

VisAlign: Dataset for Measuring the Alignment between AI and Humans in Visual Perception

Jiyoung Lee¹, Seungho Kim¹, Seunghyun Won², Joonseok Lee³, Marzyeh Ghassemi^{4,5,6}

James Thorne¹, Jaeseok Choi⁷, O-Kil Kwon⁷, Edward Choi¹

¹KAIST

²Seoul National University Bundang Hospital

³Seoul National University

⁴MIT

⁵University of Toronto

⁶Vector Institute

⁷Kangwon National University Hospital

Motivation

Why Alignment is Important in AI?

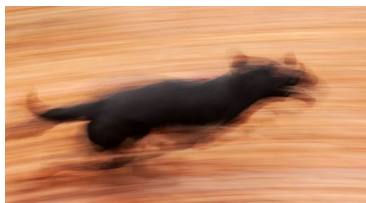
- Safety is a critical issue in AI which might cause tremendous costs.
- Ensuring deep learning safety is difficult because there is little manual control of feature interaction.
- In this project, we will evaluate **alignment** as a proxy measure for reliability.
 - Well-aligned models induce more agreeable and acceptable results.

Introduction

AI-Human Visual Alignment

- We focus on visual perception, *“AI-human visual alignment”*.

Q. Is this a dog or a cat?



Dog

Abstain

Between Dog and Abstain

Work

- Dataset: VisAlign
 - total of 8 sub-categories
 - reflects the various scenarios that can happen in the real world

- Metrics
 - Visual Alignment Metric
 - Reliability Score

VisAlign: AI-Human Visual Alignment Dataset

Must-Act

Category 1

Ordinary, generic pictures of mammals

Sources: ImageNet, images.cv



Category 2

Spurious correlations between the mammal and background

Source: Stable Diffusion



Category 3

Images that belong to **Category 1** with adversarial perturbations

Source: FGSM



Must-Abstain

Category 4

Anything other than the 10 mammals

Sources: ImageNet, Describable Textures, Caltech 10



Category 6

Close biological relatives of in-class mammals

Source: ImageNet



Category 5

Animals with multiple mammals' characteristics combined

Source: Stable Diffusion



Category 7

Animal representations that are not photo-like

Sources: DomainNet, ImageNet-R

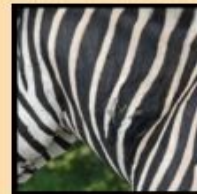


Uncertain

Category 8

Images that belong to **Category 1** with cropping or 15 corruptions applied

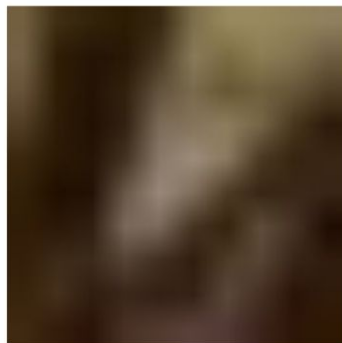
Sources: Cropping, ImageNet-C corruptions



Dataset

Uncertain Group Label Generation

- We employ 134 MTurk workers per image to classify images in Uncertain group to estimate the ground truth distribution within an error bound of 5%.



Select an option

Tiger	1
Zebra	2
Camel	3
Giraffe	4
Elephant	5
Rhino	6
Gorilla	7
Bear or Giant Panda	8
Kangaroo	9
Human	0
None of the above, uncertain, or unrecognizable	

Metrics

Alignment

We borrow Hellinger Distance to measure distance between model's probability and ground truth distributions.

$$h(P, Q) = \frac{1}{\sqrt{2}} \sum_i \|\sqrt{p_i} - \sqrt{q_i}\|_2$$

Reliability Score

Since alignment is a proxy measure for reliability, we also calculated reliability score following the below table.

