

Do Invariances in Deep Neural Networks Align with Human Perception?

Vedant Nanda^{*1,2}, Ayan Majumdar², Camila Kolling², John P. Dickerson¹, Krishna P. Gummadi²,
Bradley C. Love^{3,4}, Adrian Weller^{3,5},

¹University of Maryland ²MPI-SWS

³The Alan Turing Institute ⁴University College London ⁵University of Cambridge

Abstract

An evaluation criterion for safe and trustworthy deep learning is how well the invariances captured by representations of deep neural networks (DNNs) are shared with humans. We identify challenges in measuring these invariances. Prior works used gradient-based methods to generate *identically represented inputs* (IRIs), *i.e.*, inputs which have identical representations (on a given layer) of a neural network, and thus capture invariances of a given network. One necessary criterion for a network’s invariances to align with human perception is for its IRIs look “similar” to humans. Prior works, however, have mixed takeaways; some argue that later layers of DNNs do not learn human-like invariances ([12]) yet others seem to indicate otherwise ([36]). We argue that the loss function used to generate IRIs can heavily affect takeaways about invariances of the network and is the primary reason for these conflicting findings. We propose an *adversarial* regularizer on the IRI-generation loss that finds IRIs that make any model appear to have very little shared invariance with humans. Based on this evidence, we argue that there is scope for improving models to have human-like invariances, and further, to have meaningful comparisons between models one should use IRIs generated using the *regularizer-free* loss. We then conduct an in-depth investigation of how different components (*e.g.* architectures, training losses, data augmentations) of the deep learning pipeline contribute to learning models that have good alignment with humans. We find that architectures with residual connections trained using a (self-supervised) contrastive loss with l_p ball adversarial data augmentation tend to learn invariances that are most aligned with humans. Code: github.com/nvedant07/Human-NN-Alignment

1 Introduction

The ability to train deep neural networks (DNNs) which learn useful features and representations is key for their widespread use [3, 33]. In domains where DNNs are used for tasks that previously required human intelligence (*e.g.* image classification) and where safety and trustworthiness are important considerations, it is helpful to assess the alignment of the learned representations with human perception. Such assessments can help in understanding and diagnosing issues such as lack of robustness to distribution shifts [48, 60], adver-

sarial attacks [44, 17] or using undesirable features for a downstream task [2, 53, 50, 4, 24].

One test of human-machine alignment is whether different images that map to identical internal network representation are also judged as identical by humans. To study alignment with human perception, prior works have used the approach of *representation inversion* [36]. The key idea is the following: given an input to a neural network, the approach first finds *identically represented inputs* (IRIs), *i.e.* inputs which have similar representations on some given layer(s) of the neural network. In the second step, the inputs that are perceived similarly by the neural network are checked by humans for visual similarity. Thus, the approach relies on estimating whether a transformation of the inputs which is representation invariant to a neural network is also an invariant transformation to the human eye, *i.e.* it checks whether models and humans have shared or aligned invariances.

Prior works use gradient-based methods to generate IRIs for a given target input starting with a random seed input. These works revealed exciting insights: (a) Feather et al. studied representational invariance for different layers of DNNs trained over ImageNet data (using the standard cross-entropy loss). They showed that while later layer representations of DNNs do not share any invariances with human perception, the earlier layers are somewhat better aligned with human perception [12]. (b) Engstrom et al. found that, unlike standard DNNs, adversarially robust DNNs, *i.e.*, DNNs trained using adversarial training [35], learn representations that are well aligned with human perception, even in later layers [10]. This was also confirmed by other works [26, 55]. However, some of these findings are contradicted when differently regularized methods are used for generating IRIs, which show that even later layers of DNNs learn human aligned representations [36, 42, 41].

We seek to make sense of these confusing earlier results, and thereby to better understand alignment. We show that when we evaluate alignment of DNNs’ invariances and human perception using IRIs generated using different loss functions, we can arrive at very different conclusions. For example, Fig 1 shows how visual similarity of IRIs can vary massively across different categories of losses.

We group existing IRI generation processes into two broad categories: *regularizer-free* (as in [12]), where the goal is to find an IRI without any additional constraints; and *human-*

*Correspondence to: vnanda@mpi-sws.org
Copyright © 2023, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

leaning (as in [42, 41, 36]), where the goal is to find an IRI that is also visually human-comprehensible. Additionally, we propose and explore a new (third) broad category, *adversarial*, where the goal is to find an IRI that is visually (from a human perception perspective) far apart from the target input.

We find that compared to the regularizer-free IRI generation approach, the human-leaning IRI generation approach applies strong constraints on the kind of IRIs generated and thus limits the ability to freely explore the large space of possible IRIs. On the other hand, our proposed adversarial approach shows that in the worst case, all models have close to zero alignment, suggesting that there is scope for improvement in designing models that have human-like invariances (as shown in Fig 1 and Table 2). Based on this evidence, we argue that in order to have meaningful comparisons between models, one should measure alignment using the regularizer-free loss for IRI generation.

Many prior works do not formally define a measure that can quantify alignment with human perception beyond relying on visual inspection of the images by the authors [42, 41, 36, 10]. We show how alignment can be quantified reliably by designing simple visual perception tests that can be crowdsourced, *i.e.* used in human surveys. We also show how one can leverage widely used measures of perceptual distance [66] to automate our human surveys, which allows us to obtain insights at a scale not possible in previous works.

Next, inspired by the prior works that suggest that changes in the model training pipeline (as in training adversarially robust DNNs [10, 26]) can lead to human-like invariant representations, we conduct an in-depth investigation to understand which parts of the deep learning pipeline are critical in helping DNNs better learn human-like invariances. We find that certain choices in the deep learning pipeline can significantly help learn representation that have human-like invariances. For example, we show that residual architectures (*e.g.*, ResNets [21]), when trained with a self-supervised contrastive loss (*e.g.*, SimCLR [5]), using ℓ_2 ball adversarial data augmentations (*e.g.*, as in RoCL [27]); the learned representations – while typically having lower accuracies than their fully supervised counterparts – have higher alignment of invariances with human perception. We highlight the following contributions:

- We show how different losses used for generating IRIs lead to different conclusions about a model’s shared invariances with human perception, thus leading to seemingly contradictory findings in prior works.
- We propose an adversarial IRI generation loss, using which we show empirically that we can almost always discover invariances of DNNs that do not align with human perception, thus suggesting that there is scope to design better mechanisms to learn representations that are more aligned with human perception.
- We conduct an in-depth study of how loss functions, architectures, data augmentations and training paradigms lead to learning human-like shared invariances.

2 Measuring Shared Invariance with Human Perception

Measuring the extent to which invariances learned by DNNs are shared by humans is a two step process. We first generate IRIs, *i.e.*, inputs that are mapped to identical representations by the DNN. IRIs give us an estimate about the invariances of the DNN. Then, we assess if these inputs are also considered identical by humans. More concretely, if invariances of a given DNN (g_{model}) are shared by humans (g_{human}) on a set of n d -dimensional samples $X \in \mathbb{R}^{n \times d}$, then:

$$g_{\text{human}}(X^i) \approx g_{\text{human}}(X^j) \forall (X^i, X^j) \in \mathcal{S} \times \mathcal{S};$$

$$\mathcal{S} = \{X\} \cup \{X^i \mid g_{\text{model}}(X^i) \approx g_{\text{model}}(X)\}.$$

\mathcal{S} denotes the IRIs for g_{model} . There are three major challenges here:

- Access to representations in the brain, *i.e.*, g_{human} is not available.
- Due to the highly non-linear nature of DNNs, \mathcal{S} can be very hard to obtain.
- The fine-grained input space implies very many inputs n , making the choice of X hard.

We address each of these below. We also show how prior works that do not directly engage with these points can miss important issues in their conclusions about shared invariances of DNNs and humans.

