accuracy (in percent) for networks with potentially distorted training data. Rows show different test conditions at an intermediate difficulty (exact condition indicated in brackets, units as in Figure 3). Columns correspond to differently trained networks (leftmost column: human observers for comparison; no human data available for salt-and-pepper noise). All of the networks were trained from scratch on (a potentially manipulated version of) 16-class-ImageNet. Manipulations included in the training data are indicated by a red rectangle; additionally ‘greyscale’ is underlined if it was part of the training data because a certain distortion encompasses greyscale images at full contrast. Models **A1 to A9**: ResNet-50 trained on a single distortion (100 epochs). Models **B1 to B9**: ResNet-50 trained on uniform noise plus one other distortion (200 epochs). Models **C1 & C2**: ResNet-50 trained on all but one distortion (200 epochs). Chance performance is at $\frac{1}{16} = 6.25\%$ accuracy.

noise and blur [13, 14, 59–61]. Overall, DNNs seem to have much more problems generalising to weaker signals than humans, across a wide variety of image distortions. While the human visual system has been exposed to a number of distortions during evolution and lifetime, we clearly had no exposure whatsoever to many of the exact image manipulations that we tested here. Thus, our human data show that a high level of generalisation is, in principle, possible. There may be many different reasons for the discrepancy between human and DNN generalisation performance that we find: Are there limitations in terms of the currently used network architectures (as hypothesised by [60]), which may be inferior to the human brain’s intricate computations? Is it a problem of the training data (as suggested by e.g. [61]), or are today’s training methods / optimisers not sufficient to solve robust and general object recognition? In order to shed light on the dissimilarities we found, we performed a second batch of experiments by training networks directly on distorted images.

4 Training DNNs directly on distorted images

We trained one network per distortion directly and from scratch on (potentially manipulated) 16-class-ImageNet images. The results of this training are visualised in Figure 4 (models A1 to A9). We find that these specialised networks consistently outperformed human observers, by a large margin, on the image manipulation they were trained on (as indicated by strong network performance on the diagonal). This is a strong indication that currently employed architectures (such as ResNet-50) and training methods (standard optimiser and training procedure) are sufficient to ‘solve’ distortions under i.i.d. train/test conditions. We were able to not only close the human-DNN performance gap that was observed by [13] (who fine-tuned networks on distortions, reporting improved but not human-level DNN performance) but to surpass human performance in this respect. While the human visual system certainly has a much more complicated structure [24], this does not seem to be necessary to deal with even strong image manipulations of the type employed here.

However, as noted earlier, robust generalisation is primarily not about solving a specific problem known exactly in advance. We therefore tested how networks trained on a certain distortion type perform when tested on other distortions. These results are visualised in Figure 4 by the off-diagonal

cells of models A1 to A9. Overall, we find that training on a certain distortion slightly improves performance on other distortions in a few instances, but is detrimental in other cases (when compared to a vanilla ResNet-50 trained on colour images, model A1 in the figure).⁶ Performance on salt-and-pepper noise as well as uniform noise was close to chance level for all networks, even for a network trained directly on the respectively other noise model. This may be surprising given that these two types of noise do not seem very different to a human eye (as indicated in Figure 1c). Hence, training a network on one distortion does not generally lead to improvements on other distortions.

Since training on a single distortion alone does not seem to be sufficient to evoke robust generalisation performance in DNNs, we also trained the same architecture (ResNet-50) on two additional settings. Models B1 to B9 in Figure 4 show performance for training on one particular distortion in combination with uniform noise (training consisted of 50% images from each manipulation). Uniform noise was chosen since it seemed to be one of the hardest distortions for all networks, and hence they might benefit from including this particular distortion in the training data. Furthermore, we trained models C1 and C2 on all but one distortion (either uniform or salt-and-pepper noise was left out).

We find that object recognition performance of models B1 to B9 is improved compared to models A1 to A9, both on the distortions they were actually trained on (diagonal entries with red rectangles in Figure 4) as well as on a few of the distortions that were not part of the training data. However, this improvement may be largely due to the fact that models B1 to B9 were trained on 200 epochs instead of 100 epochs as for models A1 to A9, since the accuracy of model B9 (trained & tested on uniform noise, 200 epochs) also shows an improvement towards model A9 (trained & tested on uniform noise, 100 epochs). Hence, in the presence of heavy distortions, training longer may go a long way but incorporating other distortions in the training does not seem to be generally beneficial to model performance. Furthermore, we find that it is possible even for a single model to reach high accuracies on all of the eight distortions it was trained on (models C1 and C2), however for both left-out uniform and salt-and-pepper noise, object recognition accuracy stayed around 11 to 14%, which is by far closer to chance level (approx. 6%) than to the accuracy reached by a specialised network trained on this exact distortion (above 70%, which serves as a lower bound on the achievable performance).

Taken together, these findings indicate that data augmentation with distortions alone may be insufficient to overcome the generalisation problem that we find. It may be necessary to move from asking “why are DNNs generalising so well (under i.i.d. settings)?” [29] to “why are DNNs generalising so poorly (under non-i.i.d. settings)?”. It is up to future investigations to determine how DNNs that are currently being handled as computational models of human object recognition can solve this challenge. At the exciting interface between cognitive science / visual perception and deep learning, inspiration and ideas may come from both fields: While the computer vision sub-area of domain adaptation (see [63] for a review) is working on robust machine inference in spite of shifts in the input distribution, the human vision community is accumulating strong evidence for the benefits of local gain control mechanisms. These normalisation processes seem to be crucial for many aspects of robust animal and human vision [46], are predictive for human vision data [21, 64] and have proven useful in the context of computer vision [65, 66]. It could be an interesting avenue for future research to determine whether there is a connection between neural normalisation processes and DNN generalisation performance.

5 Conclusion

We conducted a behavioural comparison of human and DNN object recognition robustness against twelve different image distortions. In comparison to human observers, we find the classification performance of three well-known DNNs trained on ImageNet—ResNet-152, GoogLeNet and VGG-19—to decline rapidly with decreasing signal-to-noise ratio under image distortions. Additionally, we find progressively diverging patterns of classification errors between humans and DNNs with weaker signals. Our results, based on 82,880 psychophysical trials under well-controlled lab conditions, demonstrate that there are still marked differences in the way humans and current DNNs process

⁶The no free lunch theorem [62] states that better performance on some input is necessarily accompanied by worse performance on other input; however we here are only interested in a very narrow subset of the possible input space—namely, natural images corrupted by distortions. The high accuracies of human observers across distortions indicate that it is, in principle, possible to achieve good performance on many distortions simultaneously.

object information. These differences, in our setting, cannot be overcome by training on distorted images (i.e., data augmentation): While DNNs cope perfectly well with the exact distortion they were trained on, they still show a strong generalisation failure towards previously unseen distortions. Since the space of possible distortions is literally unlimited (both theoretically and in real-world applications), it is not feasible to train on all of them. DNNs have a generalisation problem when it comes to settings that go beyond the usual (yet often unrealistic) i.i.d. assumption. We believe that solving this generalisation problem will be crucial both for robust machine inference and towards better models of human object recognition, and we envision that our findings as well as our carefully measured and freely available behavioural data⁷ may provide a new useful benchmark for improving DNN robustness and a motivation for neuroscientists to identify mechanisms in the brain that may be responsible for this remarkable robustness.

Author contributions

The initial project idea of comparing humans against DNNs was developed by F.A.W. and R.G. All authors jointly contributed towards designing the study and interpreting the data. R.G. and C.R.M.T. developed the image manipulations and acquired the behavioural data with input from H.H.S. and F.A.W.; J.R. trained networks on distortions; experimental data and networks were evaluated by C.R.M.T., R.G. and J.R. with input from H.H.S, M.B. and F.A.W.; R.G. and C.R.M.T. worked on making our work reproducible (data, code and materials openly accessible; writing supplementary material); R.G. wrote the paper with significant input from all other authors.

Acknowledgments

This work has been funded, in part, by the German Federal Ministry of Education and Research (BMBF) through the Bernstein Computational Neuroscience Program Tübingen (FKZ: 01GQ1002) as well as the German Research Foundation (DFG; Sachbeihilfe Wi 2103/4-1 and SFB 1233 on “Robust Vision”). The authors thank the International Max Planck Research School for Intelligent Systems (IMPRS-IS) for supporting R.G. and J.R.; J.R. acknowledges support by the Bosch Forschungsstiftung (Stifterverband, T113/30057/17); M.B. acknowledges support by the Centre for Integrative Neuroscience Tübingen (EXC 307) and by the Intelligence Advanced Research Projects Activity (IARPA) via Department of Interior/Interior Business Center (DoI/IBC) contract number D16PC00003.

