

# Do Invariances in Deep Neural Networks Align with Human Perception?

AAAI 2023



**Vedant Nanda**  
PhD Student  
University of Maryland & MPI-SWS



Camila  
Kolling



Ayan  
Majumdar



Krishna  
Gummadi



John  
Dickerson



Adrian  
Weller



Bradley  
Love

# Invariances are Crucial for Robust Deep Learning

We need to make sure models learn correct invariances

$$f(\text{6}) = f(\text{6}) \quad \checkmark$$

$$f(\text{6}) = f(\text{9}) \quad \times$$

Lack of human-like invariances => Models fail in unexpected ways!

Measuring alignment of invariances is a fundamental measure of robustness

# Robustness Evaluation Today

Accuracy under adversarial perturbations ([Carlini et al., 2019](#); [Madry et al., 2018](#))

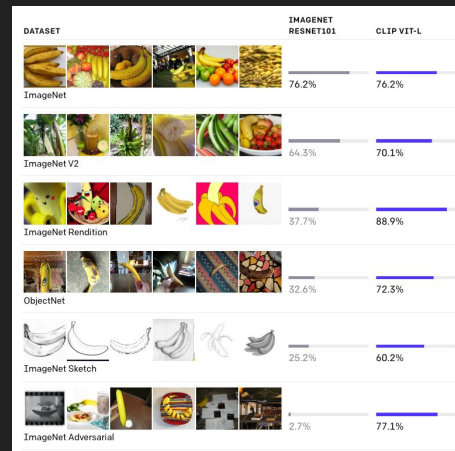
- Evaluate accuracy under worst case perturbation in a given threat model (eg:  $\ell_p$ , patch etc)

Accuracy under various distribution shifts

- ImageNetV2 ([Recht et al., 2019](#))
- ImageNet-R ([Hendrycks et al., 2021](#))
- ImageNet-C, ImageNet-P ([Hendrycks et al., 2019](#))
- ObjectNet ([Barbu et al., 2019](#))
- ImageNet-Sketch ([Wang et al., 2019](#))
- ImageNet-A ([Hendrycks et al., 2019](#))



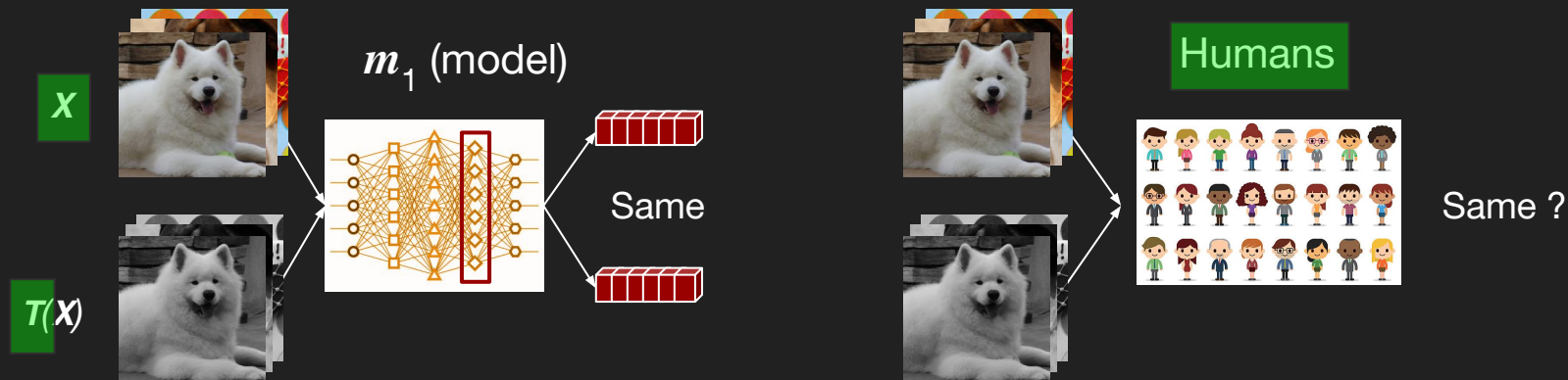
<https://robustbench.github.io> (Croce et al., 2021)



<https://openai.com/blog/clip> (Radford et al., 2021)

# The Other Direction of Robustness Evaluation

Do Invariances in DNNs align with Human Perception?