2.1 Approximating g_{human}

Assuming we have a set of images with identical representations (\mathcal{S} ; how we obtain this is discussed in Section 2.2), we must check if humans also perceive these images to be identical. The extent to which humans think this set of images is identical defines how aligned the invariances learned by the DNN are with human perception. In prior works this has been done by either eyeballing IRIs [10] or by asking annotators to assign class labels to IRIs [12]; both approaches do not scale well. Additionally, assigning class labels to IRIs limits X to being samples from a data distribution containing human-recognizable images (*i.e.*, X cannot be sampled from any arbitrary distribution) with only a few annotations (*e.g.*, asking annotators to assign one class label out of 1000 ImageNet classes is not feasible). To address the issues of scalability and class labels, we propose the following as a measure of alignment between DNN and human invariances:

$$\text{Alignment} = \frac{|\mathcal{A}|}{\sum_{x_t \in X} |\mathcal{S}_{x_t}|}, \text{ where} \quad (1)$$

$$\mathcal{A} = \{x_{r_i} \mid \|g_{\text{human}}(x_t) - g_{\text{human}}(x_{r_i})\| < \|g_{\text{human}}(x_0) - g_{\text{human}}(x_{r_i})\| \forall x_t \in X, x_{r_i} \in \mathcal{S}_{x_t}\},$$

$$\mathcal{S}_{x_t} = \{x_{r_i} \mid g_{\text{model}}(x_{r_i}) \approx g_{\text{model}}(x_t) \forall x_t \in X\},$$

where x_0 is the starting point for Eq 2 sampled from $\mathcal{N}(0, 1)$. In Section 2.4 we see how alignment is robust to the choice of x_0 . By directly looking for perceptual similarity

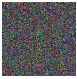
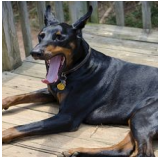
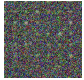




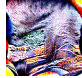
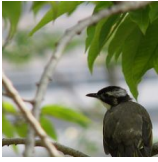






Seed (x_0)		Regularizer-free	Human-leaning Regularizer	Adversarial Regularizer
Target (x_t)	Model	Result (x_r)		
	Standard			
	AT $\ell_2 \epsilon = 1$			
	Standard			
	AT $\ell_2 \epsilon = 1$			

Figure 1: **[Representation Inversion for different kinds of \mathcal{R} ; For ImageNet trained ResNet50]** For the standard ResNet50 (trained using cross-entropy loss), with regularizer-free and adversarial inversion, x_r looks perceptually much closer to x_0 than x_t , even though from the model’s point of view, x_r and x_t are the same. However, with the human-leaning regularizer, we see that x_r contains some information like color patterns of x_t . For adversarially robust ResNet50 [35, 54] even though regularizer-free and human-leaning inversions look perceptually similar to x_t , for the adversarial regularizer even these models produce x_r that looks nothing like x_t . Images are generated by starting from x_0 and solving Eq 2 with different kinds of regularizers.

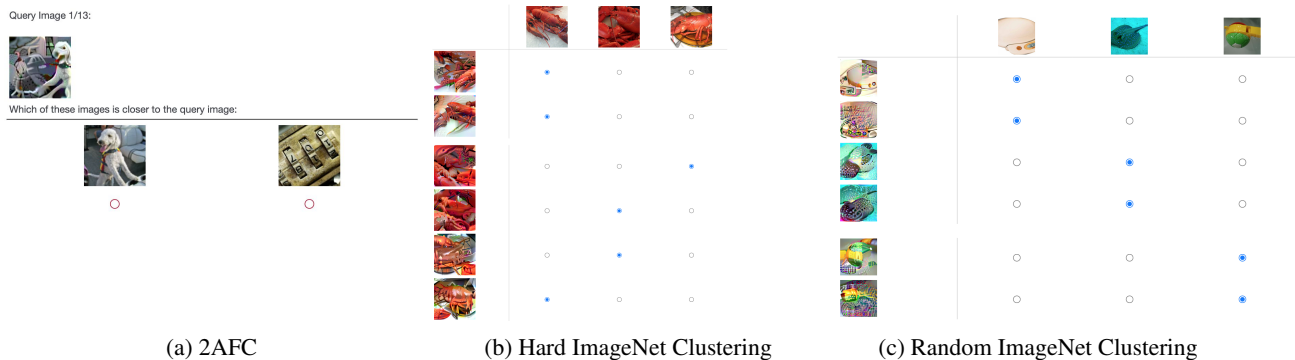


Figure 2: **[Survey Prompts for AMT workers]** In the 2AFC (left) setting we ask the annotator to choose which of the two images (x_t or x_0) is perceptually closer to the query image (x_r). In the clustering setting (center and right) we show 3 images from the dataset (target images, x_t) in the columns and for each of these, we generate $x_{r_1} \in \mathcal{S}_{x_t}$ and $x_{r_2} \in \mathcal{S}_{x_t}$. Each of these is shown across the rows. The task here is to match each image on the row with the corresponding target image on the column.

of IRIs (captured by \mathcal{A}), we get past the issue of assigning class labels to IRIs. The comparison used to generate \mathcal{A} is referred to as the 2 alternative forced choice test (2AFC) which is commonly used to assess sensitivity of humans to stimuli [13]. In order to compute \mathcal{A} , we estimate perceptual distance $d(x_i, x_j) = \|g_{\text{human}}(x_i) - g_{\text{human}}(x_j)\|$ between two inputs. Ideally, we would like to measure $d(x_i, x_j)$ by directly asking for human annotations, however, this approach is expensive and does not scale when we wish to evaluate many models. To address scalability, we use LPIPS [66] which is a commonly used measure for perceptual distance and thus can be used to approximate $d(x_i, x_j)$ ¹. While LPIPS is by no means a perfect approximation, it allows us to gain insights

¹For all evaluations we report the average over 4 different backbones used to calculate LPIPS including the finetuned weights released by the authors. More details in Appendix A.2

at a scale not possible in prior works.

To ensure the efficacy of LPIPS as a proxy for human judgements, we deploy two types of surveys on Amazon Mechanical Turk (AMT) to also elicit human similarity judgements. Prompts for these surveys are shown in Fig. 2. We received approval from the Ethical Review Board of our institute for this survey. Each survey consists of 100 images plus some attention checks to ensure the validity of our responses. The survey was estimated to take 30 minutes (even though on average our annotators took less than 20 minutes), and we paid each worker 7.5 USD per survey.

Clustering In this setting, we ask humans to match the IRIs (x_{r_i}) on the row to the most perceptually similar image (x_t) on the column (each row can only be matched to one column). A prompt for this type of a task is shown in Fig. 2b & 2c. With these responses, we calculate a quantitative mea-

CIFAR10						
	MODEL	HUMAN 2AFC	LPIPS 2AFC	HUMAN CLUSTERING	LPIPS CLUSTERING	
AT ℓ_2 $\epsilon = 1$	RESNET18	96.00 \pm 2.55	87.25 \pm 9.52	97.48 \pm 1.80		88.13 \pm 6.57
	VGG16	38.83 \pm 7.59	4.00 \pm 3.86	55.39 \pm 5.63		46.09 \pm 3.84
	INCEPTIONV3	82.00 \pm 8.44	54.12 \pm 19.23	84.47 \pm 6.32		74.87 \pm 6.74
	DENSENET121	98.67 \pm 0.24	91.75 \pm 8.2	97.64 \pm 2.08		91.92 \pm 6.13
STANDARD	RESNET18	0.17 \pm 0.24	0.0 \pm 0.0	38.55 \pm 1.19		35.35 \pm 3.27
	VGG16	0.17 \pm 0.24	0.0 \pm 0.0	33.84 \pm 2.70		32.58 \pm 1.04
	INCEPTIONV3	0.17 \pm 0.24	0.38 \pm 0.41	38.38 \pm 4.06		36.62 \pm 3.08
	DENSENET121	9.83 \pm 9.97	0.12 \pm 0.22	42.42 \pm 5.02		37.12 \pm 3.54
IMAGENET						
	MODEL	HUMAN 2AFC	LPIPS 2AFC	HUMAN CLUSTERING	HUMAN CLUSTERING HARD	LPIPS CLUSTERING
AT ℓ_2 $\epsilon = 3$	RESNET18	93.17 \pm 5.95	53.37 \pm 20.19	96.00 \pm 3.59	87.75 \pm 7.60	65.28 \pm 10.58
	RESNET50	99.50 \pm 0.00	53.63 \pm 20.64	99.49 \pm 0.71	97.06 \pm 3.47	71.21 \pm 9.93
	VGG16	95.50 \pm 2.12	59.38 \pm 21.48	91.75 \pm 5.22	90.69 \pm 3.13	70.33 \pm 9.78
STANDARD	RESNET18	0.00 \pm 0.00	1.12 \pm 1.67	33.33 \pm 0.00	-	34.60 \pm 0.56
	RESNET50	5.33 \pm 7.54	0.38 \pm 0.41	38.38 \pm 2.53	-	35.35 \pm 0.62
	VGG16	0.00 \pm 0.00	0.00 \pm 0.00	33.96 \pm 2.00	-	34.47 \pm 1.49

Table 1: [CIFAR10 and ImageNet Surveys To Confirm Efficacy of LPIPS] We use LPIPS to simulate a human in both 2AFC and Clustering setups described in Section 2.1 and compare it with AMT worker’s responses. A value close to 33% for clustering means random assignment and indicates no alignment. We see that LPIPS and humans rank models similarly in both 2AFC and clustering setups, thus showing that LPIPS is a reliable proxy for judging perceptual similarity of IRIs. These experiments were conducted on IRIs generated using regularizer-free loss in Eq 2. The variance reported for LPIPS is for different backbone networks that are available for LPIPS.

sure of alignment by measuring the fraction of x_{r_i} that were correctly matched to their respective x_t . For ImageNet, we observed that a random draw of three images (e.g., Fig. 2c) can often be easy to match to based on how different the drawn images (x_t) are. Thus, we additionally construct a “hard” version of this task by ensuring that the three images are very “similar” (as shown in Fig. 2b). We leverage human annotations of ImageNet-HSJ [49] to draw these similar images. More details can be found in Appendix A.

