

# Generalisation report

## Generalisation report produced by model-vs-human

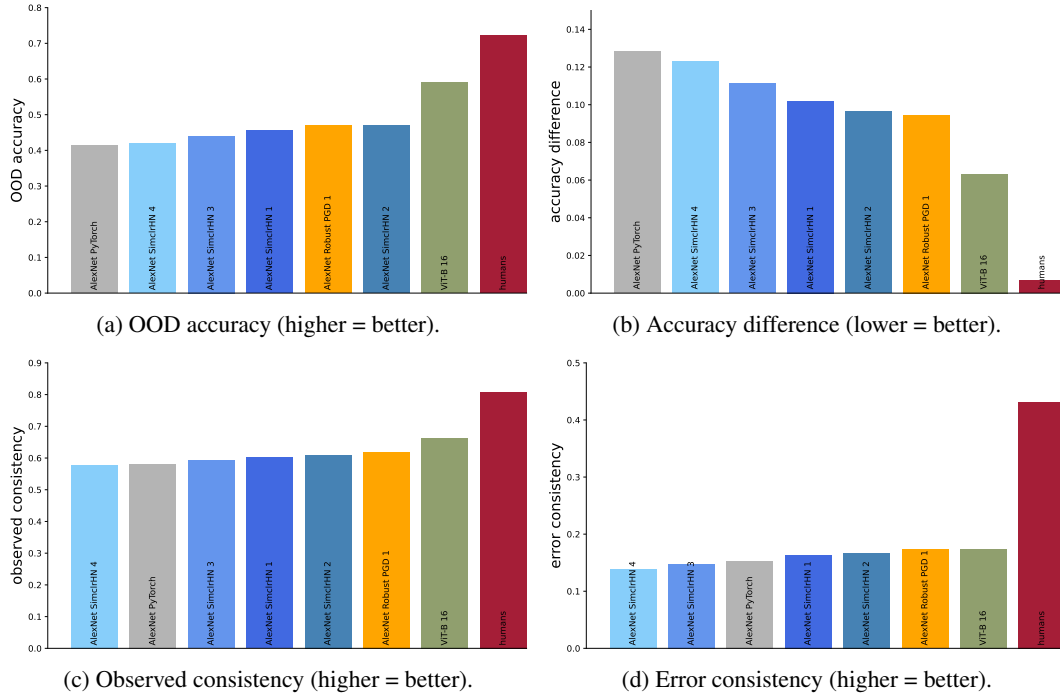


Figure 1: Benchmark results for different models, aggregated over datasets.

Table 1: Benchmark table of model results for most human-like behaviour. The three metrics “accuracy difference” “observed consistency” and “error consistency” (plotted in Figure ??) each produce a different model ranking. The mean rank of a model across those three metrics is used to rank the models on our benchmark.

model	accuracy diff. ↓	obs. consistency ↑	error consistency ↑	mean rank ↓
ViT-B 16	<b>0.063</b>	<b>0.664</b>	<b>0.175</b>	<b>1.000</b>
AlexNet Robust PGD 1	0.095	0.619	0.173	2.000
AlexNet SimclrHN 2	0.096	0.610	0.167	3.000
AlexNet SimclrHN 1	0.102	0.602	0.163	4.000
AlexNet SimclrHN 3	0.111	0.592	0.148	5.333
AlexNet PyTorch	0.128	0.581	0.153	6.000
AlexNet SimclrHN 4	0.123	0.579	0.138	6.667