We would like to thank David Janssen for his invaluable contributions in shaping the early stage of this project. Furthermore, we are very grateful to Tom Wallis for providing the MATLAB source code of one of his experiments, and for allowing us to use and modify it; Silke Gramer for administrative and Uli Wannek for technical support, as well as Britta Lewke for the method of creating response icons and Patricia Rubisch for help with testing human observers. Moreover, we would like to thank Nikolaus Kriegeskorte, Jakob Macke and Tom Wallis for helpful feedback, and three anonymous reviewers for constructive suggestions.

References

- [1] Irving Biederman. Recognition-by-components: a theory of human image understanding. *Psychological Review*, 94(2):115–147, 1987.
- [2] James J DiCarlo, Davide Zoccolan, and Nicole C Rust. How does the brain solve visual object recognition? *Neuron*, 73(3):415–434, 2012.
- [3] Mary C Potter. Short-term conceptual memory for pictures. *Journal of Experimental Psychology: human learning and memory*, 2(5):509, 1976.
- [4] Simon Thorpe, Denis Fize, and Catherine Marlot. Speed of processing in the human visual system. *Nature*, 381(6582):520–522, 1996.
- [5] Melvyn A. Goodale and A. David Milner. Separate visual pathways for perception and action. *Trends in Neurosciences*, 15(1):20–25, 1992.
- [6] A. Krizhevsky, I. Sutskever, and G. E. Hinton. ImageNet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems*, pages 1097–1105, 2012.

⁷<https://github.com/rgeirhos/generalisation-humans-DNNs>

- [7] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on ImageNet classification. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1026–1034, 2015.
- [8] David Silver, Aja Huang, Chris J. Maddison, Arthur Guez, Laurent Sifre, George van den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, Sander Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy Lillicrap, Madeleine Leach, Koray Kavukcuoglu, Thore Graepel, and Demis Hassabis. Mastering the game of Go with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016. ISSN 0028-0836.
- [9] Charles F. Cadieu, H. Hong, D. L. K. Yamins, N. Pinto, D. Ardila, E. A. Solomon, N. J. Majaj, and J. J. DiCarlo. Deep neural networks rival the representation of primate IT cortex for core visual object recognition. *PLoS Computational Biology*, 10(12), 2014.
- [10] Daniel LK Yamins, Ha Hong, Charles F Cadieu, Ethan A Solomon, Darren Seibert, and James J DiCarlo. Performance-optimized hierarchical models predict neural responses in higher visual cortex. *Proceedings of the National Academy of Sciences*, 111(23):8619–8624, 2014.
- [11] Radoslaw Martin Cichy, Aditya Khosla, Dimitrios Pantazis, and Aude Oliva. Dynamics of scene representations in the human brain revealed by magnetoencephalography and deep neural networks. *NeuroImage*, 2016.
- [12] Saeed Reza Kheradpisheh, Masoud Ghodrati, Mohammad Ganjtabesh, and Timothée Masquelier. Deep networks resemble human feed-forward vision in invariant object recognition. arXiv preprint arXiv:1508.03929, 2016.
- [13] Samuel Dodge and Lina Karam. A study and comparison of human and deep learning recognition performance under visual distortions. arXiv preprint arXiv:1705.02498, 2017.
- [14] Samuel Dodge and Lina Karam. Can the early human visual system compete with deep neural networks? arXiv preprint arXiv:1710.04744, 2017.
- [15] Ron Dekel. Human perception in computer vision. arXiv preprint arXiv:1701.04674, 2017.
- [16] RT Pramod and SP Arun. Do computational models differ systematically from human object perception? In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1601–1609, 2016.
- [17] Hamid Karimi-Rouzbahani, Nasour Bagheri, and Reza Ebrahimpour. Invariant object recognition is a personalized selection of invariant features in humans, not simply explained by hierarchical feed-forward vision models. *Scientific reports*, 7(1):14402, 2017.
- [18] Amir Rosenfeld, Markus D Solbach, and John K Tsotsos. Totally looks like-how humans compare, compared to machines. arXiv preprint arXiv:1803.01485, 2018.
- [19] Alban Flachot and Karl R Gegenfurtner. Processing of chromatic information in a deep convolutional neural network. *JOSA A*, 35(4):B334–B346, 2018.
- [20] Thomas SA Wallis, Christina M Funke, Alexander S Ecker, Leon A Gatys, Felix A Wichmann, and Matthias Bethge. A parametric texture model based on deep convolutional features closely matches texture appearance for humans. *Journal of vision*, 17(12):5–5, 2017.
- [21] Alexander Berardino, Valero Laparra, Johannes Ballé, and Eero Simoncelli. Eigen-distortions of hierarchical representations. In *Advances in Neural Information Processing Systems*, pages 3533–3542, 2017.
- [22] Kamila M Jozwik, Nikolaus Kriegeskorte, Katherine R Storrs, and Marieke Mur. Deep convolutional neural networks outperform feature-based but not categorical models in explaining object similarity judgments. *Frontiers in psychology*, 8:1726, 2017.
- [23] Brenden M Lake, Tomer D Ullman, Joshua B Tenenbaum, and Samuel J Gershman. Building machines that learn and think like people. *Behavioral and Brain Sciences*, 40, 2017.
- [24] Tim C Kietzmann, Patrick McClure, and Nikolaus Kriegeskorte. Deep neural networks in computational neuroscience. bioRxiv, 2017.
- [25] Rodney J Douglas and Kevan A C Martin. Opening the grey box. *Trends in Neurosciences*, 14(7):286–293, 1991.

- [26] George E P Box. Science and statistics. *Journal of the American Statistical Association*, 71(356):791–799, 1976.
- [27] Nikolaus Kriegeskorte. Deep neural networks: A new framework for modeling biological vision and brain information processing. *Annual Review of Vision Science*, 1(15):417–446, 2015.
- [28] Zhiyuan Chen and Bing Liu. *Lifelong Machine Learning*. Morgan & Claypool Publishers, 2016. ISBN 1627055010, 9781627055017.
- [29] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. arXiv preprint arXiv:1611.03530, 2016.
- [30] Kenji Kawaguchi, Leslie Pack Kaelbling, and Yoshua Bengio. Generalization in deep learning. arXiv preprint arXiv:1710.05468, 2017.
- [31] Behnam Neyshabur, Srinadh Bhojanapalli, David McAllester, and Nati Srebro. Exploring generalization in deep learning. In *Advances in Neural Information Processing Systems*, pages 5949–5958, 2017.
- [32] Ravid Shwartz-Ziv and Naftali Tishby. Opening the black box of deep neural networks via information. arXiv preprint arXiv:1703.00810, 2017.
- [33] Matthias Kümmerer, Thomas SA Wallis, and Matthias Bethge. Deepgaze II: Reading fixations from deep features trained on object recognition. arXiv preprint arXiv:1610.01563, 2016.
- [34] Hong-Wei Ng, Viet Dung Nguyen, Vassilios Vonikakis, and Stefan Winkler. Deep learning for emotion recognition on small datasets using transfer learning. In *Proceedings of the 2015 ACM on international conference on multimodal interaction*, pages 443–449. ACM, 2015.
- [35] Hayit Greenspan, Bram van Ginneken, and Ronald M Summers. Guest editorial deep learning in medical imaging: Overview and future promise of an exciting new technique. *IEEE Transactions on Medical Imaging*, 35(5):1153–1159, 2016.
- [36] Jeff Donahue, Yangqing Jia, Oriol Vinyals, Judy Hoffman, Ning Zhang, Eric Tzeng, and Trevor Darrell. Decaf: A deep convolutional activation feature for generic visual recognition. In *International conference on machine learning*, pages 647–655, 2014.
- [37] Sebastian Thrun. Is learning the n-th thing any easier than learning the first? In *Advances in Neural Information Processing Systems*, pages 640–646. The MIT Press, 1996.
- [38] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–9, 2015.
- [39] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556, 2015.
- [40] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016.
- [41] Robert Geirhos, David HJ Janssen, Heiko H Schütt, Jonas Rauber, Matthias Bethge, and Felix A Wichmann. Comparing deep neural networks against humans: object recognition when the signal gets weaker. arXiv preprint arXiv:1706.06969, 2017.
- [42] Jacob Nachmias and R V Sansbury. Grating contrast: Discrimination may be better than detection. *Vision Research*, 14(10):1039–1042, 1974.
- [43] Denis G Pelli and Bart Farell. Why use noise? *Journal of the Optical Society of America A*, 16(3):647–653, 1999.
- [44] Felix A Wichmann. *Some Aspects of Modelling Human Spatial Vision: Contrast Discrimination*. PhD thesis, The University of Oxford, 1999.
- [45] G Bruce Henning, C M Bird, and Felix A Wichmann. Contrast discrimination with pulse trains in pink noise. *Journal of the Optical Society of America A*, 19(7):1259–1266, 2002.
- [46] Matteo Carandini and David J Heeger. Normalization as a canonical neural computation. *Nature Reviews Neuroscience*, 13(1):51–62, 2012.
- [47] Matteo Carandini, David J Heeger, and J Anthony Movshon. Linearity and normalization in simple cells of the macaque primary visual cortex. *The Journal of Neuroscience*, 17(21):8621–8644, 1997.