2AFC This is the exact test used to generate \mathcal{A} . In this setting we show the annotator a reconstructed image (x_r) and ask them to match it to one of the two images shown in the options. The images shown in the options are the seed (x_0 , i.e., starting value of x in Eq. 2) and the original image (x_t). Since x_r and x_t are IRIs for the model (by construction), alignment would imply humans also perceive x_r and x_t similarly. See Fig. 2a for an example of this type of survey.

2.2 Generating IRIs

Even if we assume a finite sampled set $X \sim \mathcal{D}$ (discussed in Section 2.3), there can be many samples in \mathcal{S} due to the highly non-linear nature of DNNs. However, we draw on the insight that there is often some structure to the set of IRIs, that is heavily dependent on the IRI generation process. Prior work on understanding shared invariance between DNNs and humans [e.g., 10, 12] has used representation inversion [36] to generate IRIs. However, IRIs generated this way depend heavily on the loss function used in representation inversion, as demonstrated by [42]. Fig. 1 shows how different loss functions can lead to very different looking IRIs. We group these losses previously used in the literature to generate IRIs into two broad types: **regularizer-free** (used by [10, 12]), and

human-leaning (used by many works on interpretability of DNNs including [36, 42, 41, 38, 40]). We also explore a third kind of **adversarial** regularizer, that aims to generate *controversial stimuli* [16] between a DNN and a human.

Representation inversion is the task of starting with a random seed image x_0 to reconstruct a given image $x_t \in X$ from its representation $g(x_t)$ where $g(\cdot)$ is the trained DNN. The reconstructed image (x_r) is same as x_t from the DNN’s point of view, i.e., $g(x_t) \approx g(x_r)$. This is achieved by performing gradient descent on x_0 (in our experiments we use SGD with a learning rate of 0.1) to minimize a loss of the following general form:

$$\mathcal{L}_x = \frac{\|g(x_t) - g(x)\|_2}{\|g(x_t)\|_2} + \lambda * \mathcal{R}(x) \quad (2)$$

where λ is an appropriate scaling constant for regularizer \mathcal{R} . All of these reconstructions induce representations in the DNN that are very similar to the given image (x_t), as measured using ℓ_2 norm. Depending on the choice of seed x_0 and the choice of \mathcal{R} , we get different reconstructions of x_t thus giving us a set of inputs $\{x_t, x_{r_1}, \dots, x_{r_k}\}$ that are all mapped to similar representations by $g(\cdot)$. Doing this for all $x_t \in X$, we get the IRIs, $\mathcal{S} = \{X, X^{r_1}, \dots, X^{r_k}\}$.

In practice we find that the seed x_0 does not have any significant impact on the measurement of shared invariance. However, the choice of \mathcal{R} *does* significantly impact the invariance measurement (as also noted by [42]). We identify the following distinct categories of IRIs based on the choice of \mathcal{R} .

Regularizer-free. These methods do not use a regularizer, i.e., $\mathcal{R}(x) = 0$.

human-leaning regularizer. This kind of a regularizer

purposefully puts constraints on x such that the reconstruction has some “meaningful” features. [36] use $\mathcal{R}(x) = TV(x) + \|x\|_p$ where TV is the total variation in the image. Intuitively this penalizes high frequency features and smoothens the image to make it look more like natural images. [40] achieve a similar kind of high frequency penalization by blurring x before each optimization step. We combine both these frequency-based regularizers with pre-conditioning in the Fourier domain [42] and robustness to small transformations [40]. More details can be found in Appendix A.4. Intuitively a regularizer from this category generates IRIs that have been “biased” to look meaningful to humans.

Adversarial regularizer. We propose a new regularizer to generate IRIs while intentionally making them look *perceptually dissimilar* from the target, *i.e.*, $\mathcal{R} = -\|g_{\text{human}}(x_t) - g_{\text{human}}(x)\|$ (negative sign since we want to maximize perceptual distance between x and x_t). We leverage LPIPS (Learned Perceptual Image Patch Similarity) [66], a widely used *perceptual distance* measure, to approximate $\|g_{\text{human}}(x_t) - g_{\text{human}}(x)\|$. LPIPS uses initial layers of an ImageNet trained model (finetuned on a dataset of human similarity judgements) to approximate perceptual distance between images which makes it differentiable and thus can be easily plugged into Eq. 2. Thus, the regularizer used is $\mathcal{R}(x) = -\text{LPIPS}(x, x_t)$. IRIs generated using this regularizer can be thought of as *controversial stimuli* [16] – they’re similar from the DNN’s perspective, but distinct from a human’s perspective.

2.3 Choice of inputs X

In order to overcome the challenge of choosing X , we assume X to be sampled from a given data distribution \mathcal{D} . In our experiments, we try out many different distributions, including the training data distribution and random noise distributions, and find that takeaways about a alignment of model’s invariances with humans do not depend heavily on the choice of \mathcal{D} . Some examples of X sampled from the data distribution and noise distributions (two random Gaussian distributions, $\mathcal{N}(0, 1)$ and $\mathcal{N}(0.5, 2)$), along with the corresponding IRIs are shown in Fig. 4, Appendix A.3. Interestingly, the human-leaning regularizer, which explicitly tries to remove high-frequency features from x_r , fails to reconstruct an x_t that itself consists of high-frequency features.

2.4 Evaluation and Takeaways

For each model, we randomly picked 100 images from the data distribution along with a seed image with random pixel values. For each of the 100 images, we do representation inversion using one regularizer each from *regularizer-free*, *human-leaning*, and *adversarial*.

Reliability of using LPIPS Table 1 shows the results for the surveys conducted with AMT workers². Each survey was completed by 3 workers. For a well aligned model, the scores under 2AFC and Clustering should be close to 1, while for a non-aligned model scores under 2AFC should be close to 0, and scores under Clustering should be close to a

²This was conducted only using IRIs from regularizer-free inversion.

random guess (*i.e.*, about 33%). We see that LPIPS (with different backbone nets, e.g., AlexNet, VGG) orders models similar to human annotators for both the survey setups, thus showing that it’s a reliable proxy.

Reliability of Human Annotators In Table 1, we make three major observations: 1) variance between different annotators is very low; 2) scores under Human 2AFC and Human Clustering order different models similarly; and finally, 3) even though accuracy drops for the “hard” version of ImageNet task, the relative ordering of models remains the same. These observations indicate that alignment can be reliably measured by generating IRIs and does not depend on bias in annotators. Note that AMT experiments were only performed on IRIs generated using the regularizer-free loss in Eq 2.

Impact of regularizer Table 2 shows the results of Alignment (Eq 1) for different regularizers for IRI generation. We evaluated multiple architectures of both standard and adversarially trained [35] CIFAR10 and ImageNet models. We find that under different types of regularizers, the alignment of models can look very different. We also see that adversarial regularizer makes alignment bad for almost all models, thus showing that for the worst pick of IRIs the alignment between learned invariances and human invariances has a lot of room for improvement. Conversely, the human-leaning regularizer overestimates the alignment.

Impact of X In the case of OOD targets (x_t) we see that humans are still able to faithfully judge similarity, yielding the same ranking of models as in-distribution targets. Some results for human judgements about similarity of IRIs for out of distribution samples are shown in Table 4, Appendix A.3. As seen in Fig 4 (Appendix A.3), human-leaning regularizer does not work well for reconstructing noisy targets. This is because such regularizers explicitly remove high-frequency features from reconstructions [42] and thus struggle to meaningfully reconstruct targets that contain high-frequency features. Hence, all results in Table 4, Appendix A.3 are reported on IRIs generated using regularizer-free loss.

Impact of x_0 We repeat some of the experiments with other starting points for Eq 2 and find that results are generally not sensitive to the choice of x_0 . Results are included in Appendix A.5.

3 What Contributes to Learning Invariances Aligned with Humans

There have been many enhancements in the deep learning pipeline that have lead to remarkable generalization performance [31, 21, 59, 57, 5, 28, 25]. In recent years there have been efforts to understand how invariances in representations learnt by such networks align with those of humans [15, 22, 12]. However, how individual components of the deep learning pipeline affect the invariances learned is still not well understood. Prior works claim that adversarially robust models tend to learn representations with a “human prior” [26, 10, 55]. This leads to the question: how do other factors such as architecture, training paradigm, and data augmentation affect the invariances of representations?