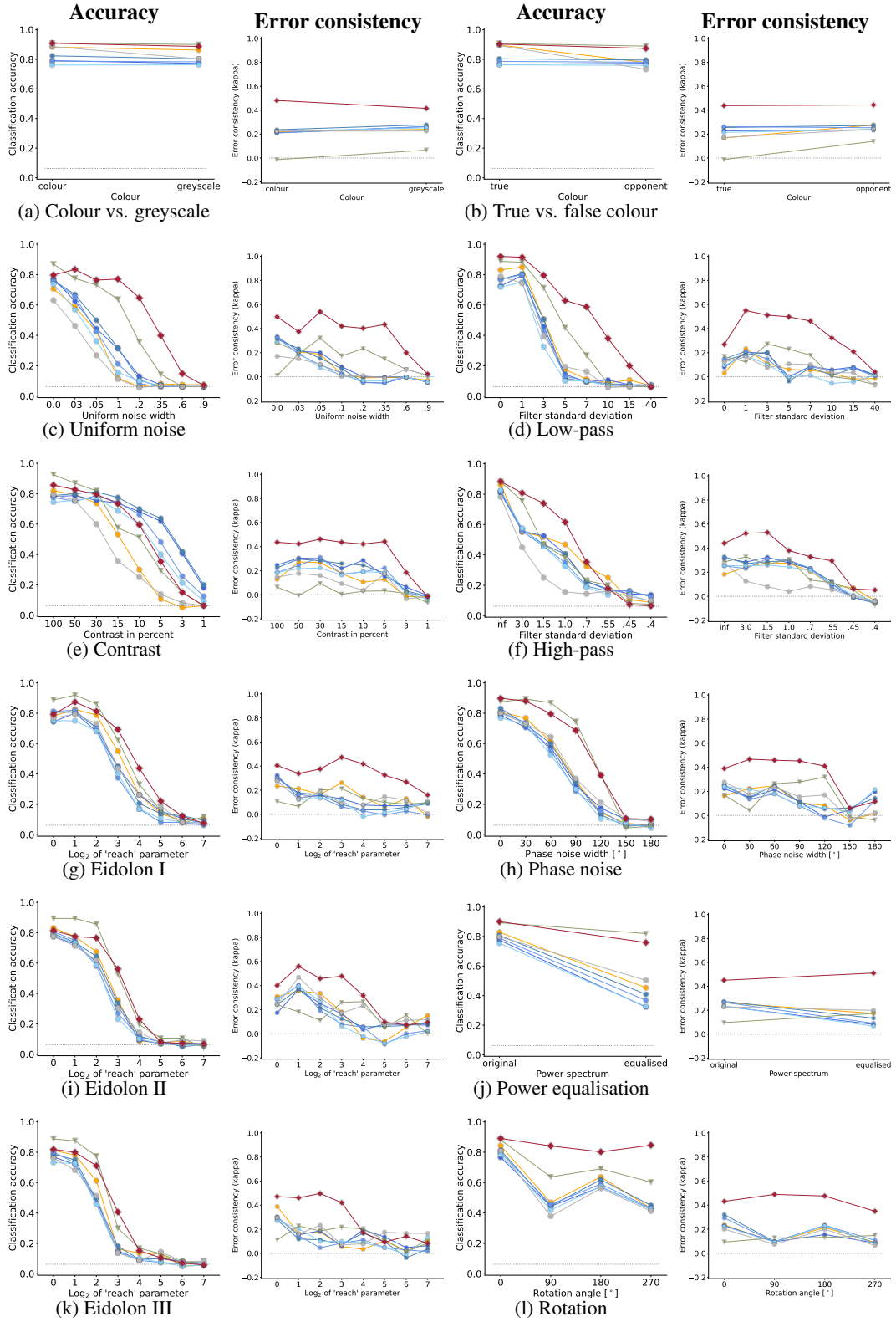


Figure 2: OOD accuracy and error consistency.

Table 2: Benchmark table of model results for highest out-of-distribution robustness.

model	OOD accuracy $\uparrow$	rank $\downarrow$
ViT-B 16	<b>0.590</b>	<b>1.000</b>
AlexNet SimclrHN 2	0.470	2.000
AlexNet Robust PGD 1	0.470	3.000
AlexNet SimclrHN 1	0.455	4.000
AlexNet SimclrHN 3	0.440	5.000
AlexNet SimclrHN 4	0.420	6.000
AlexNet PyTorch	0.415	7.000