- [48] Arnaud Delorme, Guillaume Richard, and Michele Fabre-Thorpe. Ultra-rapid categorisation of natural scenes does not rely on colour cues: a study in monkeys and humans. *Vision Research*, 40(16):2187–2200, 2000.
- [49] Felix A Wichmann, David HJ Janssen, Robert Geirhos, Guillermo Aguilar, Heiko H Schütt, Marianne Maertens, and Matthias Bethge. Methods and measurements to compare men against machines. *Electronic Imaging, Human Vision and Electronic Imaging*, 2017(14):36–45, 2017.
- [50] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015.
- [51] George A Miller. Wordnet: a lexical database for English. *Communications of the ACM*, 38(11):39–41, 1995.
- [52] Victor AF Lamme, Hans Super, and Henk Spekreijse. Feedforward, horizontal, and feedback processing in the visual cortex. *Current opinion in neurobiology*, 8(4):529–535, 1998.
- [53] Olaf Sporns and Jonathan D Zwi. The small world of the cerebral cortex. *Neuroinformatics*, 2(2):145–162, 2004.
- [54] Wulfram Gerstner. How can the brain be so fast? In J. Leo van Hemmen and Terrence J Sejnowski, editors, *23 Problems in Systems Neuroscience*, pages 135–142. Oxford University Press, 2005.
- [55] 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.
- [56] Felix A Wichmann, Doris I Braun, and Karl R Gegenfurtner. Phase noise and the classification of natural images. *Vision Research*, 46(8):1520–1529, 2006.
- [57] Jan Koenderink, Matteo Valsecchi, Andrea van Doorn, Johan Wagemans, and Karl Gegenfurtner. Eidolons: Novel stimuli for vision research. *Journal of Vision*, 17(2):7–7, 2017.
- [58] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser, M. Kudlur, J. Levenberg, D. Mane, R. Monga, S. Moore, D. Murray, C. Olah, M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viegas, O. Vinyals, P. Warden, M. Wattenberg, M. Wicke, Y. Yu, and X. Zheng. TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed Systems. arXiv preprint arXiv:1603.04467, 2016.
- [59] I. Vasiljevic, A. Chakrabarti, and G. Shakhnarovich. Examining the Impact of Blur on Recognition by Convolutional Networks. arXiv preprint arXiv:1611.05760, 2016.
- [60] Samuel Dodge and Lina Karam. Understanding how image quality affects deep neural networks. In *Quality of Multimedia Experience (QoMEX), 2016 Eighth International Conference on*, pages 1–6. IEEE, 2016.
- [61] Yiren Zhou, Sibong Song, and Ngai-Man Cheung. On classification of distorted images with deep convolutional neural networks. In *Acoustics, Speech and Signal Processing (ICASSP), 2017 IEEE International Conference on*, pages 1213–1217. IEEE, 2017.
- [62] David H Wolpert and William G Macready. No free lunch theorems for optimization. *IEEE transactions on evolutionary computation*, 1(1):67–82, 1997.
- [63] Vishal M Patel, Raghuraman Gopalan, Ruonan Li, and Rama Chellappa. Visual domain adaptation: A survey of recent advances. *IEEE signal processing magazine*, 32(3):53–69, 2015.
- [64] Heiko H Schütt and Felix A Wichmann. An image-computable psychophysical spatial vision model. *Journal of vision*, 17(12):12–12, 2017.
- [65] Kevin Jarrett, Koray Kavukcuoglu, Yann LeCun, et al. What is the best multi-stage architecture for object recognition? In *Computer Vision, 2009 IEEE 12th International Conference on*, pages 2146–2153. IEEE, 2009.
- [66] Mengye Ren, Renjie Liao, Raquel Urtasun, Fabian H Sinz, and Richard S Zemel. Normalizing the normalizers: Comparing and extending network normalization schemes. arXiv preprint arXiv:1611.04520, 2016.
- [67] R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2016. URL <https://www.R-project.org/>.

- [68] David H Brainard. The psychophysics toolbox. *Spatial Vision*, 10:433–436, 1997.
- [69] Mario Kleiner, David Brainard, Denis Pelli, Allen Ingling, Richard Murray, and Christopher Broussard. What’s new in Psychtoolbox-3. *Perception*, 36(14):1, 2007.
- [70] Eleanor Rosch. Principles of categorization. *Concepts: core readings*, 189, 1999.
- [71] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft COCO: Common objects in context. In *European Conference on Computer Vision*, pages 740–755. Springer, 2015.
- [72] Stefan Van der Walt, Johannes L Schönberger, Juan Nunez-Iglesias, François Boulogne, Joshua D Warner, Neil Yager, Emmanuelle Gouillart, and Tony Yu. scikit-image: image processing in Python. *PeerJ*, 2:e453, 2014.
- [73] A. M. Derrington, J. Krauskopf, and P. Lennie. Chromatic mechanisms in lateral geniculate nucleus of macaque. *The Journal of Physiology*, 357:241–265, 1984.
- [74] Andrew Stockman and Lindsay T Sharpe. The spectral sensitivities of the middle-and long-wavelength-sensitive cones derived from measurements in observers of known genotype. *Vision research*, 40(13): 1711–1737, 2000.
- [75] David H Brainard. *Human color vision*, chapter Cone Contrast and Opponent Modulation Color Spaces. Optical Society of America, Washington, DC, 2 edition, 1996.
- [76] A. van der Schaaf and J.H. van Hateren. Modelling the power spectra of natural images: Statistics and information. *Vision Research*, 36(17):2759 – 2770, 1996. ISSN 0042-6989. doi: [http://dx.doi.org/10.1016/0042-6989\(96\)00002-8](http://dx.doi.org/10.1016/0042-6989(96)00002-8).
- [77] Felix A Wichmann, Jan Drewes, Pedro Rosas, and Karl R Gegenfurtner. Animal detection in natural scenes: Critical features revisited. *Journal of Vision*, 10(4:6):1–27, 2010.
- [78] Yarín Gal and Zoubin Ghahramani. Dropout as a Bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059, 2016.

Supplementary material

While the main aspects of employed paradigm, procedure, observers and DNNs were already mentioned earlier, this section aims at providing exhaustive and reproducible experimental details. Furthermore, Figure 6 examines how network uncertainty develops as a function of signal strength, and Figure 7 shows the classification accuracy of networks trained on distortions across all conditions. All data, if not stated otherwise, were analysed using R version 3.2.3 [67].

Paradigm & procedure

A schematic of a typical trial is shown in Figure 5. Prior to starting the experiment, all participants were shown the response screen and asked to name all categories to ensure that the task was fully clear. They were instructed to click on the category that they thought resembles the image best, and to guess if they were unsure. They were allowed to change their choice within the 1500 ms response interval; the last click on a category icon of the response screen was counted as the answer. The experiment was not self-paced, i.e. the response screen was always visible for 1500 ms and thus, each experimental trial lasted exactly 2200 ms (300 ms + 200 ms + 200 ms + 1500 ms). During the whole experiment, the screen background was set to a grey value of 0.454 in the [0, 1] range, corresponding to the mean greyscale value of all images in the dataset (41.17 cd/m²).

On separate days we conducted twelve different experiments. The number of trials per experiment is reported in Table 1. For each experiment, we randomly chose between 70 and 80 images per category from the pool of images without replacement (i.e., no observer ever saw an image more than once throughout the entire experiment). Within each category, all conditions were counterbalanced. Random stimulus selection was done individually for each participant to reduce the influence of any accidental bias in the image selection. Images within the experiments were presented in randomised order. After 256 trials (colour, uniform noise and eidolon experiments), 128 trials (contrast experiment) and 160 trials (remaining experiments), the mean performance of the last block was displayed on the screen, and observers were free to take a short break. Ahead of each experiment, all observers conducted approximately 10 minutes of practice trials to gain familiarity with the task and the position of the categories on the response screen. Trials in which human observers failed to click on any category were recorded as an incorrect answer in the data analysis, and are shown as a separate category (top row) in the confusion matrices (DNNs, obviously, never fail to respond). Such a failure to respond occurred, on average, in only 1.91% of trials per experiment—one of the advantages of controlled laboratory studies ($SD = 0.69\%$).