We explore these questions in this section. All evaluations in this section are based on regularizer-free IRIs. We chose

CIFAR10						
TRAINING	MODEL	ALIGNMENT			CLEAN ACC.	ROBUST ACC.
		REG.-FREE	HUMAN-ALIGNED	ADVER-SARIAL		
AT $\ell_2, \epsilon = 1$	RESNET18	63.25 \pm 26.23	79.00 \pm 21.94	0.33 \pm 0.47	80.77	50.92
	VGG16	0.25 \pm 0.43	41.41 \pm 16.74	1.00 \pm 1.41	79.84	48.36
	INCEPTIONV3	23.25 \pm 25.56	64.75 \pm 24.17	3.00 \pm 4.24	81.57	51.02
	DENSENET121	82.75 \pm 20.07	86.25 \pm 14.50	1.33 \pm 1.89	83.22	52.86
STANDARD	RESNET18	0.00 \pm 0.00	21.09 \pm 13.51	1.33 \pm 1.89	94.94	0.00
	VGG16	0.00 \pm 0.00	21.88 \pm 14.82	0.00 \pm 0.00	93.63	0.00
	INCEPTIONV3	0.00 \pm 0.00	21.88 \pm 17.54	0.33 \pm 0.47	94.59	0.00
	DENSENET121	0.00 \pm 0.00	26.56 \pm 16.90	0.00 \pm 0.00	95.30	0.00

IMAGENET						
TRAINING	MODEL	ALIGNMENT			CLEAN ACC.	ROBUST ACC.
		REG.-FREE	HUMAN-ALIGNED	ADVER-SARIAL		
AT $\ell_2, \epsilon = 3$	RESNET18	42.00 \pm 38.33	46.75 \pm 39.37	0.33 \pm 0.47	53.12	31.02
	RESNET50	51.00 \pm 34.89	45.75 \pm 37.39	14.00 \pm 3.74	62.83	38.84
	VGG16	55.50 \pm 34.14	55.50 \pm 38.29	11.00 \pm 3.74	56.79	34.46
STANDARD	RESNET18	0.00 \pm 0.00	17.00 \pm 28.30	0.00 \pm 0.00	69.76	0.01
	RESNET50	0.00 \pm 0.00	16.25 \pm 26.42	0.00 \pm 0.00	76.13	0.00
	VGG16	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	73.36	0.16

Table 2: [CIFAR10 and ImageNet Model Alignment Results for Different Regularizers] For different regularizers, we see that ranking of models can look very different. For example, for Adversarially Trained (AT) Resnet18 vs InceptionV3 on CIFAR10, we see that regularizer-free inversion leads to Resnet18 being significantly more aligned, but the trend is much less pronounced for the human-leaning regularizer. We also find that alignment can vary quite a bit between different architectures – all of which achieve similar clean and robust accuracies.

regularizer-free loss over the adversarial loss as the latter shows worst case alignment for all models, which is not useful for understanding the effect of various factors in the deep learning pipeline (Appendix A.4 shows more results using the adversarial regularizer). Similarly, we preferred regularizer-free over human-leaning loss as the latter has a strong ‘bias’ enforced by the regularizer. While our approach generalizes to any layer, unless stated otherwise, all measurements of alignment are on the penultimate layer of the network.

3.1 Architectures and Loss Functions

We test the alignment of different DNNs trained using various loss functions – standard cross-entropy loss, adversarial training (AT), and variants of AT (TRADES [65], MART [63]). Both TRADES and MART have two loss terms – one each for clean and adversarial samples, which are balanced via a hyperparameter β . We report results for multiple values of β in Fig 3a and find that the alignment of standard models (blue squares) is considerably worse than the robust ones (triangles and circles). However, the effect is also influenced by the choice of model architecture, *e.g.*, for CIFAR10, for all robust training losses, VGG16 has significantly lower alignment than other architectures.

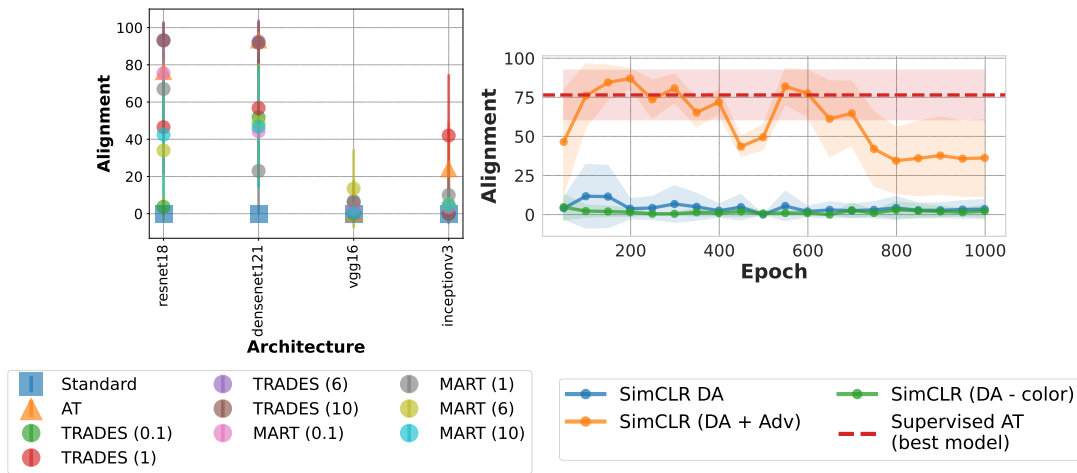
3.2 Data Augmentation

Hand-crafted data augmentations are commonly used in deep learning pipelines. If adversarial training – which augments adversarial samples during training – generally leads to better aligned representations, then how do hand-crafted data augmentations affect invariances of learned representations? For adversarially trained models, we try with and without the usual data augmentation (horizontal flip, color jitter, and

rotation). Since standard models trained with usual data augmentation show poor alignment (Section 3.1), we try stronger data augmentation (a composition of random flip, color jitter, grayscale and gaussian blur, as used in SimCLR [5]) to see if hand-crafted data augmentations can improve alignment. Table 5 Appendix C shows how hand-crafted data augmentation can be crucial in learning aligned representations for some models (*e.g.*, adversarially trained ResNet18 benefits greatly from data augmentation). In other cases data augmentation never hurts the alignment. We also see that standard models do not gain alignment even with stronger hand-crafted data augmentations. CIFAR100 and ImageNet results can be found in Table 6 Appendix C with similar takeaways.

3.3 Learning Paradigm

Since data augmentations (both adversarial and hand-crafted) along with residual architectures (like Resnet18) help alignment, self-supervised learning (SSL) models – which explicitly rely on data augmentations – should learn well aligned representations. This leads to a natural question: how do SSL models compare with the alignment of supervised models? SimCLR [5] is a widely used contrastive SSL method that learns ‘meaningful’ representations without using any labels. Recent works have built on SimCLR to also include adversarial data augmentations [27, 6]. We train both the standard version of SimCLR (using a composition of transforms, as suggested in [6]) and the one with adversarial augmentation on CIFAR10 and compare their alignment with the supervised counterparts. More training details are included in Appendix C. Additionally we also train SimCLR without the color distortion transforms – which were identified as key transforms by the authors [5] – to see how transforms



(a) While adversarially robust models generally have high alignment, we see that different architectures. For example, VGG16 has very low levels of alignment despite being trained using robust training losses. Results on CIFAR100 and Imagenet in Appendix C.

(b) Combining SimCLR’s data augmentations (DA) with adversarial augmentations (Adv) leads to best alignment (in the early and mid epochs) – in some cases even surpassing the best supervised adversarially robust model. Results for more models and datasets in Appendix C.

Figure 3: Role of Loss Function in Alignment (left); Role of Training Paradigm in Alignment (right); ResNet18, CIFAR10

that are crucial for generalization affect alignment. Fig 3b shows the results when comparing self-supervised and supervised learning. We see that SimCLR when trained with both hand-crafted and adversarial augmentations has the best alignment, even outperforming the best adversarially trained supervised model in initial and middle epochs of training. We also see that removing color based augmentations (DA - color) does not have a significant impact on alignment, thus showing that certain DA can be crucial for generalization but not necessarily for alignment.

Summary We find that there are three key components that lead to good alignment: architectures with residual connections, adversarial data augmentation using ℓ_2 threat model, and a (self-supervised) contrastive loss. We leave a more comprehensive study of the effects of these training parameters on alignment for future work.

