

# Generalisation in humans and deep neural networks

NIPS 2018

Presented by Zubia Mansoor

August 20, 2020

# Overview

- 1 Motivation
- 2 Methods
- 3 Experiments
- 4 Results
- 5 Summary
- 6 Questions?

Motivation

## Motivating Example

Consider a motivating scenario in radiology



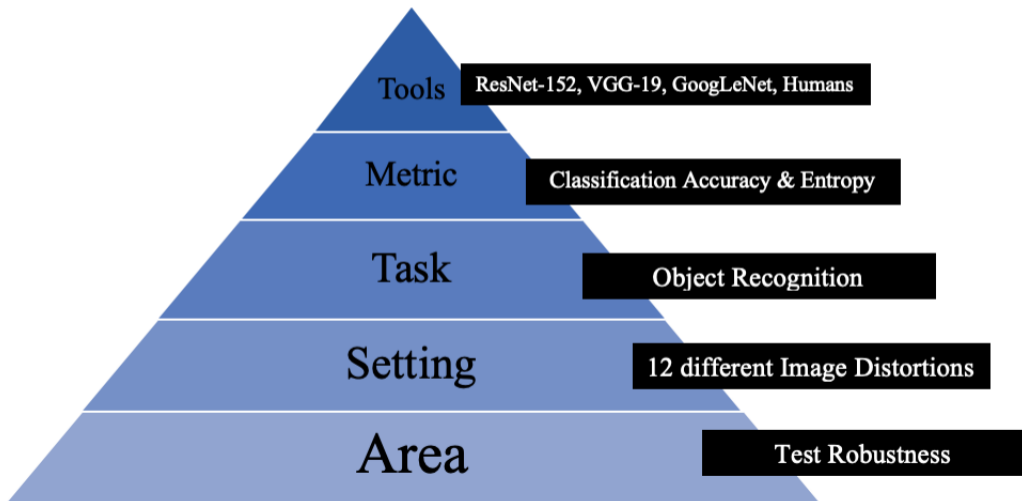
- What happens when the model encounters something it hasn't seen before?
  - For instance, if the X-Ray copies are blurry and noisy
- Changes in the training and test distribution pose a serious challenge to deep learning vision systems

### General

- Conduct a behavioral comparison of human and DNNs to test for robustness on image distortions

### Specific

- Test the performance of DNNs on the *exact distortion types they were trained on*
- Test the performance of DNNs on *previously unseen distortions*



Methods

- Controlled lab experiment
- ImageNet data
- Design Challenge 1
  - Problem: DNNs are usually trained for fine-grained categories as compared to humans
  - Solution: Distill over 20,000 categories of ImageNet into 16 entry-level categories using WordNet
- Design Challenge 2
  - Problem: DNNs only use feedforward computation versus recurrent connections used by human brain
  - Solution: Each image is followed by a 200 ms presentation of a pink noise mask



## Types of Distortion

Name	Data type
Colour vs. greyscale	Dichotomous
True vs. false colour	Dichotomous
Power equalisation	Dichotomous
Rotation	4 levels
Phase noise	7 levels
Uniform noise	8 levels
Low pass	8 levels
High pass	8 levels
Contrast	100% to 0%
Eidolon I,II,III	8 levels
Salt-and-Pepper	100% to 0%

## Types of Distortion Contd.

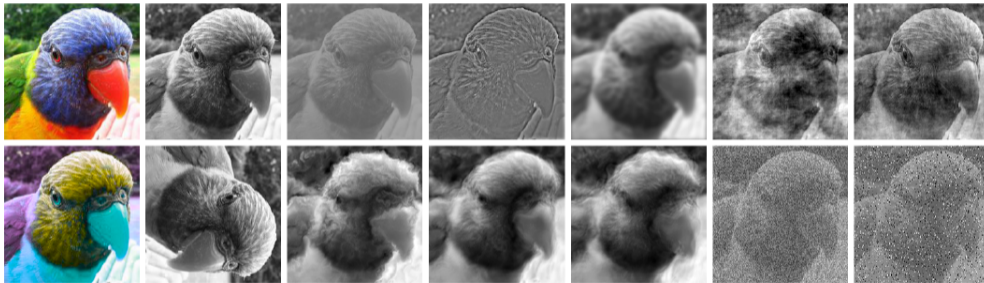


Figure: Sample image of class *bird* across all distortion types

Experiments

### Goals

- How well do DNNs generalise as compared to humans on previously unseen distortions?

### Goals

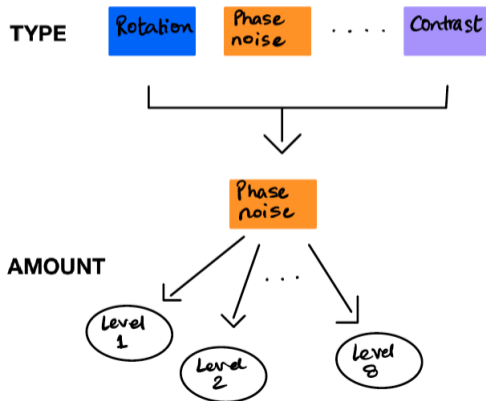
- How well do DNNs generalise as compared to humans when trained directly on the same distortions?
- To what extent does training on one distortion type helps the network cope with other distortion types?

## More on the training process...

- Experiment I
  - Standard pre-trained DNNs
- Experiment II
  - Trained on a subset of the standard ImageNet dataset (16-class-ImageNet)
  - Weight each sample in the loss function to correct for imbalanced classes
  - ResNet 50- like architecture

## More on the training process...(1)

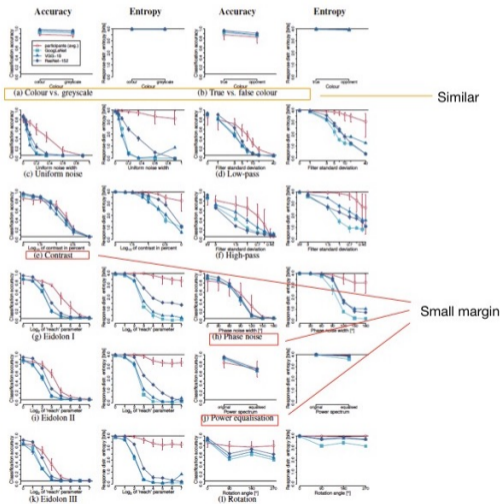
- Experiment II
  - Type and amount of distortion selected uniformly



Results



# Results: Experiment I



## Findings

- Humans and DNNs exhibit *similar performance* on colour-related distortions
- Human appear to be more robust *by a small margin* for low-contrast, power equalisation and phase noise
- Human appear to be more robust *by a large margin* for uniform noise, low-pass, high-pass, rotation and all eidolon experiments

## Results: Experiment II

Evaluation condition	Model																				
	human observers	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	94.1	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: Classification Accuracy (in %)

- **Models A1-A9:** Training on a certain distortion improves the performance greater for that same distortion but only slightly for other distortions
- **Models B1-B9:** Training on a specific distortion combined with uniform noise improves more compared to models A1 to A9
- **Models C1-C2:** Training either without uniform or salt-and-pepper noise leads to poor and closer to chance level performance

Summary

- Human visual system appears to be more robust than DNNs for the most part
- Diverging classification error-patterns between humans and DNNs as signal weakens
- DNNs surpass human performance only when trained on the exact distortions type they are later tested on
- Benchmark dataset of 83,000 human psychophysical trials

Questions?