- How to choose  $X$ ?
- [Choosing  $T$ ] Infinitely many  $T$ . How to pick appropriate  $T$ ?
- [Humans] No access to representations in human brain

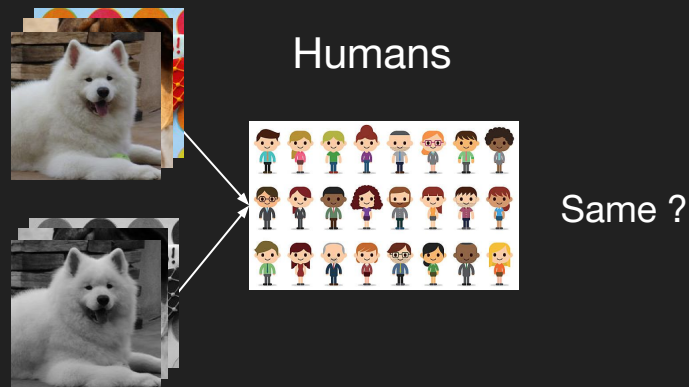
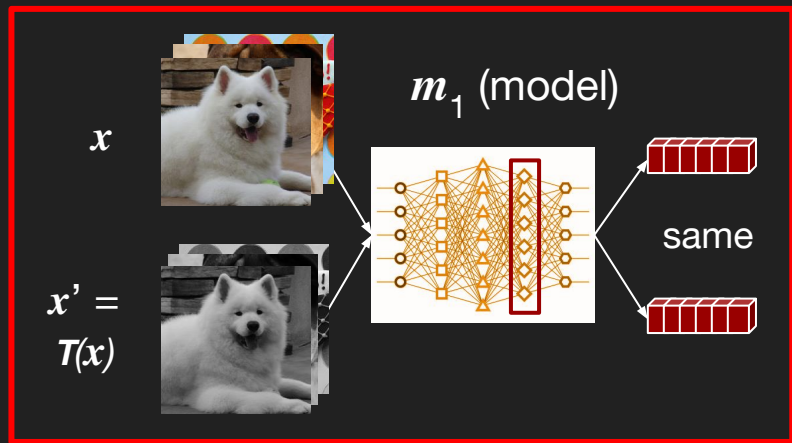
# Our Contribution

- [Choosing T] Highlight the role of loss function used in finding invariant transforms
  - Reconcile seemingly contradictory takeaways in prior work
- [Humans] Provide an improved way of measuring alignment with human perception
  - Does not require labelled data
  - Scalable
- Analyze how architectures, losses, data augmentations affect alignment

# Our Contribution

- [Choosing T] Highlight the role of loss function used in finding invariant transforms
  - Reconcile seemingly contradictory takeaways in prior work
- [Humans] Provide an improved way of measuring alignment with human perception
  - Does not require labelled data
  - Scalable
- Analyze how architectures, losses, data augmentations affect alignment

# Robustness Evaluation

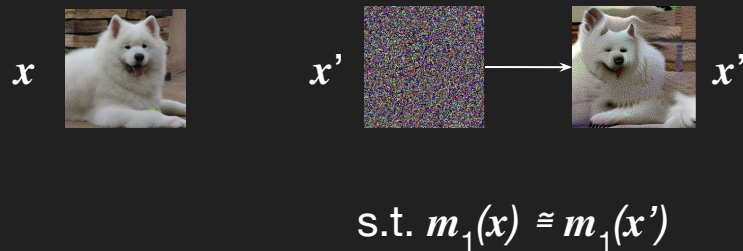


1. Find *Identically Represented Inputs (IRIs)*  $x, x'$  such that  $m_1(x) \approx m_1(x')$

Representation Inversion ([Mahendran & Vedaldi, CVPR 2015](#))