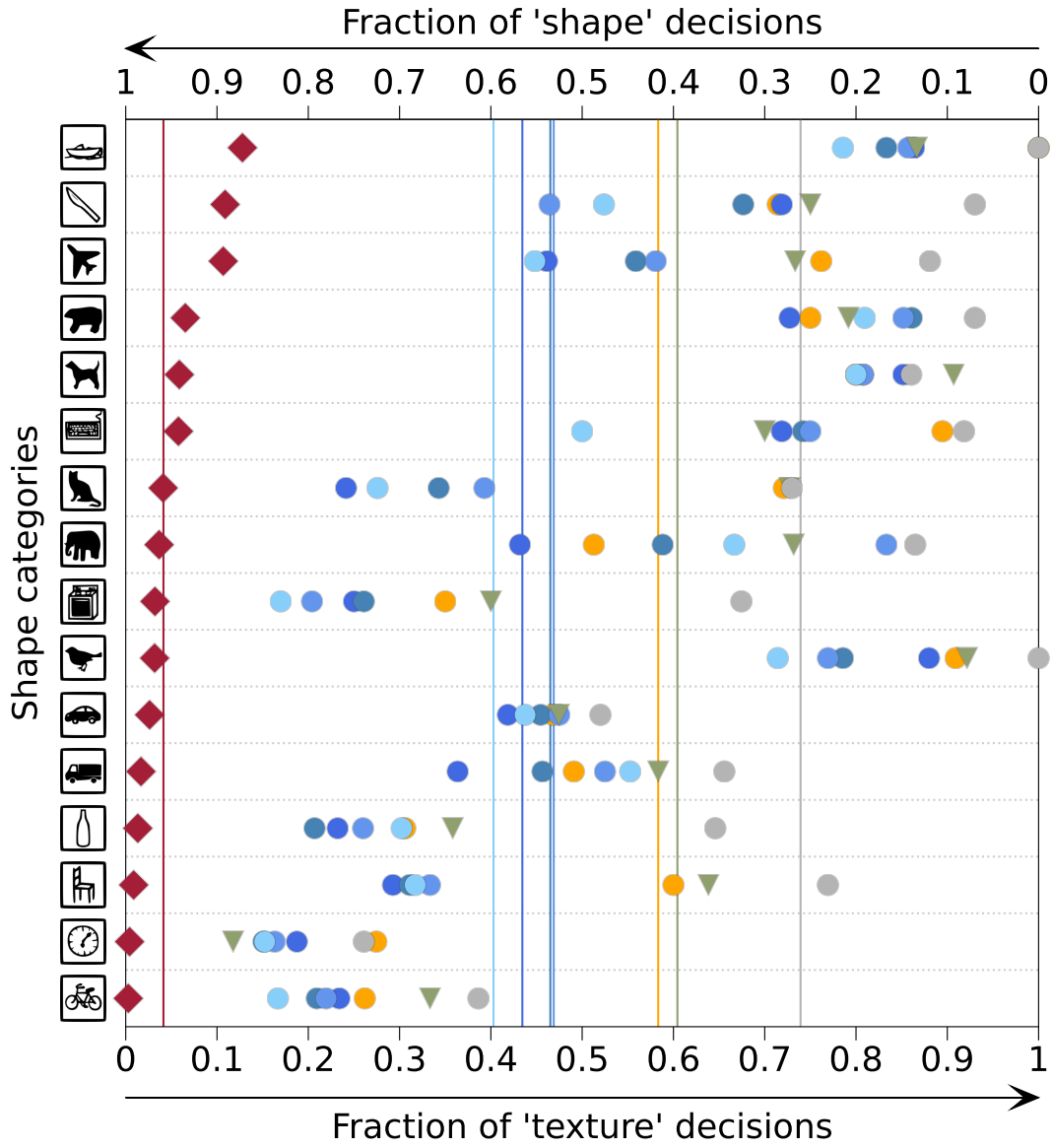


Figure 3: Shape vs. texture bias: category-level plot.