4 Related Work

Robust Models Several methods have been introduced to make deep learning models robust against adversarial attacks [45, 43, 51, 18, 61, 35, 65, 63, 7]. These works try to model a certain type of human invariance (small change to input that does not change human perception) and make the model also learn such an invariance. Our work, on the other hand, aims to evaluate what invariances have already been learned by a model and how they align with human perception. **Representation Similarity** There has been a long standing interest in comparing neural representations [32, 47, 37, 29, 62, 34, 39, 8]. While these works are related to ours in that they compare two systems of cognition, they assume complete white-box access to both neural networks. In our work, we wish to compare a DNN and a human, with only black-box access to the latter. **DNNs and Human Perception** Neural networks have been used to model many perceptual properties such as quality [1, 14] and

closeness [66] in the image space. Recently there has been interest in measuring the alignment of human and neural network perception. Roads et al. do this by eliciting similarity judgements from humans on ImageNet inputs and comparing it with the outputs of neural nets [49]. Our work, however, explores alignment in the opposite direction, *i.e.*, we measure if inputs that a network sees the same are also the same for humans. [12, 10] are closest to our work as they also evaluate alignment from model to humans, however as discussed in Section 2, unlike our work, their approaches are not scalable, they do not discuss the effects of loss function used to generate IRIs, and they do not contribute to an understanding of what components in the deep learning pipeline lead to learning human-like invariances.

5 Conclusion and Broader Impacts

Our work offers insights into how measures of alignment can vary based on different loss functions used to generate IRIs. We believe that when it is done carefully, measuring alignment is a useful model evaluation tool that provides insights beyond those offered by traditional metrics such as clean and robust accuracy, enabling better alignment of models with humans. We recognize that there are potentially worrying use cases against which we must be vigilant, such as taking advantage of alignment to advance work on deceiving humans. Human perception is complex, nuanced and discontinuous [58], which poses many challenges in measuring the alignment of DNNs with human perception [19]. In this work, we take a step toward defining and measuring the alignment of DNNs with human perception. Our proposed method is a necessary but not sufficient condition for alignment and, thus, must be used carefully and supplemented with other checks, including domain expertise. By presenting this method, we hope for better design, understanding, and auditing of DNNs.

Acknowledgements

AW acknowledges support from a Turing AI Fellowship under grant EP/V025279/1, The Alan Turing Institute, and the Leverhulme Trust via CFI. VN, AM, CK, and KPG were supported in part by an ERC Advanced Grant “Foundations for Fair Social Computing” (no. 789373). VN and JPD were supported in part by NSF CAREER Award IIS-1846237, NSF D-ISN Award #2039862, NSF Award CCF-1852352, NIH R01 Award NLM013039-01, NIST MSE Award #20126334, DARPA GARD #HR00112020007, DoD WHS Award #HQ003420F0035, ARPA-E Award #4334192, ARL Award W911NF2120076 and a Google Faculty Research Award. BCL was supported by Wellcome Trust Investigator Award WT106931MA and Royal Society Wolfson Fellowship 183029. All authors would like to thank Nina Grgić-Hlača for help with setting up AMT surveys.

References

- [1] Amirshahi, S. A.; Pedersen, M.; and Yu, S. X. 2016. Image quality assessment by comparing CNN features between images. *Journal of Imaging Science and Technology*.
- [2] Beery, S.; Van Horn, G.; and Perona, P. 2018. Recognition in terra incognita. In *ECCV*.
- [3] Bengio, Y.; Courville, A.; and Vincent, P. 2013. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8): 1798–1828.
- [4] Buolamwini, J.; and Gebru, T. 2018. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on fairness, accountability and transparency*.
- [5] Chen, T.; Kornblith, S.; Norouzi, M.; and Hinton, G. 2020. A simple framework for contrastive learning of visual representations. In *ICML*.
- [6] Chen, T.; Liu, S.; Chang, S.; Cheng, Y.; Amini, L.; and Wang, Z. 2020. Adversarial robustness: From self-supervised pre-training to fine-tuning. In *CVPR*.
- [7] Cohen, J.; Rosenfeld, E.; and Kolter, Z. 2019. Certified adversarial robustness via randomized smoothing. In *ICML*.
- [8] Ding, F.; Denain, J.-S.; and Steinhardt, J. 2021. Grounding representation similarity with statistical testing. *arXiv:2108.01661*.
- [9] Engstrom, L.; Ilyas, A.; Salman, H.; Santurkar, S.; and Tsipras, D. 2019. Robustness (Python Library).
- [10] Engstrom, L.; Ilyas, A.; Santurkar, S.; Tsipras, D.; Tran, B.; and Madry, A. 2019. Adversarial Robustness as a Prior for Learned Representations. *arXiv:1906.00945*.
- [11] Falcon, W. 2019. PyTorch Lightning. <https://github.com/PyTorchLightning/pytorch-lightning>.
- [12] Feather, J.; Durango, A.; Gonzalez, R.; and McDermott, J. 2019. Metamers of neural networks reveal divergence from human perceptual systems. In *NeurIPS*.
- [13] Fechner, G. T. 1948. Elements of psychophysics, 1860.
- [14] Gao, F.; Wang, Y.; Li, P.; Tan, M.; Yu, J.; and Zhu, Y. 2017. Deepsim: Deep similarity for image quality assessment. *Neurocomputing*.
- [15] Geirhos, R.; Rubisch, P.; Michaelis, C.; Bethge, M.; Wichmann, F. A.; and Brendel, W. 2018. ImageNet-trained CNNs are biased towards texture; increasing shape bias improves accuracy and robustness. *arXiv:1811.12231*.
- [16] Golan, T.; Raju, P. C.; and Kriegeskorte, N. 2020. Controversial stimuli: Pitting neural networks against each other as models of human cognition. *PNAS*.
- [17] Goodfellow, I. J.; Shlens, J.; and Szegedy, C. 2015. Explaining and Harnessing Adversarial Examples. *arXiv:1412.6572*.
- [18] Gu, S.; and Rigazio, L. 2014. Towards deep neural network architectures robust to adversarial examples. *arXiv:1412.5068*.
- [19] Guest, O.; and Love, B. C. 2017. What the success of brain imaging implies about the neural code. *Elife*.
- [20] Harris, C. R.; Millman, K. J.; van der Walt, S. J.; Gommers, R.; Virtanen, P.; Cournapeau, D.; Wieser, E.; Taylor, J.; Berg, S.; Smith, N. J.; Kern, R.; Picus, M.; Hoyer, S.; van Kerkwijk, M. H.; Brett, M.; Haldane, A.; del Río, J. F.; Wiebe, M.; Peterson, P.; Gérard-Marchant, P.; Sheppard, K.; Reddy, T.; Weckesser, W.; Abbasi, H.; Gohlke, C.; and Oliphant, T. E. 2020. Array programming with NumPy. *Nature*.
- [21] He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual learning for image recognition. In *CVPR*.
- [22] Hermann, K.; Chen, T.; and Kornblith, S. 2020. The origins and prevalence of texture bias in convolutional neural networks. *NeurIPS*.
- [23] Hunter, J. D. 2007. Matplotlib: A 2D graphics environment. *Computing in Science & Engineering*.
- [24] Ilyas, A.; Santurkar, S.; Tsipras, D.; Engstrom, L.; Tran, B.; and Madry, A. 2019. Adversarial examples are not bugs, they are features. In *NeurIPS*.
- [25] Ioffe, S.; and Szegedy, C. 2015. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*.
- [26] Kaur, S.; Cohen, J. M.; and Lipton, Z. C. 2019. Are Perceptually-Aligned Gradients a General Property of Robust Classifiers?
- [27] Kim, M.; Tack, J.; and Hwang, S. J. 2020. Adversarial self-supervised contrastive learning. *arXiv:2006.07589*.
- [28] Kingma, D. P.; and Ba, J. 2014. Adam: A method for stochastic optimization. *arXiv:1412.6980*.
- [29] Kornblith, S.; Norouzi, M.; Lee, H.; and Hinton, G. 2019. Similarity of neural network representations revisited. In *ICML*.
- [30] Krizhevsky, A.; Hinton, G.; et al. 2009. Learning multiple layers of features from tiny images.
- [31] Krizhevsky, A.; Sutskever, I.; and Hinton, G. E. 2012. Imagenet classification with deep convolutional neural networks. *NeurIPS*.