$$\operatorname{argmin}_{x'} \mathcal{L}(x')$$

$$x' = x' - \alpha \nabla_{x'} \mathcal{L}$$

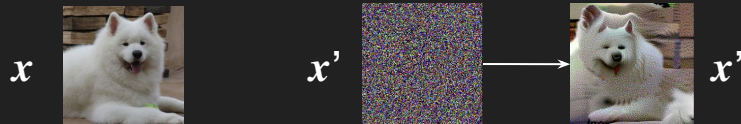


# Loss Used to Generate IRIs

1. Find *Identically Represented Inputs (IRIs)*  $\mathbf{X}, \mathbf{X}'$  such that  $m_1(\mathbf{X}) \cong m_1(\mathbf{X}')$

Representation Inversion ([Mahendran & Vedaldi, CVPR 2015](#))

$$\operatorname{argmin}_{x'} \mathcal{L}(x')$$
$$x' = x' - \alpha \nabla_{x'} \mathcal{L}$$



$$\mathcal{L}(x') = \|m_1(x) - m_1(x')\|_2 + \lambda * R(x')$$

Regularizer-free

$$R(x') = 0$$

Human-Aligned

$$R(x') = TV(x') + \|x'\|_p$$

Removes high-frequency components from  $x'$

Adversarial

$$R(x') = -1 * LPIPS(x, x')$$

Makes  $x$  and  $x'$  perceptually distant  
([Zhang et al., 2018](#))

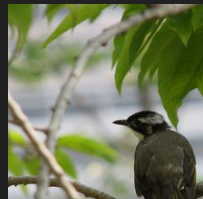


# Loss Used to Generate IRIs

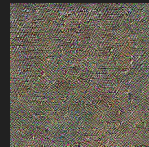
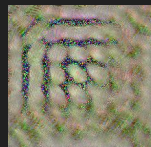
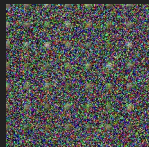
Regularizer-free  $R$   
 $R(x') = 0$

Human-Aligned  
 $R(x') = TV(x') + \|x'\|_p$

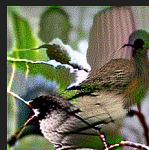
Adversarial  
 $R(x') = -1 * LPIPS(x, x')$



Standard



AT  $\ell_2 \epsilon = 1$



# Regularizer impacts takeaways about alignment!

Engstrom et al., 2019

Adversarially trained models induce a “human prior” over learned representations

Feather et al., 2019

Invariances in later layers diverge from human perception

CIFAR10				
TRAINING	MODEL	ALIGNMENT		
		REG.-FREE	HUMAN-ALIGNED	ADVERSARIAL
AT $\ell_2, \epsilon = 1$	RESNET18	63.25 $\pm$ 26.23	79.00 $\pm$ 21.94	0.33 $\pm$ 0.47
	VGG16	0.25 $\pm$ 0.43	41.41 $\pm$ 16.74	1.00 $\pm$ 1.41
	INCEPTIONV3	23.25 $\pm$ 25.56	64.75 $\pm$ 24.17	3.00 $\pm$ 4.24
	DENSENET121	82.75 $\pm$ 20.07	86.25 $\pm$ 14.50	1.33 $\pm$ 1.89
STANDARD	RESNET18	0.00 $\pm$ 0.00	21.09 $\pm$ 13.51	1.33 $\pm$ 1.89
	VGG16	0.00 $\pm$ 0.00	21.88 $\pm$ 14.82	0.00 $\pm$ 0.00
	INCEPTIONV3	0.00 $\pm$ 0.00	21.88 $\pm$ 17.54	0.33 $\pm$ 0.47
	DENSENET121	0.00 $\pm$ 0.00	26.56 $\pm$ 16.90	0.00 $\pm$ 0.00
IMAGENET				
TRAINING	MODEL	ALIGNMENT		
		REG.-FREE	HUMAN-ALIGNED	ADVERSARIAL
AT $\ell_2, \epsilon = 3$	RESNET18	42.00 $\pm$ 38.33	46.75 $\pm$ 39.37	0.33 $\pm$ 0.47
	RESNET50	51.00 $\pm$ 34.89	45.75 $\pm$ 37.39	14.00 $\pm$ 3.74
	VGG16	55.50 $\pm$ 34.14	55.50 $\pm$ 38.29	11.00 $\pm$ 3.74
STANDARD	RESNET18	0.00 $\pm$ 0.00	17.00 $\pm$ 28.30	0.00 $\pm$ 0.00
	RESNET50	0.00 $\pm$ 0.00	16.25 $\pm$ 26.42	0.00 $\pm$ 0.00
	VGG16	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00