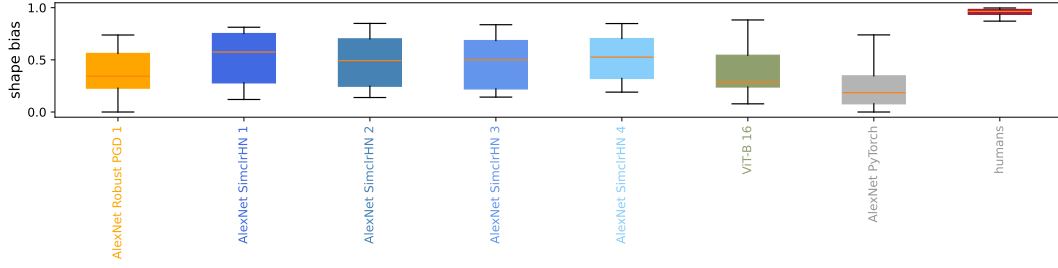


Figure 4: Shape vs. texture bias: boxplot.

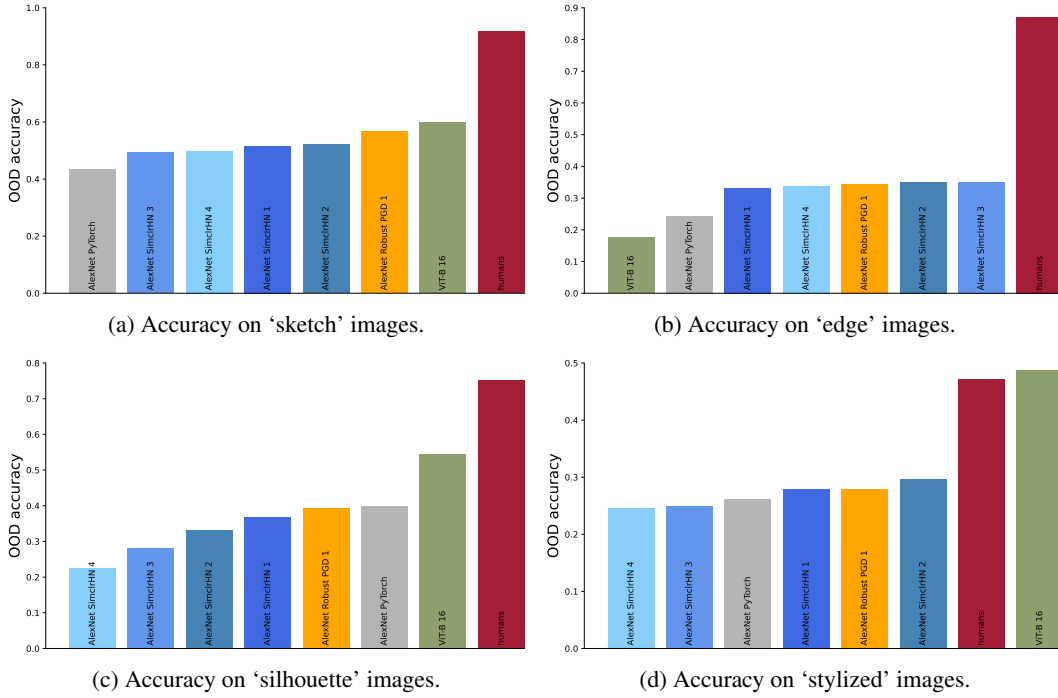


Figure 5: OOD accuracy on four nonparametric datasets (i.e., datasets with only a single corruption type and strength).