- [32] Laakso, A.; and Cottrell, G. 2000. Content and cluster analysis: assessing representational similarity in neural systems. *Philosophical psychology*, 13(1): 47–76.
- [33] LeCun, Y.; Bengio, Y.; and Hinton, G. 2015. Deep learning. *nature*, 521(7553): 436–444.
- [34] Li, Y.; Yosinski, J.; Clune, J.; Lipson, H.; and Hopcroft, J. 2016. Convergent Learning: Do different neural networks learn the same representations? [arXiv:1511.07543](https://arxiv.org/abs/1511.07543).
- [35] Madry, A.; Makelov, A.; Schmidt, L.; Tsipras, D.; and Vladu, A. 2019. Towards Deep Learning Models Resistant to Adversarial Attacks. [arXiv:1706.06083](https://arxiv.org/abs/1706.06083).
- [36] Mahendran, A.; and Vedaldi, A. 2014. Understanding Deep Image Representations by Inverting Them. [arXiv:1412.0035](https://arxiv.org/abs/1412.0035).
- [37] Morcos, A. S.; Raghu, M.; and Bengio, S. 2018. Insights on representational similarity in neural networks with canonical correlation. In *NeurIPS*.
- [38] Mordvintsev, A.; Olah, C.; and Tyka, M. 2015. Inceptionism: Going Deeper into Neural Networks.
- [39] Nanda, V.; Speicher, T.; Kolling, C.; Dickerson, J. P.; Gummadi, K.; and Weller, A. 2022. Measuring Representational Robustness of Neural Networks Through Shared Invariances. In *ICML*.
- [40] Nguyen, A.; Yosinski, J.; and Clune, J. 2015. Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. In *CVPR*.
- [41] Olah, C.; Cammarata, N.; Schubert, L.; Goh, G.; Petrov, M.; and Carter, S. 2020. Zoom In: An Introduction to Circuits. *Distill*.
- [42] Olah, C.; Mordvintsev, A.; and Schubert, L. 2017. Feature Visualization. *Distill*.
- [43] Papernot, N.; McDaniel, P.; Goodfellow, I.; Jha, S.; Celik, Z. B.; and Swami, A. 2017. Practical black-box attacks against machine learning. In *Asia CCS*.
- [44] Papernot, N.; McDaniel, P.; Jha, S.; Fredrikson, M.; Celik, Z. B.; and Swami, A. 2016. The limitations of deep learning in adversarial settings. In *2016 EuroS&P*, 372–387. IEEE.
- [45] Papernot, N.; McDaniel, P.; Wu, X.; Jha, S.; and Swami, A. 2016. Distillation as a defense to adversarial perturbations against deep neural networks. In *S&P*.
- [46] Paszke, A.; Gross, S.; Massa, F.; Lerer, A.; Bradbury, J.; Chanan, G.; Killeen, T.; Lin, Z.; Gimelshein, N.; Antiga, L.; Desmaison, A.; Kopf, A.; Yang, E.; DeVito, Z.; Raison, M.; Tejani, A.; Chilamkurthy, S.; Steiner, B.; Fang, L.; Bai, J.; and Chintala, S. 2019. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *NeurIPS*.
- [47] Raghu, M.; Gilmer, J.; Yosinski, J.; and Sohl-Dickstein, J. 2017. SVCCA: Singular Vector Canonical Correlation Analysis for Deep Learning Dynamics and Interpretability. In *NeurIPS*.
- [48] Recht, B.; Roelofs, R.; Schmidt, L.; and Shankar, V. 2019. Do imagenet classifiers generalize to imagenet? In *ICML*.
- [49] Roads, B. D.; and Love, B. C. 2021. Enriching ImageNet With Human Similarity Judgments and Psychological Embeddings. In *CVPR*.
- [50] Rosenfeld, A.; Zemel, R.; and Tsotsos, J. K. 2018. The elephant in the room. [arXiv:1808.03305](https://arxiv.org/abs/1808.03305).
- [51] Ross, A. S.; and Doshi-Velez, F. 2018. Improving the adversarial robustness and interpretability of deep neural networks by regularizing their input gradients. In *AAAI*.
- [52] Russakovsky, O.; Deng, J.; Su, H.; Krause, J.; Satheesh, S.; Ma, S.; Huang, Z.; Karpathy, A.; Khosla, A.; Bernstein, M.; et al. 2015. Imagenet large scale visual recognition challenge. *IJCV*.
- [53] Sagawa, S.; Koh, P. W.; Hashimoto, T. B.; and Liang, P. 2019. Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization. [arXiv:1911.08731](https://arxiv.org/abs/1911.08731).
- [54] Salman, H.; Ilyas, A.; Engstrom, L.; Kapoor, A.; and Madry, A. 2020. Do adversarially robust imagenet models transfer better? In *NeurIPS*.
- [55] Santurkar, S.; Ilyas, A.; Tsipras, D.; Engstrom, L.; Tran, B.; and Madry, A. 2019. Image Synthesis with a Single (Robust) Classifier. In *NeurIPS*.
- [56] Shafahi, A.; Najibi, M.; Ghiasi, A.; Xu, Z.; Dickerson, J.; Studer, C.; Davis, L. S.; Taylor, G.; and Goldstein, T. 2019. Adversarial training for free! [arXiv:1904.12843](https://arxiv.org/abs/1904.12843).
- [57] Simonyan, K.; and Zisserman, A. 2014. Very deep convolutional networks for large-scale image recognition. [arXiv:1409.1556](https://arxiv.org/abs/1409.1556).
- [58] Stankiewicz, B. J.; and Hummel, J. E. 1996. Categorical relations in shape perception. *Spatial vision*.
- [59] Szegedy, C.; Liu, W.; Jia, Y.; Sermanet, P.; Reed, S.; Anguelov, D.; Erhan, D.; Vanhoucke, V.; and Rabinovich, A. 2015. Going deeper with convolutions. In *CVPR*.
- [60] Taori, R.; Dave, A.; Shankar, V.; Carlini, N.; Recht, B.; and Schmidt, L. 2020. Measuring robustness to natural distribution shifts in image classification. In *NeurIPS*.
- [61] Tramèr, F.; Kurakin, A.; Papernot, N.; Goodfellow, I.; Boneh, D.; and McDaniel, P. 2017. Ensemble adversarial training: Attacks and defenses. [arXiv:1705.07204](https://arxiv.org/abs/1705.07204).
- [62] Wang, L.; Hu, L.; Gu, J.; Wu, Y.; Hu, Z.; He, K.; and Hopcroft, J. 2018. Towards understanding learning representations: To what extent do different neural networks learn the same representation. [arXiv:1810.11750](https://arxiv.org/abs/1810.11750).
- [63] Wang, Y.; Zou, D.; Yi, J.; Bailey, J.; Ma, X.; and Gu, Q. 2019. Improving adversarial robustness requires revisiting misclassified examples. In *ICLR*.
- [64] Wightman, R. 2019. PyTorch Image Models. <https://github.com/rwightman/pytorch-image-models>.
- [65] Zhang, H.; Yu, Y.; Jiao, J.; Xing, E.; El Ghaoui, L.; and Jordan, M. 2019. Theoretically principled trade-off between robustness and accuracy. In *ICML*.
- [66] Zhang, R.; Isola, P.; Efros, A. A.; Shechtman, E.; and Wang, O. 2018. The Unreasonable Effectiveness of Deep Features as a Perceptual Metric. In *CVPR*.

A Measuring Human Alignment via Representation Inversion

A.1 Measuring Human Perception Similarity

We recruited AMT workers with completion rate $\geq 95\%$ and who spoke English. To further ensure that the workers understood the task, we added attention checks. For 2AFC task, this meant making the query image as the same image as one of the images in the option. In clustering setting, this meant making the image on the row same as one of the images in the columns. All the workers who took our survey passed the attention checks.

We estimated a completion time of about 30 minutes for each survey and thus paid each worker 7.5\$. We allotted 60 minutes per survey, so workers are not forced to rush through the survey. Most of the workers were able to complete the task in less than 20 minutes. Our study was approved by the Ethical Review Board of our institute.

ImageNet Clustering Hard In order to create the hard ImageNet clustering task we use the human annotations of similarity between ImageNet images collected by ImageNet-HSJ authors [49]. This contains a matrix (M) of similarity scores for each image in ImageNet validation set, where M_{ij} is an indication for similarity between i^{th} and j^{th} images. For each image i (randomly picked), we sample two more images that are the most similar to i as per M_i . This creates a task that is much harder to perform for human annotators since the images on the columns look perceptually very similar (See Fig 2b for an example).

A.2 Using LPIPS as a proxy for g_{human}

In order to ensure that LPIPS [66] is a reliable proxy to simulate human perception we measured if LPIPS could simulate human annotators on two perceptual similarity task setups: 2AFC and Clustering, as described in Section 2.1. For 2AFC, this meant using LPIPS to measure the distance between the query image and the two images shown in the options and then matching the query image to the one with lesser LPIPS distance. And similarly in Clustering, for each image on the row, we used LPIPS to measure its distance from each of the 3 images in the column and then matched it to the closest one. Our results (Table 1 in main paper) show that this can serve as a good proxy for human perception similarity. For all our experiments, we report average over 4 different LPIPS backbones: ImageNet trained Alexnet & VGG16, and both of the Imagenet trained Alexnet & VGG16 finetuned by the authors for perceptual similarity <https://github.com/richzhang/PerceptualSimilarity>.

A.3 Role of Input Distribution

Fig 4 shows some examples of inputs sampled from different gaussians (the lighter ones are sampled from $\mathcal{N}(0.5, 2)$ and darker ones from $\mathcal{N}(0, 1)$). We find that completing 2AFC and Clustering tasks on inputs that look like noise to humans is a qualitatively harder task than when doing this on in-distribution target samples.