Olah et al., 2017

Parts of DNNs encode human-like concepts

**Prior works do not directly engage with the choice of regularizer and hence make incomplete conclusions**

Under the pessimistic lens of adversarial regularizer all models are poorly aligned

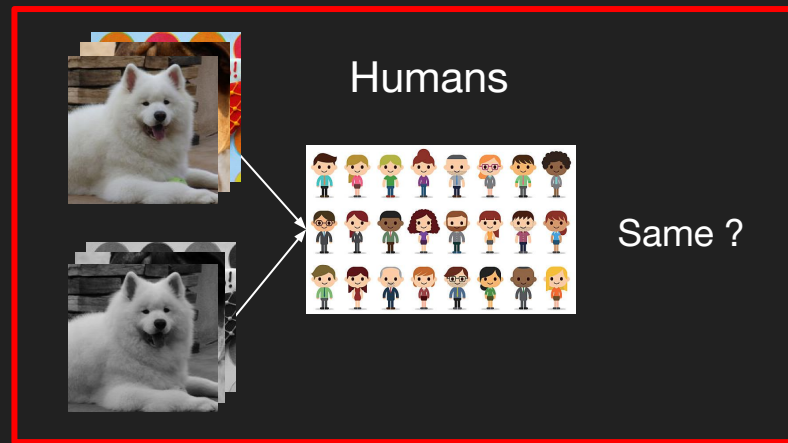
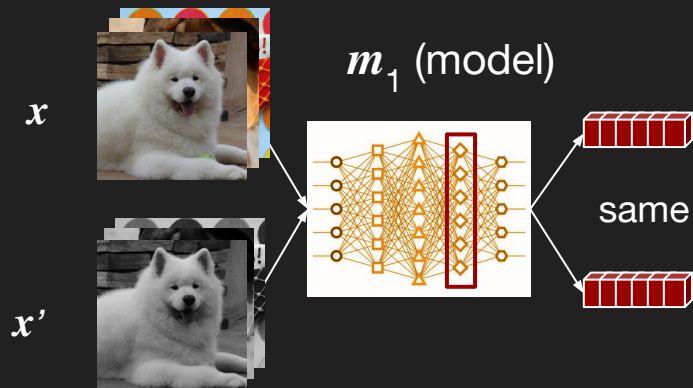
# Our Contribution

- [Choosing T] Highlight the role of loss function used in finding invariant transforms
  - Reconcile seemingly contradictory takeaways in prior work – choice of regularizer impacts takeaways
- [Humans] Provide an improved way of measuring alignment with human perception
  - Does not require labelled data
  - Scalable
- Analyze how architectures, losses, data augmentations affect alignment

# Our Contribution

- [Choosing  $T$ ] Highlight the role of loss function used in finding invariant transforms
  - Reconcile seemingly contradictory takeaways in prior work – choice of regularizer impacts takeaways
- [Humans] Provide an improved way of measuring alignment with human perception
  - Does not require labelled data
  - Scalable
- Analyze how architectures, losses, data augmentations affect alignment

# Robustness Evaluation




We need to *reliably* and *scalably* check if Humans perceive  $x$  and  $x'$  similarly

# Check if humans perceive these inputs similarly

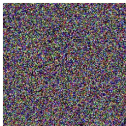
## 2AFC

Query Image 1/13:



Final  $x'$









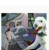
Which of these images is closer to the query image:

Target $x$			Initial $x'$
<input type="radio"/>		<input type="radio"/>	

## Clustering

Query 1/2:

For each query image (shown in each row), choose the most perceptually similar image from the columns:

			
	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

Scalability: since these tests are based on comparisons, we can use perceptual distance measures like LPIPS to simulate humans ([Zhang et al., 2018](#))

# Evaluation: Reliability



			CIFAR10		
TRAINING	MODEL	HUMAN 2AFC	HUMAN CLUSTERING		
ROBUST	RESNET18	96.00 $\pm$ 2.55	97.48 $\pm$ 1.80		Both clustering and 2AFC achieve similar ranking among models
	VGG16	38.83 $\pm$ 7.59	55.387 $\pm$ 5.629		
	INCEPTIONV3	82.00 $\pm$ 8.44	84.47 $\pm$ 6.32		
	DENSENET121	98.67 $\pm$ 0.24	97.64 $\pm$ 2.08		
STANDARD	RESNET18	0.17 $\pm$ 0.24	38.55 $\pm$ 1.19		Low variance among annotators
	VGG16	0.17 $\pm$ 0.24	33.84 $\pm$ 2.70		
	INCEPTIONV3	0.17 $\pm$ 0.24	38.38 $\pm$ 4.06		
	DENSENET121	9.83 $\pm$ 9.97	42.42 $\pm$ 5.02		
			IMAGENET		
TRAINING	MODEL	HUMAN 2AFC	HUMAN CLUSTERING		
ROBUST	RESNET18	93.17 $\pm$ 5.95	96.00 $\pm$ 3.59		=> Humans can determine alignment reliably
	RESNET50	99.50 $\pm$ 0.00	99.49 $\pm$ 0.71		
	VGG16	95.50 $\pm$ 2.12	91.75 $\pm$ 5.22		
STANDARD	RESNET18	0.00 $\pm$ 0.00	33.33 $\pm$ 0.00		
	RESNET50	5.33 $\pm$ 7.54	38.38 $\pm$ 2.53		
	VGG16	0.00 $\pm$ 0.00	33.96 $\pm$ 2.00		

# Evaluation: Scalability



CIFAR10					
TRAINING	MODEL	HUMAN 2AFC	LPIPS 2AFC	HUMAN CLUSTERING	LPIPS CLUSTERING
ROBUST	RESNET18	96.00 $\pm$ 2.55	98.0	97.48 $\pm$ 1.80	91.41
	VGG16	38.83 $\pm$ 7.59	11.5	55.387 $\pm$ 5.629	46.97
	INCEPTIONV3	82.00 $\pm$ 8.44	87.5	84.47 $\pm$ 6.32	76.77
	DENSENET121	98.67 $\pm$ 0.24	100.0	97.64 $\pm$ 2.08	97.98
STANDARD	RESNET18	0.17 $\pm$ 0.24	0.0	38.55 $\pm$ 1.19	31.31
	VGG16	0.17 $\pm$ 0.24	0.0	33.84 $\pm$ 2.70	31.31
	INCEPTIONV3	0.17 $\pm$ 0.24	1.0	38.38 $\pm$ 4.06	38.38
	DENSENET121	9.83 $\pm$ 9.97	0.5	42.42 $\pm$ 5.02	36.87
IMAGENET					
TRAINING	MODEL	HUMAN 2AFC	LPIPS 2AFC	HUMAN CLUSTERING	LPIPS CLUSTERING
ROBUST	RESNET18	93.17 $\pm$ 5.95	82.00	96.00 $\pm$ 3.59	76.77
	RESNET50	99.50 $\pm$ 0.00	85.00	99.49 $\pm$ 0.71	82.83
	VGG16	95.50 $\pm$ 2.12	82.00	91.75 $\pm$ 5.22	78.79
STANDARD	RESNET18	0.00 $\pm$ 0.00	0.50	33.33 $\pm$ 0.00	34.85
	RESNET50	5.33 $\pm$ 7.54	0.50	38.38 $\pm$ 2.53	35.86
	VGG16	0.00 $\pm$ 0.00	0.00	33.96 $\pm$ 2.00	34.34