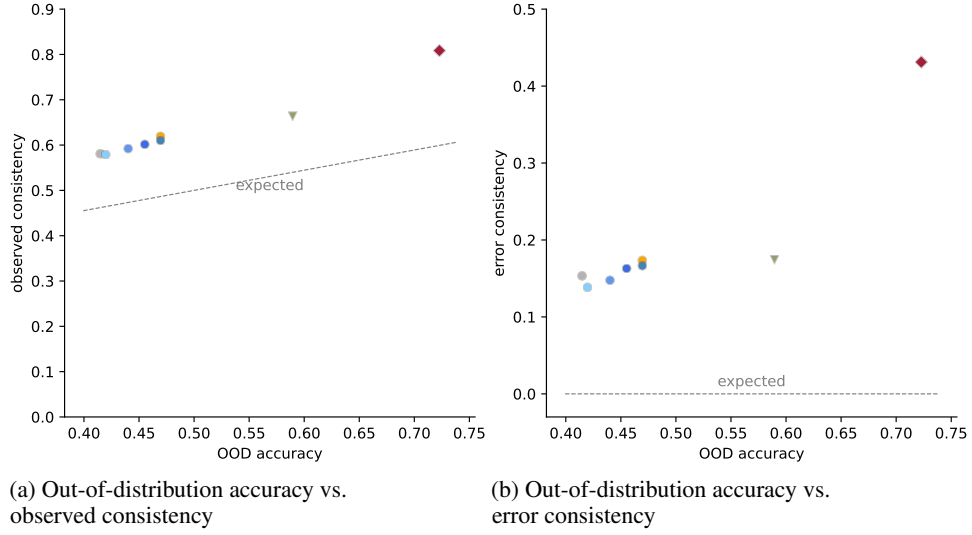


Figure 6: Observed consistency and error consistency between models and humans as a function of out-of-distribution (OOD) accuracy. Dotted lines indicate consistency expected by chance.

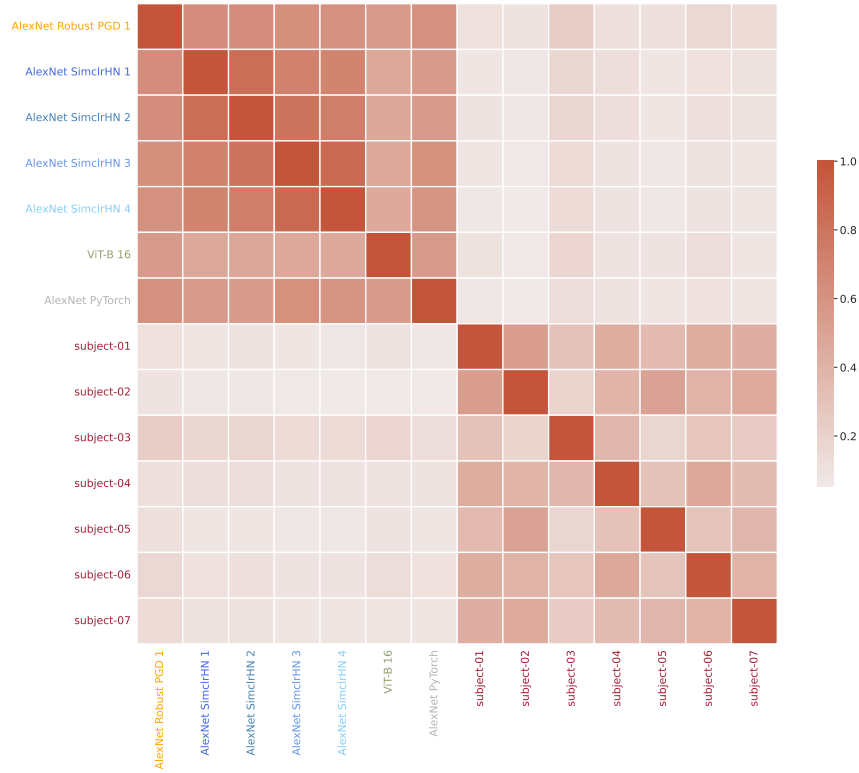


Figure 7: Error consistency for 'sketch' images.