However, remarkably, we observe that humans are still able to bring out the differences between different models, even when given (a harder) task of matching re-constructed

noisy inputs. Table 4 shows the results of surveys conducted with noisy target samples. It’s worth noting that the accuracy of humans drop quite a bit from in-distribution targets, thus indicating that this is indeed a harder task.

A.4 Regularizers

For the human-aligned regularizers, we use the ones discussed in [42]. These fall into three broad categories: *frequency penalization*, *transformation robustness*, and *pre-conditioning*.

- **Frequency Penalization:** The goal is to explicitly penalize high-frequency features in the reconstruction (x_r). This is done by adding a regularizer of the form $\mathcal{R}(x) = TV(x) + \|x\|_p$, where TV is the total variation and $p = 1$ [36]. A similar effect of frequency penalization can also be done by ensuring robustness of x_r to blurring, *i.e.*, $\mathcal{R}(x) = \|x - \text{StopGradient}(\text{blur}(x))\|_2^2$ [40].
- **Transformation Robustness:** This ensures that x_r is such that the representation is same even if we slightly transform x_r . This is achieved by replacing x with $T(x)$ in Eq 2. We use T as a composition of color jitter, random scaling, and random rotation.
- **Pre-conditioning:** This involves taking gradient steps in the fourier domain, which decorrelates the pixels in x_r .

For our experiments we find that transformation robustness generates the best looking x_r and thus we report results under *human-learning regularizer* based on x_r generated using *transformation robustness* during representation inversion.

For adversarial regularizer, we report results in Table 3 and find that such a regularizer can make almost all models look like they have bad alignment.

A.5 Role of x_0

We additionally report results for the adversarial regularizer where IRIs were generated from a separate seed. While experiments in the main paper reported for a seed sampled from $\mathcal{N}(0, 1)$, we report results here for a seed sampled from $\mathcal{N}(0, 0.01)$ in Table 3 and find that regardless of seed, adversarial regularizer makes all models look bad.

B Model, Code, Assets, and Compute Details

B.1 Code and Assets

In our code we make use of many open source libraries such as timm [64], pytorch [46], pytorch-lightning [11], numpy [20], robustness [9], matplotlib [23]. timm, pytorch-lightning and have an Apache 2.0 license. Numpy has a BSD 3-Clause License. Robustness has an MIT license. PyTorch’s license can be found here: <https://github.com/pytorch/pytorch/blob/master/LICENSE>, and matplotlib’s here: <https://github.com/matplotlib/matplotlib/blob/main/LICENSE/LICENSE>. All these licenses allow free use, modification and distribution. We use publicly available academic datasets CIFAR10/100 [30] and ImageNet [52].

CIFAR10		
	$x_0 \sim \mathcal{N}(0, 1)$	$x_0 \sim \mathcal{N}(0, 0.01)$
MODEL	ALIGNMENT ON ADVERSARIAL IRIS	
SUPERVISED RESNET18, STANDARD	1.33±1.89	0.00±0.00
SUPERVISED VGG16, STANDARD	0.00±0.00	0.00±0.00
SUPERVISED INCEPTIONV3, STANDARD	0.33±0.47	0.00±0.00
SUPERVISED DENSENET121, STANDARD	0.00±0.00	0.00±0.00
SUPERVISED RESNET18, AT $\epsilon(\ell_2) = 1$	0.33±0.47	1.00±0.00
SUPERVISED VGG16, AT $\epsilon(\ell_2) = 1$	1.00±1.41	1.00±0.00
SUPERVISED INCEPTIONV3, AT $\epsilon(\ell_2) = 1$	3.00±4.24	0.00±0.00
SUPERVISED DENSENET121, AT $\epsilon(\ell_2) = 1$	1.33±1.89	1.00±0.00
SUPERVISED RESNET18, SIMCLR DA	0.00±0.00	0.00±0.00
SIMCLR RESNET18, STANDARD DA	0.00±0.00	1.00±0.00
SIMCLR RESNET18, DA WITHOUT COLOR	0.00±0.00	0.00±0.00
SIMCLR RESNET18, STANDARD + ADV DA $\epsilon(\ell_2) = 1$	0.00±0.00	0.00±0.00
IMAGENET		
	$x_0 \sim \mathcal{N}(0, 1)$	$x_0 \sim \mathcal{N}(0, 0.01)$
MODEL	ALIGNMENT ON ADVERSARIAL IRIS	
SUPERVISED RESNET18, STANDARD	0.00±0.00	3.00±0.00
SUPERVISED RESNET50, STANDARD	0.00±0.00	3.50±0.50
SUPERVISED VGG16, STANDARD	0.00±0.00	3.00±2.00
SUPERVISED RESNET18, AT $\epsilon(\ell_2) = 3$	0.33±0.47	2.50±1.50
SUPERVISED RESNET50, AT $\epsilon(\ell_2) = 3$	14.00±3.74	3.50±0.50
SUPERVISED VGG16, AT $\epsilon(\ell_2) = 3$	11.00±3.74	4.50±1.50

Table 3: [Adversarial IRIs] We observe that using the adversarial regularizer (described in Section 2.2) makes alignment for all models look bad. AT = Adversarial Training, DA = Data Augmentations. For details about standard SimCLR DA and DA without color, see Section C.

CIFAR10			
	IN-DIST.	NOISE $\mathcal{N}(0, 1)$	NOISE $\mathcal{N}(0.5, 2)$
MODEL	HUMAN 2AFC	HUMAN 2AFC	HUMAN 2AFC
RESNET18	96.00±2.55	31.053±15.912	64.386±7.310
VGG16	38.83±7.59	0.351±0.496	0.351±0.496
INCEPTIONV3	82.00±8.44	28.772±20.476	2.105±1.549
DENSENET121	98.67±0.24	60.702±11.413	68.947±16.119
IMAGENET			
	IN-DIST.	NOISE $\mathcal{N}(0, 1)$	NOISE $\mathcal{N}(0.5, 2)$
MODEL	HUMAN 2AFC	HUMAN 2AFC	HUMAN 2AFC
RESNET18	93.17±5.95	28.748±34.491	65.079±5.616
RESNET50	99.50±0.00	70.745±7.409	93.617±9.027
VGG16	95.50±2.12	29.806±19.704	46.914±13.827

Table 4: [CIFAR10 and ImageNet In-Distr vs OOD Survey Results] We observe that even on OOD samples that look like noise, humans can still bring out the relative differences between models, e.g., densenet121 on CIFAR10 is still ranked best aligned model on targets sampled from both kinds of noise. Reduced accuracy of humans on noise shows that this identifying similarities between IRIs on OOD samples is a harder task than with in-distribution target samples.

B.2 Models

Supervised We used VGG16, ResNet18, Densenet121 and InceptionV3 for experiments on CIFAR10 and CIFAR100. The “robust” version of these models were trained using ad-

versarial training [35], with an ℓ_2, ϵ of 1. All these models were trained using the standard data augmentations (a composition of RandomCrop, RandomHorizontalFlip, ColorJitter, RandomRotation). For ImageNet, we

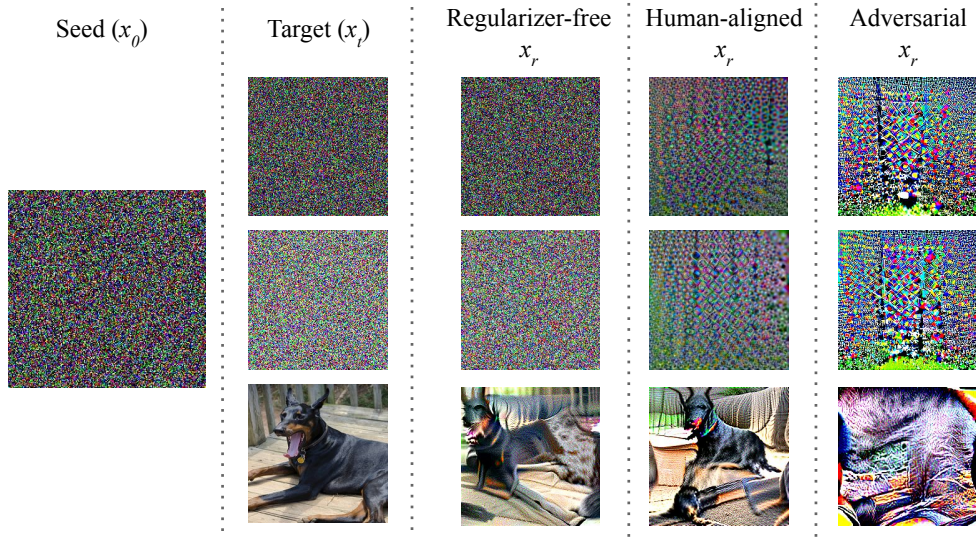


Figure 4: [In vs Out of Distribution Samples] Examples of reconstructions of in-distribution (bottom row) sample vs out-of-distribution samples for an ImageNet trained ResNet50 using the three regularizers mentioned in Section 2.2. Here the OOD targets are sampled from two separate random gaussians $\mathcal{N}(0, 1)$ (top row) and $\mathcal{N}(0.5, 2)$ (second row). We see that similar to the in-distribution sample, regularizer-free and adversarial inversions result in x_r resembling and differing from x_t respectively. Interestingly, for human-aligned regularizer, which explicitly tries to remove high-frequency features from x_r , fails to reconstruct an x_t that itself consists of high-frequency features.

used VGG16, ResNet18 and ResNet50 and the “robust” versions of these models were taken from [54] with an ℓ_2, ϵ of 3. ImageNet models used slightly different data augmentations: RandomHorizontalFlip ColorJitter Lighting

Self-supervised We used SimCLR [5] to train a ResNet18 backbone on CIFAR10 and CIFAR100. More details about different types of data augmentations in Section C.