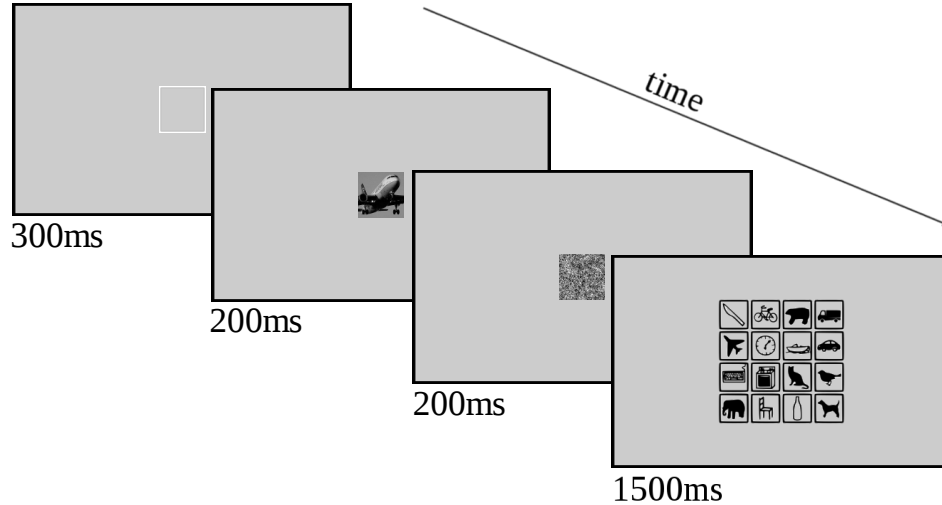


Figure 5: Schematic of a trial. After the presentation of a central fixation square (300 ms), the image was visible for 200 ms, followed immediately by a noise mask with $1/f$ spectrum (200 ms). Then, a response screen appeared for 1500 ms, during which the observer clicked on a category. Note that we increased the contrast of the noise mask in this figure for better visibility when printed. Categories row-wise from top to bottom: knife, bicycle, bear, truck, airplane, clock, boat, car, keyboard, oven, cat, bird, elephant, chair, bottle, dog. The icons are a modified version of the ones from the MS COCO website (<http://mscoco.org/explore/>).

Apparatus

All stimuli were presented on a VIEWPixx LCD monitor (VPixx Technologies, Saint-Bruno, Canada) in a dark chamber. The 22" monitor (484×302 mm) had a spatial resolution of 1920×1200 pixels at a refresh rate of 120 Hz. Stimuli were presented at the center of the screen with 256×256 pixels, corresponding, at a viewing distance of 123 cm, to 3×3 degrees of visual angle. A chin rest was used in order to keep the position of the head constant over the course of an experiment. Stimulus presentation and response recording were controlled using MATLAB (Release 2016a, The MathWorks, Inc., Natick, Massachusetts, U.S.) and the Psychophysics Toolbox extensions version 3.0.12 [68, 69] along with the iShow library⁸ on a desktop computer (12 core CPU i7-3930K, AMD HD7970 graphics card "Tahiti" by AMD, Sunnyvale, California, United States) running Kubuntu 14.04 LTS. Responses were collected with a standard computer mouse.

Observers

Three observers participated in the colour experiment (all male; 22 to 28 years; mean: 25 years) and in the contrast experiment. Six observers participated in the opponent colour, high-pass filter, low-pass filter, phase noise and power equalisation experiments (three female, three male; 20 to 25 years; mean: 22 years). In the other two experiments, five observers took part (uniform noise experiment: one female, four male; 20 to 28 years; mean: 23 years; eidolon experiments: three female, two male; 19 to 28 years; mean: 22 years). Subject-01 is an author and participated in all but the eidolon experiments. All other participants were either paid € 10 per hour for their participation or gained course credit. All observers were students and reported normal or corrected-to-normal vision.

Pre-trained deep neural networks

We used GoogLeNet [38], VGG-19 [39] and ResNet-152 [40] for our analyses. For all three networks, we used the pretrained implementations as provided by the TensorFlow-Slim framework⁹ and programmed in the TensorFlow library for machine learning [58]. The individual pretrained weights were also downloaded from the latter GitHub repository. We validated that our installation reproduced the classification accuracies provided on the website. The networks' input were 224×224 pixel RGB images. For greyscale images, we set all three channels to be equal to the greyscale image's single channel. Images were fed through the networks using a single feedforward pass.

⁸<http://dx.doi.org/10.5281/zenodo.34217>

⁹<https://github.com/tensorflow/models/tree/master/research/slim>, cloned on May 2, 2017

Table 1: Numbers of trials in the respective experiments. C. = conditions; P. = practice trials & blocks; M.= main experiment trials and blocks. The per condition column reports the number of trials per category and distortion level. The duration is reported without breaks.

Distortion type	C.	P. blocks	P. total	M. blocks	M. total	Per C.	Duration
Colour	2	2	320	5	1280	40	47 min
Uniform noise	8	2	256	5	1280	10	47 min
Contrast	8	2	256	10	1280	10	47 min
Eidolon I	8	4	384	5	1280	10	47 min
Eidolon II	8	4	384	5	1280	10	47 min
Eidolon III	8	4	384	5	1280	10	47 min
Opponent colours	2	2	224	7	1120	35	41 min
Low-pass filtering	8	2	256	8	1280	10	47 min
High-pass filtering	8	2	256	8	1280	10	47 min
Phase noise	7	2	224	7	1120	10	41 min
Power-equalisation	2	2	224	7	1120	35	41 min
Rotation	4	2	256	8	1280	20	47 min

Categories and image database

The images serving as psychophysical stimuli were images extracted from the training set of the ImageNet Large Scale Visual Recognition Challenge 2012 database [50]. This database contains millions of labeled images grouped into 1,000 very fine-grained categories (e.g., over a hundred different dog breeds). If human observers are asked to name objects, however, they most naturally categorise them into so-called basic or entry-level categories, e.g. dog rather than German shepherd [70]. The Microsoft COCO (MS COCO) database [71] is an image database structured according to 91 such entry-level categories, making it an excellent source of categories for an object recognition task. Thus for our experiments we fused the carefully selected entry-level categories in the MS COCO database with the large quantity of images in ImageNet. Using WordNet’s *hypernym* relationship (x is a hypernym of y if y is a “kind of” x , e.g., dog is a hypernym of German shepherd), we mapped every ImageNet label to an entry-level category of MS COCO in case such a relationship exists, retaining 16 clearly non-ambiguous categories with sufficiently many images within each category (see Figure 5 for a iconic representation of the 16 categories). A complete list of ImageNet labels used for the experiments can be found in our online repository.¹⁰ Since all investigated DNNs, when shown an image, output classification predictions for all 1,000 ImageNet categories, we disregarded all predictions for categories that were not mapped to any of the 16 entry-level categories. For each of those 16 categories we summed over the predictions of all ImageNet categories mapping to that particular entry-level category. Then the entry-level category with the highest summed prediction was selected as the network’s response. This way, the DNN response selection corresponds directly to the forced-choice paradigm for our human observers.

Image preprocessing and distortions

We used Python for all image preprocessing (Version 2.7.11) and for running experiments through pre-trained networks (Version 3.5). From the pool of ImageNet images of the 16 entry-level categories, we excluded all greyscale images (1%) as well as all images not at least 256×256 pixels in size (11% of non-greyscale images). We then cropped all images to a center patch of 256×256 pixels as follows: First, every image was cropped to the largest possible center square. This center square was then downsampled to the desired size with `PIL.Image.thumbnail((256, 256), Image.ANTIALIAS)`. Human observers get adapted to the mean luminance of the display during experiments, and thus images which are either very bright or very dark may be harder to recognise due to their very different perceived brightness. We therefore excluded all images which had a mean deviating more than two standard deviations from that of other images (5% of correct-sized colour-images excluded). In total we retained 213,555 images from ImageNet.