LPIPS orders  
models  
same as  
humans

=> Can  
analyze  
models at  
scale



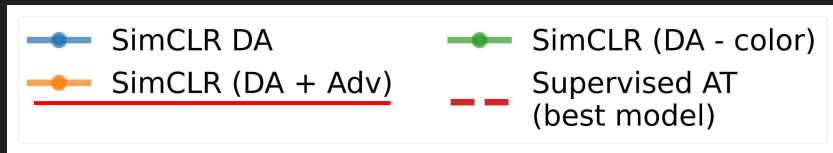
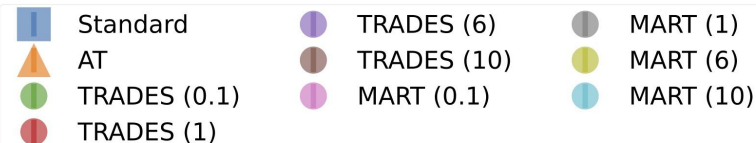
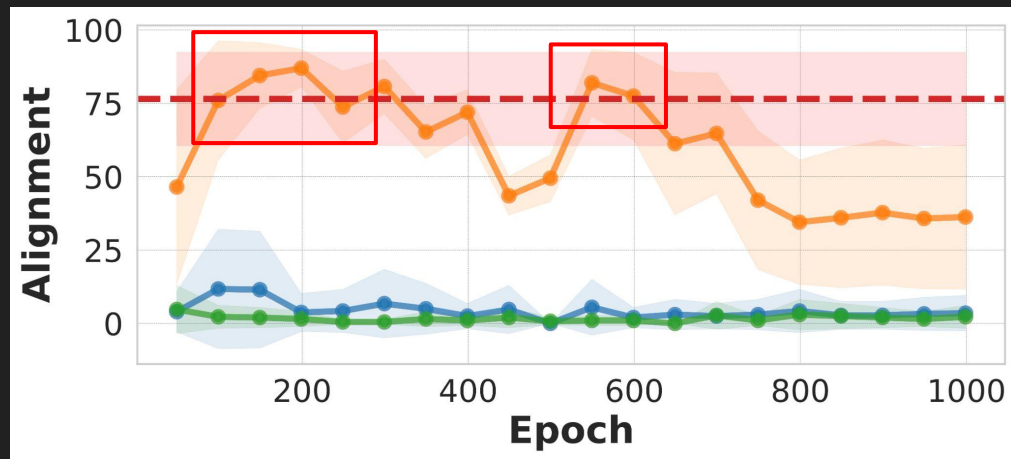
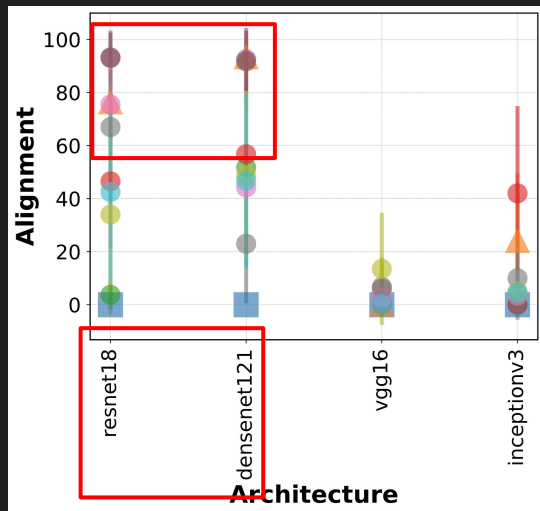
# Our Work

- [Choosing T] Highlight the role of loss function used in finding invariant transforms
  - Reconcile seemingly contradictory takeaways in prior work – choice of regularizer impacts takeaways
- [Humans] Provide an improved way of measuring alignment with human perception
  - Does not require labelled data
  - Scalable
- Analyze how architectures, losses, data augmentations affect alignment

# Our Work

- [Choosing T] Highlight the role of loss function used in finding invariant transforms
  - Reconcile seemingly contradictory takeaways in prior work – choice of regularizer impacts takeaways
- [Humans] Provide an improved way of measuring alignment with human perception
  - Does not require labelled data
  - Scalable
- Analyze how architectures, losses, data augmentations affect alignment