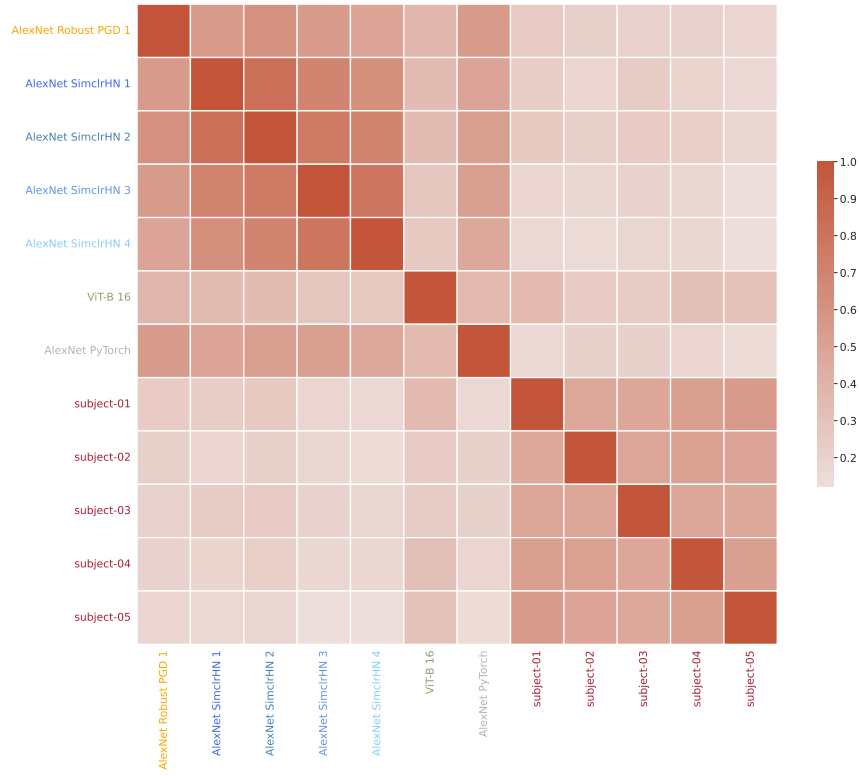


Figure 8: Error consistency for ‘stylized’ images.

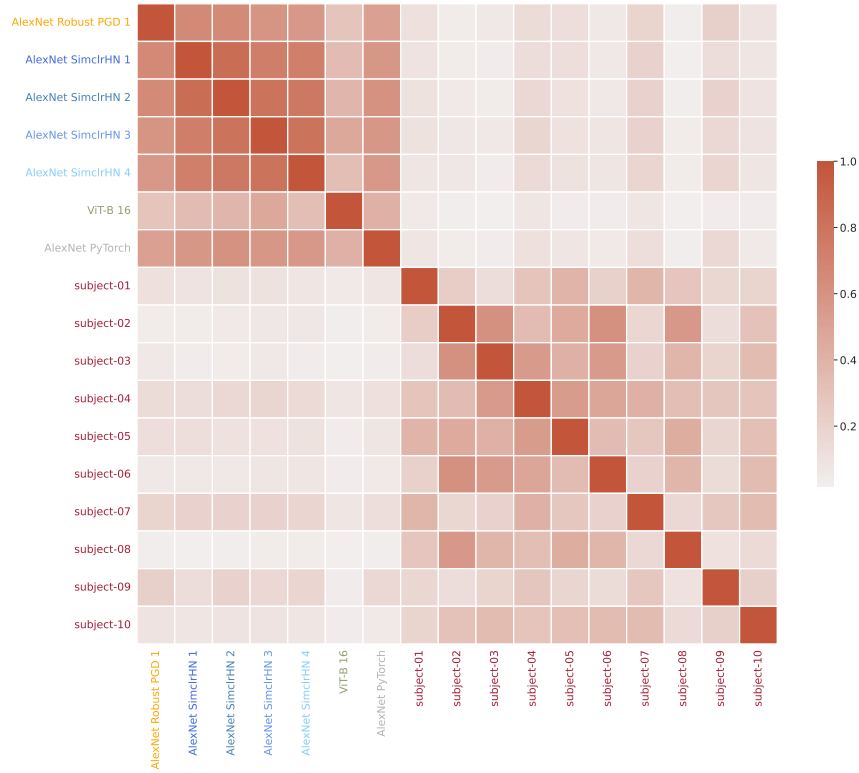


Figure 9: Error consistency for ‘edge’ images.