For the experiments using greyscale images the stimuli were converted using the `rgb2gray` method [72] in Python. This was the case for all experiments and conditions except for the ‘colour’ condition of the **colour experiment**, as well as for the opponent colour experiment. For the **contrast experiment**, we employed eight different contrast levels $c \in \{1, 3, 5, 10, 15, 30, 50, 100\%$. For an image in the $[0, 1]$ range, scaling the image to a new contrast level c was achieved by computing

$$new_value = \frac{c}{100\%} \cdot original_value + \frac{1 - \frac{c}{100\%}}{2}$$

¹⁰<https://github.com/rgeirhos/generalisation-humans-DNNs>

for each pixel. For the **uniform noise experiment**, we first scaled all images to a contrast level of $c = 30\%$. Subsequently, white uniform noise of range $[-w, w]$ was added pixelwise, $w \in \{0.0, 0.03, 0.05, 0.1, 0.2, 0.35, 0.6, 0.9\}$. In case this resulted in a value out of the $[0, 1]$ range, this value was clipped to either 0 or 1. By design, this never occurred for a noise range less or equal to 0.35 due to the reduced contrast (see above). For $w = 0.6$, clipping occurred in 17.2% of all pixels and for $w = 0.9$ in 44.4% of all pixels. See Figure 10 for example stimuli. For the **salt and pepper noise experiment**, used in DNN training experiments, we also scaled the greyscale image to a contrast level of 30% prior to adding noise in order to ensure maximal comparability with the uniform noise experiment. Salt and pepper noise, i.e. setting pixels to either black or white, was drawn pixelwise with a certain probability p , $p \in \{0, 10, 20, 35, 50, 65, 80, 95\}\%$. See Figure 14 for example salt-and-pepper stimuli at all conditions.

For the **opponent colours experiment**, our aim was to produce images that would be perceived by human observers as having exactly the opposite colours of the original, while retaining the same luminance. Therefore, we converted images to a colour space in which we could invert the colours without affecting luminance values. One such colour space is the Derrington-Krauskopf-Lennie (DKL) colour space [73]. In order to account for the nonlinearity of our experimental display monitor, we measured the emitted luminance for RGB grey values between 0 and 255. From this we built a lookup table from RGB grey values to actual emitted luminance values $f_{monitor}$. To evaluate how much the human retina’s long-, middle-, and short-wave receptors would be excited by the colours presented on the monitor, we measured the intensity of all emitted wave lengths between 390-780 nm for the RGB values (255 0 0), (0 255 0), (0 0 255), respectively. We then multiplied the respective emitted spectra between 390-780 nm with the corresponding cone sensitivities taken from the 2-deg LMS fundamentals proposed by [74] and summed over them. This resulted in a matrix C from RGB to cone activities (LMS space). Then we calculated a conversion matrix D of cone activities into the DKL colour space following the conversion example in [75]. An image was, consequently, converted from RGB to DKL by applying $f_{monitor}$ to it and subsequent multiplication with first C and then D . The DKL space has three channels reminiscent of the opponent colour process of the human visual system [75]. They are DKL_{lum} , a luminance channel, DKL_{L-M} , a channel representing the difference between long- and middle-wave receptor activation, as well as DKL_{S-lum} , a channel representing the difference between the activation of the short-wave receptor and the luminance. Since we wanted to keep the luminance unchanged, we multiplied the DKL_{L-M} and DKL_{S-lum} channels with the value ‘-1’. Subsequently, we converted the manipulated images back to RGB using the inverse matrices of D and C and then applied the inverse of $f_{monitor}$ to them. All resulting pixel values outside the range $[0, 1]$ were clipped to 0 or 1. This only happened for 0.34% of pixels with a mean clipped away value of 0.004. This corresponds to the minimal colour intensity step as $0.004 \approx \frac{1}{255}$.

For the low-pass and high-pass experiments we used the `scipy.ndimage.filters.gaussian_filter()` function. The **low-pass experiment**’s eight conditions differed in the standard deviation of the Gaussian filter. Standard deviations were 0 (original image), 1, 3, 7, 10, 15 and 40 pixels (Figure 11). We used constant padding with the mean pixel value over the testing images (0.4423) and truncation at four standard deviations. The **high-pass experiment** also had eight conditions. Standard deviations were 0.4, 0.45, 0.55, 0.7, 1, 1.5, 3 pixels and inf (original image) (Figure 11). The high-pass filtered images were produced by subtracting a low-pass filtered image as described above from the original image. However, many of the high-pass filtered images’ pixels fell outside the $[0, 1]$ range. To resolve this, we calculated the difference between the mean pixel value over all test images (0.4423) and the mean pixel value of the high-pass filtered image. That difference was added back to the image. This had the effect that images approached a uniform mean grey image of value 0.4423 for low standard deviations. For both experiments pixel values were clipped to the $[0, 1]$ range, if lying outside after the filtering. This only happened for <0.001% of pixels with a mean clipped away value of <0.001 for the both filtering experiments.

We implemented the equalisation of the power spectra and phase noise in the Fourier domain. Conversion to frequency domain was accomplished by a fast Fourier transform through the application of the `fft2()` and then `fftshift()` functions of the Python package `scipy.fftpack`. This results in a matrix of complex numbers F , which represents both the phases and amplitudes of the individual frequencies in one complex number. F is organised in symmetric pairs of complex numbers with just their imaginary part differing in its sign and cancelling each other out when reversing the Fourier transform again. When transforming F to polar coordinates, the angle represents the respective frequency’s phase and the distance from the origin represents its amplitude. Hence, we extracted the phases and amplitudes of the individual frequencies with the functions `numpy.angle(F)` and `numpy.abs(F)`, respectively. The **power equalisation experiment** had two conditions: original and power-equalised (Figure 13). For the power-equalised images, we first calculated the mean amplitude spectrum over all test images, which showed the typical $\frac{1}{f}$ shape [e.g. 76, 77]. Thereafter, we set all images amplitudes to the mean amplitude spectrum. Since the power spectrum is the square of the amplitude spectrum, the images were essentially power-equalised. There were seven conditions in the **phase noise experiment**. These were 0, 30, 60, 90, 120, 150 and 180 degrees noise width w (Figure 13). To each frequency’s phase a phase shift randomly drawn from a continuous uniform distribution over the interval $[-w, w]$ was added. To ensure that the imaginary parts would later cancel out again, we added the same phase noise to both frequencies of each symmetric pair. After performing the respective manipulations, a F_{new} was calculated by recombining the new phases and amplitudes. Then we did an inverse Fourier transform using `ifftshift()`

and then `ifft2()`. Finally we clipped all pixel values to the $[0, 1]$ range. This was the case for 0.038% of pixels with a mean clipped value of about 0.003 for the phase noise experiment and for 0.013% of pixels with a mean clipped value of 0.005 for the power-equalisation experiment.

There were four conditions for the **rotation experiment**: 0 (original), 90, 180 and 270 degrees rotation angle. Rotation by 90 degrees was accomplished by first transposing the image matrix and then reversing the column order. Rotation by 180 degrees was done by reversing both, row and column ordering. Rotation by 270 degrees was implemented by first reversing the images columns and then transposing it.

For the **eidolon experiments**, all stimuli were generated using the eidolon toolbox for Python¹¹, more specifically its `PartiallyCoherentDisarray(image, reach, coherence, grain)` function. Using a combination of the three parameters `reach`, `coherence` and `grain`, one obtains a distorted version of the original image (a so-called eidolon). The parameters `reach` and `coherence` were varied in the experiment; `grain` was held constant with a value of 10.0 throughout the experiment (`grain` indicates how fine-grained the distortion is; a value of 10.0 corresponds to a medium-grainy distortion). `Reach` $\in \{1.0, 2.0, 4.0, 8.0, 16.0, 32.0, 64.0, 128.0\}$ is an amplitude-like parameter indicating the strength of the distortion, `coherence` $\in \{0.0, 0.3, 1.0\}$ defines the relationship between local and global image structure. Those two parameters were fully crossed, resulting in $8 \cdot 3 = 24$ different eidolon conditions. A high coherence value “retains the local image structure even when the global image structure is destroyed” [57, p. 10]. A coherence value of 0.0 corresponds to ‘completely incoherent’, a value of 1.0 to ‘fully coherent’. The third value 0.3 was chosen because it produces images that perceptually lie—as informally determined by the authors—in the middle between those two extremes. See Figures 12 and 13 for example eidolon stimuli. The coherence levels of 1.0, 0.3 and 0.0 are referred to as eidolon experiment I, II and III throughout the paper.

Experimental modifications

Our psychophysical experiments were conducted in two batches and over an extended period of time. After completing the first batch of experiments (all experiments on the left half of Figure 3, i.e. a, c, e, g, i and k), we performed a number of modifications for the second batch of experiments. We here briefly list all the changes in which the second batch of experiments differed from the previously reported methods.