ImageNet We used VGG16 (with batchnorm), ResNet18 and ResNet50 for ImageNet. The “robust” versions of these models were taken from [54], who trained these models using adversarial training with an ℓ_2 epsilon of 3.

B.3 Compute Details

We used our institute’s GPU cluster to run all experiments. Since our experiments involve standard models and datasets, these can be run on any hardware supported by PyTorch. In our case, we used 5 machines with 2 V100 Nvidia Tesla GPUs (32GB each, volta architecture) and a Nvidia dgx machine with 8 Nvidia Tesla P100 GPUs (16GB each). We estimate a total of 500+ GPU hours.

C What Contributes to Good Alignment

SimCLR training details We used data augmentations shown to work best by the authors (a composition of RandomHorizontalFlip, ColorJitter, RandomGrayscale and GaussianBlur, as implemented in the original codebase <https://github.com/google-research/simclr>). We also train SimCLR models without the color augmentations (*i.e.* only RandomHorizontalFlip and GaussianBlur). Since color transforms were crucial for obtaining representations with good generalization performance, we wanted to analyze how removing augmentations crucial for generalization impacts alignment. Finally, we also train a variant of SimCLR with adversarial data augmentations, as proposed in some recent works [27, 6]. As opposed

to traditional adversarial training, here we generate adversarial data augmentations for a model (g) by solving the following maximization for each input x :

$$\operatorname{argmax}_{x'} \|g(x') - g(x)\|_2 \text{ st } \|x - x'\|_2 \leq \epsilon$$

For our experiments $\epsilon = 1$ for CIFAR10 and $\epsilon = 3$ for ImageNet (similar to supervised models).

Architectures and Loss Function, CIFAR100 & ImageNet Fig 6 shows results for CIFAR100 and ImageNet for standard and robust training. Similar to previous works, we find that robust models are better aligned with human perception. Interestingly, we find that the variance between different architectures that we observed for CIFAR10 does not exist for CIFAR100 and ImageNet, *i.e.*, regardless of architecture, robustly trained models are well aligned with human perception. Indicating that (unsurprisingly) training dataset also plays a major role in alignment.

SimCLR, CIFAR100 Fig 5 shows alignment of different types of SimCLR models throughout training. We observe a similar trend as CIFAR10, where adversarial data augmentation improves alignment.

C.1 Data Augmentation, CIFAR10, CIFAR100 & ImageNet

Since re-training ImageNet models with adversarial training is very resource intensive, we train ImageNet models using Free Adversarial Training (Free AT) [56]. Free AT only has an implementation for ℓ_{inf} threat model, hence we train these models with $\ell_{\text{inf}}, \epsilon = 4/255$ (for both with and without data augmentation). Table 6 shows that similar to CIFAR10 (Table 5), for some models, like ResNet18, data augmentation is crucial in learning aligned representations (despite being trained to be adversarially robust). For other models, data augmentation never hurts alignment (except InceptionV3 for CIFAR100).

ADV. TRAINING				
	RESNET18	DENSENET121	VGG16	INCEPTIONV3
USUAL DATA AUG	76.50 \pm 15.91	93.50 \pm 9.60	0.25 \pm 0.43	24.25 \pm 25.17
NO DATA AUG	30.00 \pm 12.02	93.75 \pm 8.20	1.00 \pm 1.73	12.25 \pm 20.08
STANDARD				
	RESNET18	DENSENET121	VGG16	INCEPTIONV3
STRONG DATA AUG	0.00 \pm 0.00	1.00 \pm 1.73	0.00 \pm 0.00	0.00 \pm 0.00
USUAL DATA AUG	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00

Table 5: **[CIFAR10 Models; Effect of Data Augmentation]** For certain models, *e.g.*, adversarially trained resnet18, data augmentation is crucial in learning aligned representations.

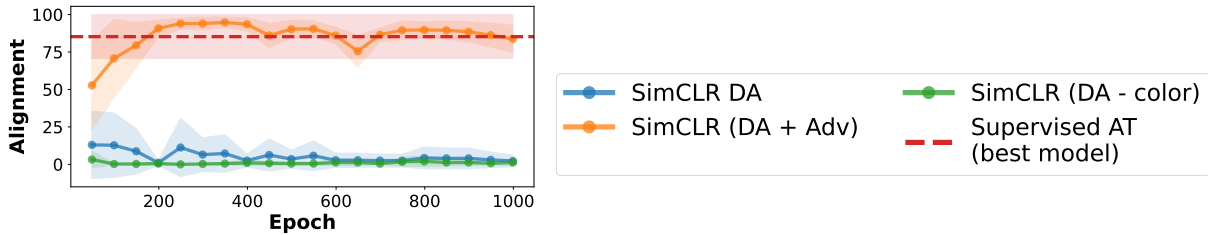


Figure 5: ResNet18 backbone trained using SimCLR on CIFAR100.

CIFAR100				
	RESNET18	DENSENET121	VGG16	INCEPTIONV3
USUAL DATA AUG	82.50 \pm 20.85	88.00 \pm 14.51	58.75 \pm 32.15	69.25 \pm 26.39
NO DATA AUG	81.50 \pm 15.58	89.75 \pm 15.01	46.25 \pm 32.25	84.75 \pm 13.70
IMAGENET				
	RESNET18		RESNET50	
USUAL DATA AUG	13.00 \pm 22.52		0.00 \pm 0.00	
NO DATA AUG	0.75 \pm 1.30		0.00 \pm 0.00	

Table 6: **[CIFAR100 & ImageNet Models all trained to be adversarially robust; Effect of Data Augmentation]** Similar to CIFAR10, data augmentation is crucial for adversarially trained resnet18 to learn aligned representations.

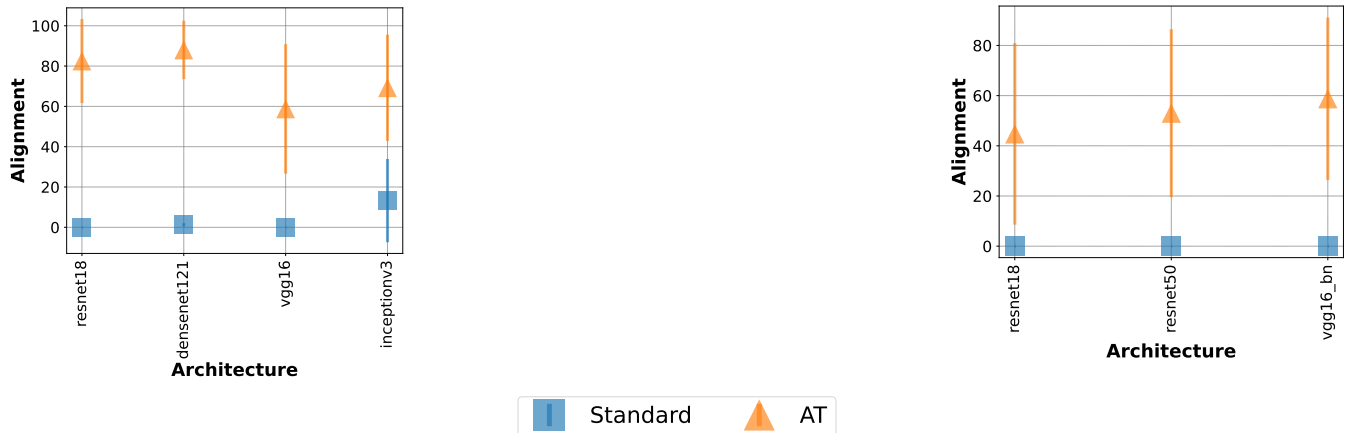


Figure 6: Role of Loss Function in Alignment; CIFAR100 (left), and ImageNet (right)