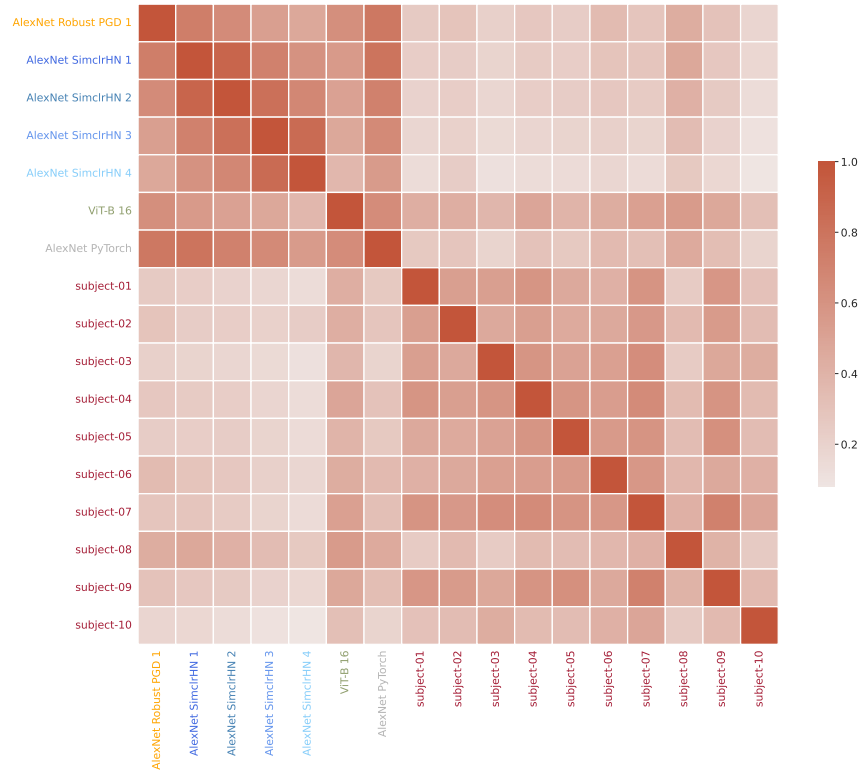


Figure 10: Error consistency for ‘silhouette’ images.

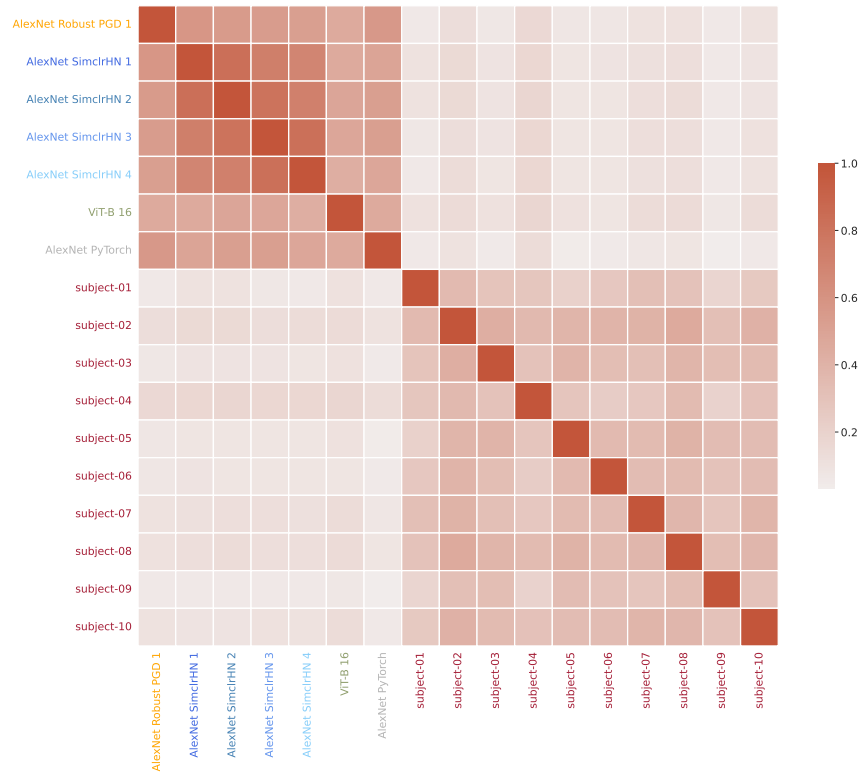


Figure 11: Error consistency for ‘cue conflict’ images.