# Results: Architectures, Losses, Data Augmentations



1. Adversarial data augmentation using  $\ell_2$  threat model

2. Architectures with residual connections

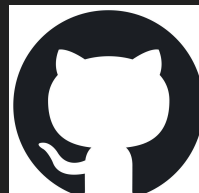
3. Self-supervised contrastive loss

# Summary

- We highlight challenges and common pitfalls in measuring alignment with human perception
- We propose an improved method to measure alignment at scale
- Using our method we show how residual connections, adversarial data augmentation and contrastive loss help in increasing alignment



[tinyurl.com/invariances-human](https://tinyurl.com/invariances-human)



[github.com/nvedant07/Human-NN-Alignment](https://github.com/nvedant07/Human-NN-Alignment)

**Thank You!**

**[vnanda@mpi-sws.org](mailto:vnanda@mpi-sws.org)**

**Poster # 111**

# Prior Works

([Colah et al., 2017](#)) Certain parts of DNNs encode human-like concepts

([Engstrom et al., 2019](#)) Adversarially trained models induce a “human prior” over learned representations

([Feather et al., 2019](#)) Invariances in later layers diverge from human perception

Contradictory takeaways. What's going on?

# Evaluation

CIFAR10							
TRAINING	MODEL	HUMAN 2AFC	LPIPS 2AFC	HUMAN CLUSTERING	LPIPS CLUSTERING	CLEAN ACC.	ROBUST ACC.
ROBUST	RESNET18	96.00 $\pm$ 2.55	98.0	97.48 $\pm$ 1.80	91.41	80.77	50.92
	VGG16	38.83 $\pm$ 7.59	11.5	55.387 $\pm$ 5.629	46.97	79.84	48.36
	INCEPTIONV3	82.00 $\pm$ 8.44	87.5	84.47 $\pm$ 6.32	76.77	81.57	51.02
	DENSENET121	98.67 $\pm$ 0.24	100.0	97.64 $\pm$ 2.08	97.98	83.22	52.86
STANDARD	RESNET18	0.17 $\pm$ 0.24	0.0	38.55 $\pm$ 1.19	31.31	94.94	0.00
	VGG16	0.17 $\pm$ 0.24	0.0	33.84 $\pm$ 2.70	31.31	93.63	0.00
	INCEPTIONV3	0.17 $\pm$ 0.24	1.0	38.38 $\pm$ 4.06	38.38	94.59	0.00
	DENSENET121	9.83 $\pm$ 9.97	0.5	42.42 $\pm$ 5.02	36.87	95.30	0.00
IMAGENET							
TRAINING	MODEL	HUMAN 2AFC	LPIPS 2AFC	HUMAN CLUSTERING	LPIPS CLUSTERING	CLEAN ACC.	ROBUST ACC.
ROBUST	RESNET18	93.17 $\pm$ 5.95	82.00	96.00 $\pm$ 3.59	76.77	53.12	31.02
	RESNET50	99.50 $\pm$ 0.00	85.00	99.49 $\pm$ 0.71	82.83	62.83	38.84
	VGG16	95.50 $\pm$ 2.12	82.00	91.75 $\pm$ 5.22	78.79	56.79	34.46
STANDARD	RESNET18	0.00 $\pm$ 0.00	0.50	33.33 $\pm$ 0.00	34.85	69.76	0.01
	RESNET50	5.33 $\pm$ 7.54	0.50	38.38 $\pm$ 2.53	35.86	76.13	0.00
	VGG16	0.00 $\pm$ 0.00	0.00	33.96 $\pm$ 2.00	34.34	73.36	0.16

# Our Work

- Reconcile difference in takeaways of prior work [Choosing T]
  - Highlight challenges and common pitfalls in measuring alignment – choice of regularizer impacts takeaways
- Provide an improved way of measuring alignment with human perception [Humans]
  - Does not require labelled data
  - Scalable
- Analyze how architectures, losses, data augmentations affect alignment
  - Architectures with residual connections,
  - Adversarial data augmentation using  $\ell_2$  threat model
  - Self-supervised contrastive loss