Noise mask: In the human experiments, each experimental image was immediately followed by a $\frac{1}{7}$ pink noise mask (cf. Figure 5). In the second batch of experiments this noise mask was enhanced to improve its masking effect. This was done by multiplying each pixel value by four. Values greater than 1 or smaller than 0 due to the multiplication were then clipped to 1 or 0.

Cropping vs. downsampling: In the first experimental batch, humans saw 256×256 images. However, DNN classification was based on those images 224×224 centre crop. Thus, humans and DNNs saw slightly different images. Therefore, we used downsampling to 224×224 for the second batch of both, human and DNN experiments. As a consequence, the mean grey pixel value over all experimental images and hence the background grey value for presenting those images changed slightly from 0.454 to 0.442 in the $[0, 1]$ range.

JPEG vs. PNG: All images of the first batch of experiments, prior to showing them to human observers or DNNs, were saved in the JPEG format using the default settings of the `skimage.io.imsave` function. The JPEG format was chosen because the image training database for all three networks, ImageNet [50], consists of JPEG images. However, as the lossy compression of JPEG may introduce artefacts, we also examined the difference in DNN results between saving to JPEG and to PNG, which is lossless up to rounding issues. We therefore ran all those DNN experiments additionally saving them in the (up to rounding issues) lossless PNG format. We did not find any noteworthy differences for the colour, noise, and eidolon experiments. However, for the contrast experiment, the networks achieved on average better results for PNG images. We therefore tested three human observers additionally on the same stimuli (PNG instead of JPEG images). In this experiment, three of the JPEG experiment’s five observers participated for maximal comparability.¹² We found human observers to be better for PNG images as well. In absolute terms, participants were 2.68% better on average. In order to disentangle the influence of JPEG compression and image manipulations, we used PNG images for all other experiments, that is for the false colour, phase noise, power equalisation, rotation, high-pass and low-pass experiments as well as for the DNN training experiments.

Python Version: The second batch used Python Version 3.5 instead of Python 2.7 for image preprocessing.

¹¹<https://github.com/gestaltrevision/Eidolon>

¹²A time gap of approximately six months between both experiments should minimise memory effects; furthermore, human participants were not shown any feedback (correct / incorrect classification choice) during the experiments.

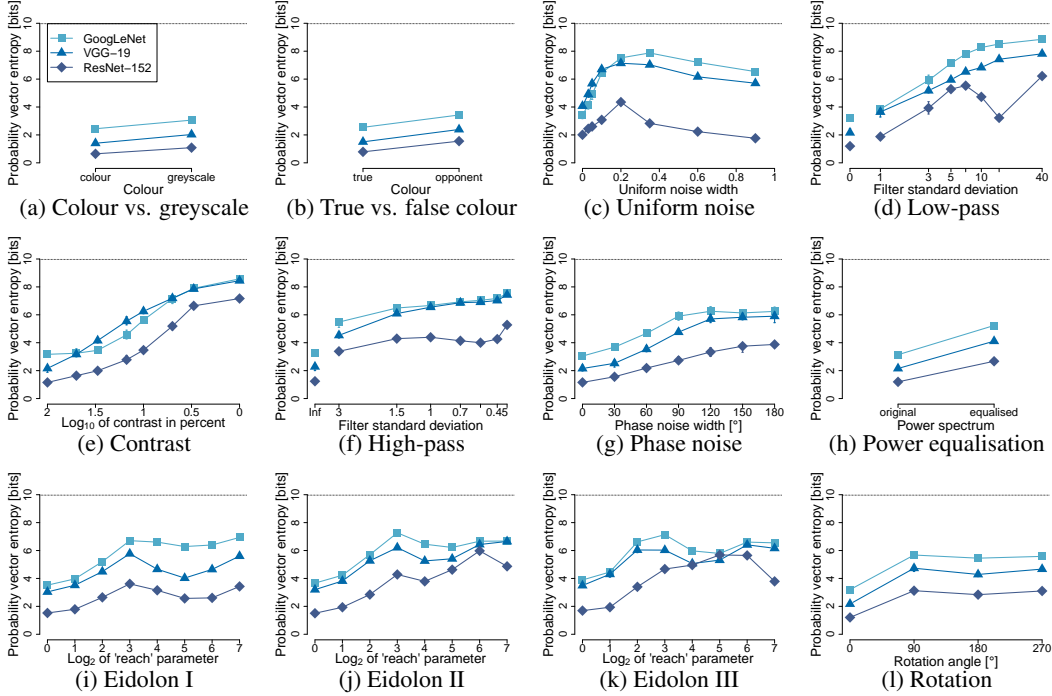


Figure 6: Mean entropy of the probabilities for the 1000 ILSVRC classes for [GoogLeNet](#), [VGG-19](#) and [ResNet-152](#). Dotted line indicates the maximum possible entropy. This is a measure of network ‘uncertainty’.

Error bars & entropy

When showing *accuracy* in any of the plots, the error bars provided report the range of the data observed for different observers (not the often shown S.E. of the means, which would be much smaller). To produce a comparable measure of uncertainty for the DNNs, we computed seven runs with different subsets of the data, with each run consisting of the same number of images per category and condition that a single human observer was exposed to and report the range of accuracies observed in these runs. Seven runs are the maximum possible number of runs without ever showing an image to a DNN more than once per experiment.

For all response distribution entropy results (Figures 3 and 7), we calculated the entropy as the average of individual participants’ entropies: otherwise, if the entropy was calculated over the aggregated human trials, individual differences might cancel each other out, which would lead to a higher human response distribution entropy.

Prediction uncertainty

Figure 6 shows the entropy of the networks’ predictions over the 1000 ILSVRC12 classes as a measure of the networks’ ‘uncertainty’. In principle, the more uncertain a network is in its predictions the more evenly it will distribute its softmax activations between classes and thus the higher the entropy will be. For all experiments, uncertainty roughly increases with distortion strength as reported by previous studies [59, 60]. However, for uniform noise and the eidolon distortions all networks become more certain again for some higher distortion levels. Furthermore, ResNet-152 also becomes more certain for stronger distortions in the low-pass experiment and is consistently more confident in its predictions than GoogleNet and VGG-19. Thus there are distortions for which all networks and especially ResNet-152 fail to capture that the input signal is becoming worse. Instead they limit their predictions to only a few classes (cf. Figure 3) with high certainty. This result is in line with previous reports stating that uncertainty in standard discriminative deep neural networks is not well represented [e.g. 78].

Network training results across all distortion conditions

While Figure 4 shows performance of networks trained on distortions for a single distortion level per manipulation, we here report the performance across all stimulus levels. In Figure 7, the performance of a vanilla ResNet-50