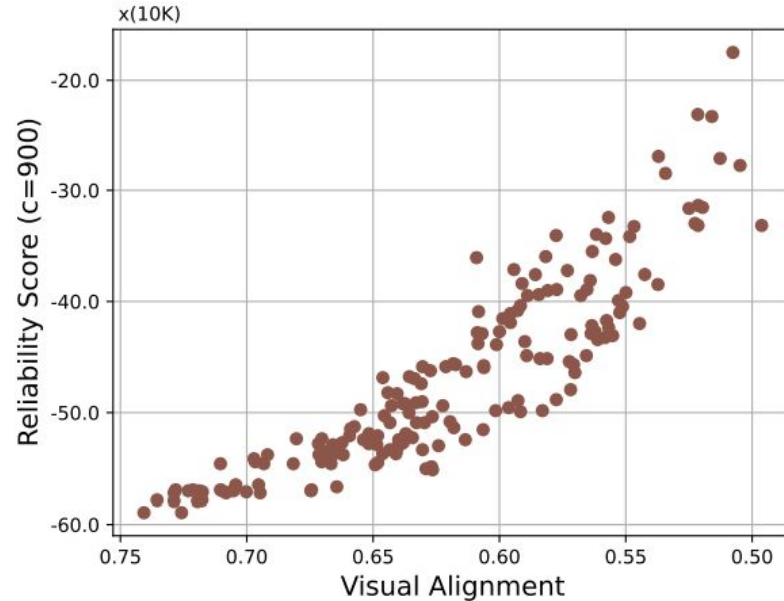
Sample Type	Model Action	$RS_c(x)$
Must-Act	Correct Prediction	+1
	Incorrect Prediction	-c
	Abstention	0
Must-Abstain	Original Label Prediction*	0
	Other Prediction	-c
	Abstention	+1

Experiment Results

	Visual Alignment (\downarrow)								Reliability score (\uparrow)			
	Must-Act			Must-Abstain				Uncertain	Average	RS_0	RS_{450}	RS_{900}
	Category 1	Category 2	Category 3	Category 4	Category 5	Category 6	Category 7	Category 8				
ViT [11]												
SP	0.261 \pm 0.051	0.556 \pm 0.029	0.367 \pm 0.038	0.793 \pm 0.057	0.808 \pm 0.057	0.787 \pm 0.056	0.792 \pm 0.059	0.671 \pm 0.032	0.629 \pm 0.021	313	-245837	-491987
ASP	0.208 \pm 0.036	0.514 \pm 0.033	0.325 \pm 0.022	1.000 \pm 0.000	1.000 \pm 0.000	1.000 \pm 0.000	1.000 \pm 0.000	0.767 \pm 0.010	0.727 \pm 0.007	253	-285047	-570347
MD [36]	0.390 \pm 0.030	0.658 \pm 0.025	0.485 \pm 0.023	0.725 \pm 0.021	0.721 \pm 0.023	0.726 \pm 0.023	0.664 \pm 0.025	0.623 \pm 0.012	0.624 \pm 0.005	270	-275580	-551430
KNN [70]	0.382 \pm 0.047	0.634 \pm 0.029	0.484 \pm 0.033	0.679 \pm 0.058	0.696 \pm 0.050	0.679 \pm 0.049	0.674 \pm 0.067	0.612 \pm 0.034	0.605 \pm 0.020	282	-264768	-529818
TAPUDD [13]	0.375 \pm 0.070	0.628 \pm 0.073	0.468 \pm 0.074	0.809 \pm 0.079	0.809 \pm 0.084	0.835 \pm 0.065	0.768 \pm 0.089	0.678 \pm 0.024	0.671 \pm 0.017	253	-285047	-570347
OpenMax [3]	0.238 \pm 0.027	0.536 \pm 0.033	0.344 \pm 0.022	0.804 \pm 0.050	0.816 \pm 0.037	0.804 \pm 0.059	0.766 \pm 0.055	0.696 \pm 0.025	0.626 \pm 0.020	335	-229165	-458665
MC-Dropout [16]	0.210 \pm 0.036	0.516 \pm 0.032	0.326 \pm 0.022	0.968 \pm 0.009	0.970 \pm 0.010	0.968 \pm 0.009	0.968 \pm 0.010	0.749 \pm 0.014	0.709 \pm 0.005	253	-285047	-570347
Deep Ensemble [33]	0.305	0.571	0.400	0.712	0.732	0.705	0.713	0.628	0.596	376	-205274	-410924

- Distance-based functions (MD, KNN, and TAPUDD) exhibits better visual alignment for Must-Act.
 - SP aligns better in Must-Abstain.
 - No current method performs well in Uncertain.
- There is no method that performs well in all categories.

Experiment Results



From the figure, we can see a strong correlation between visual alignment and reliability.

Therefore, we can assess visual alignment as a proxy measure for reliability.



NEURAL INFORMATION
PROCESSING SYSTEMS

Thank you