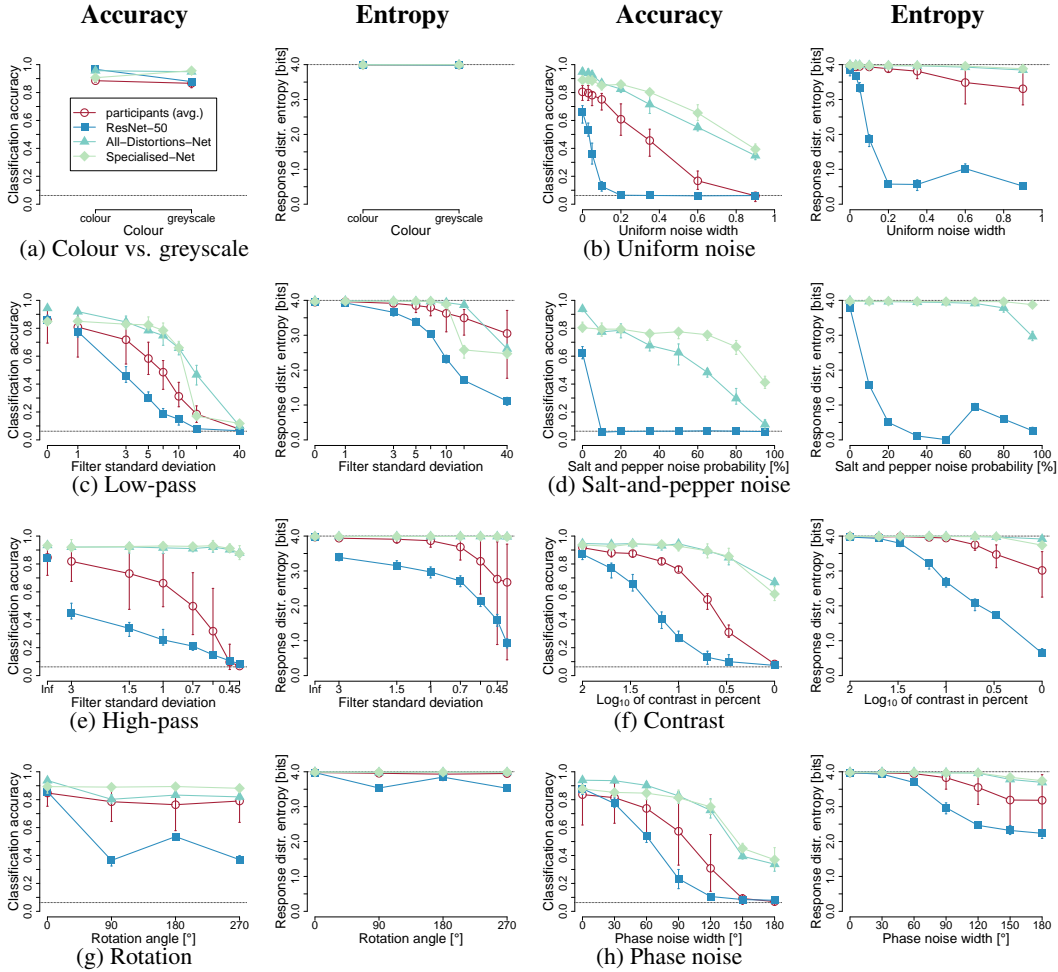


Figure 7: Classification accuracy and response distribution entropy for **human observers**, **ResNet-50** as well as an **All-Distortions-Net** and a **Specialised-Net**. All networks are trained from scratch; the Specialised-Net in every plot is trained on a single distortion (models A1 to A9 in Figure 4) whereas the All-Distortions-Net is trained on a number of distortions simultaneously. This corresponds to models C1 and C2 in Figure 4: for subplot 7d, salt-and-pepper noise, performance of model C1 is shown. For subplot 7b, uniform noise, performance of C2 is shown. For all other plots, performance of the All-Distortions-Net is shown as the mean of performance for models C1 and C2.

is compared against a network with the same architecture that is trained on a distortion directly (referred to as ‘Specialised-Net’), as well as to a network that is trained on all distortions simultaneously (named ‘All-Distortions-Net’). Object recognition accuracy shows a relatively consistent pattern across experiments: human performance is better than the performance of a vanilla ResNet-50. However, both an All-Distortions-Net and a Specialised-Net reach extremely high accuracies, with the Specialised-Net being either on par with or slightly better than the All-Distortions-Net. Interestingly, the response distribution entropy of those two networks is largely human-like, i.e. close to 4 bits of entropy (or no bias towards a certain category), even for conditions where the overall accuracy is low (e.g. for the difficult conditions of uniform and salt-and-pepper noise).

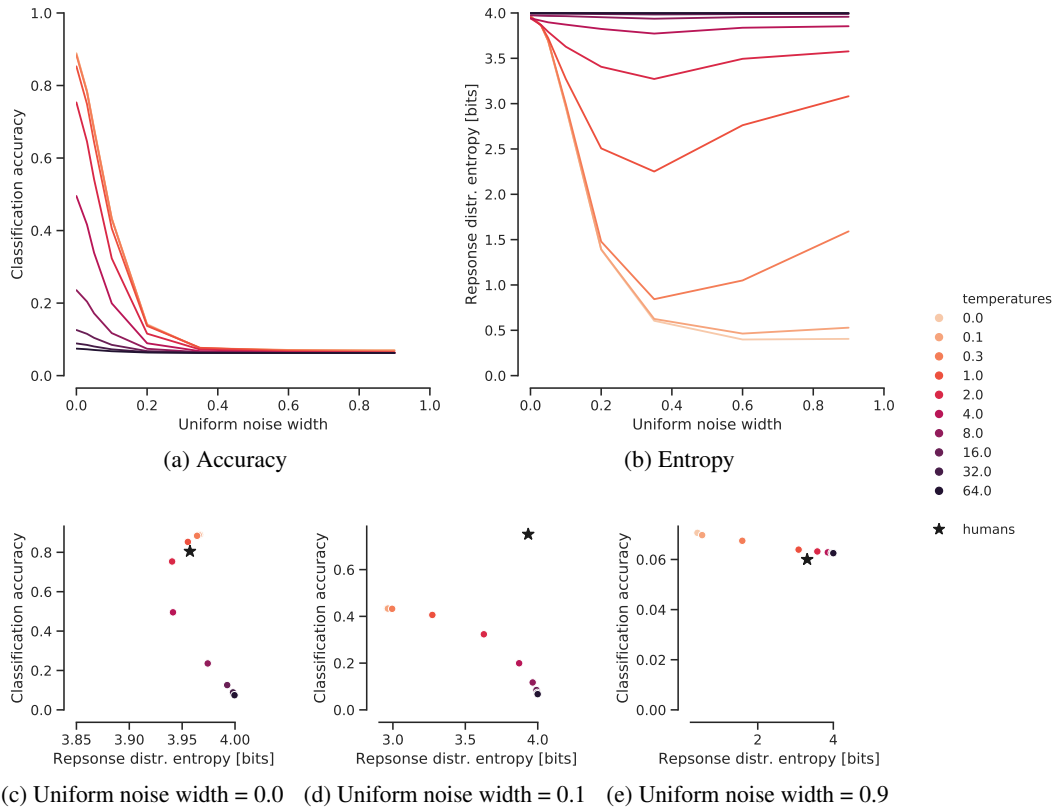


Figure 8: Classification accuracy **(a)** and response distribution entropy **(b)** as well as the trade-off between accuracy and entropy **(c, d, e)** for different softmax temperatures when the decision of a ResNet-50 model is sampled from its distribution over classes (softmax output) rather than taking the argmax of the distribution (which is equivalent to sampling with temperature $\rightarrow 0$). While increasing the temperature does increase the response distribution entropy of ResNet-50, it simultaneously decreases the classification accuracy. For uniform noise with a width of 0.1 **(d)**, increasing the temperature to match the response distribution entropy of humans reduces the accuracy of ResNet-50 below 0.1 whereas human accuracy is at 0.75.

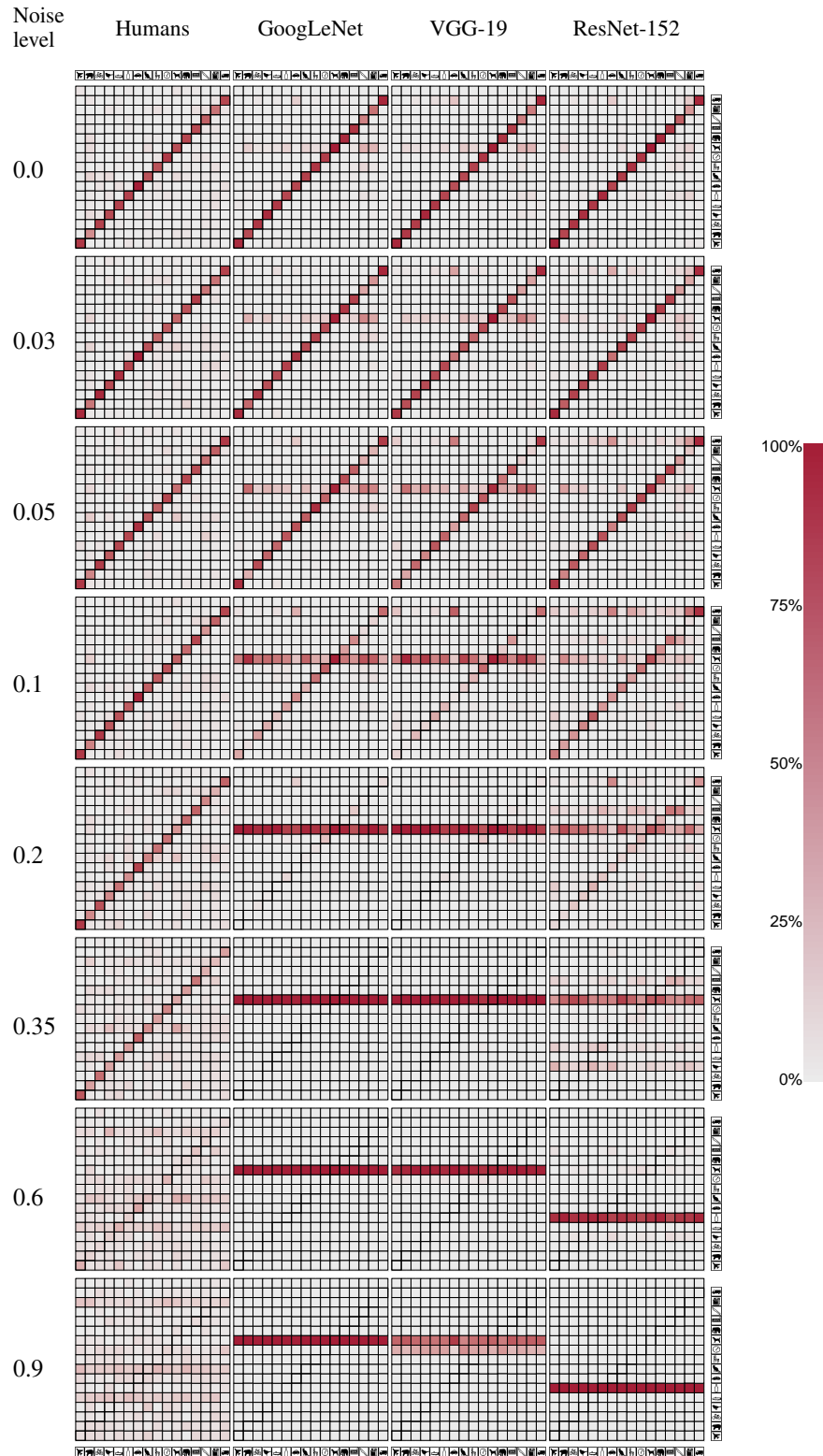


Figure 9: Confusion matrices for additive uniform noise. Columns within a confusion matrix indicate correct category, rows the given response / classification decision. Responses on the diagonal indicate correct responses. The top row of each confusion matrix indicates the fraction of failures to respond (i.e. if human observers failed to click on a category). Within each column, the colours show how often humans / DNNs responded with a certain category in percent. ‘Noise level’ indicates the bounds of additive uniform noise as visualised in Figure 10.

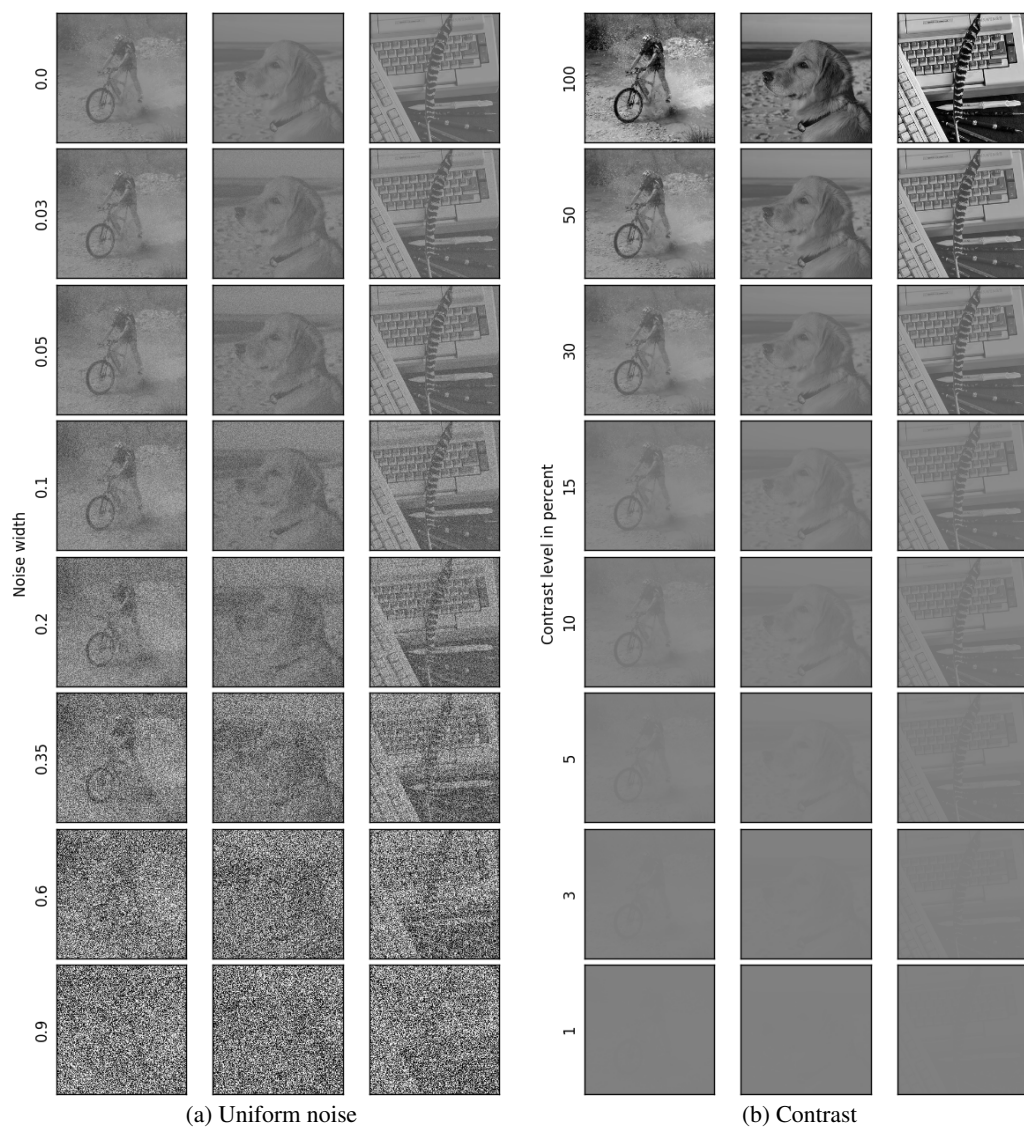


Figure 10: Three example stimuli for different conditions of uniform noise and contrast experiments. The three images (categories bicycle, dog and keyboard) were drawn randomly from the pool of images used in the experiments. Best viewed on screen.



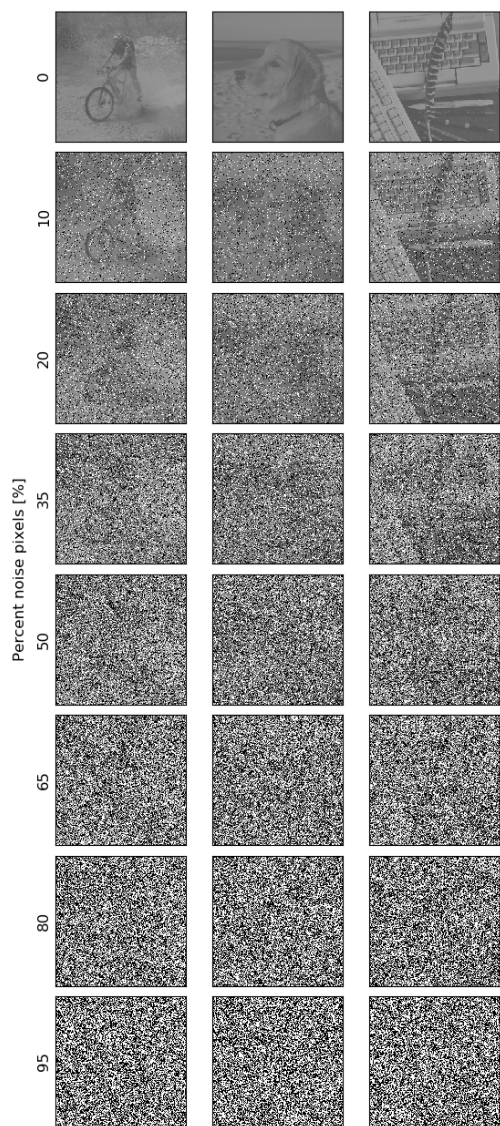
Figure 11: Three example stimuli for different conditions of low-pass and high-pass experiments. The three images (categories bicycle, dog and keyboard) were drawn randomly from the pool of images used in the experiments. Best viewed on screen.



Figure 12: Three example stimuli for different conditions of eidolon I and eidolon II experiments. The three images (categories bicycle, dog and keyboard) were drawn randomly from the pool of images used in the experiments. Best viewed on screen.



Figure 13: Three example stimuli for different conditions of Eidolon III, phase noise, false colour and power equalisation experiments. The three images (categories bicycle, dog and keyboard) were drawn randomly from the pool of images used in the experiments. Best viewed on screen.



(a) Salt and pepper noise

Figure 14: Three example stimuli for different conditions of salt and pepper noise. The three images (categories `bicycle`, `dog` and `keyboard`) were drawn randomly from the pool of images used in the experiments. Best viewed on screen. Salt and pepper noise was used in DNN training experiments with the conditions depicted in the figure above.