
Joint rotational invariance and adversarial training of a dual-stream Transformer yields state of the art Brain-Score for Area V4

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Modern high-scoring models of vision in the brain score competition do not stem
2 from Vision Transformers. However, in this paper, we provide evidence against
3 the unexpected trend of Vision Transformers (ViT) being not perceptually aligned
4 with human visual representations by showing how a dual-stream Transformer, a
5 CrossViT *a la* Chen et al. (2021), under a joint rotationally-invariant and adver-
6 sarial optimization procedure yields 2nd place in the aggregate Brain-Score 2022
7 competition (Schrimpf et al., 2020b) averaged across all visual categories, and at
8 the time of the competition held 1st place for the highest explainable variance of
9 area V4. In addition, our current Transformer-based model also achieves greater
10 explainable variance for areas V4, IT and Behavior than a biologically-inspired
11 CNN (ResNet50) that integrates a frontal V1-like computation module (Dapello
12 et al., 2020). To assess the contribution of the optimization scheme with respect
13 to the CrossViT architecture, we perform several additional experiments on differ-
14 ently optimized CrossViT’s regarding adversarial robustness, common corruption
15 benchmarks, mid-ventral stimuli interpretation and feature inversion. Against our
16 initial expectations, our family of results provides tentative support for an “*All
17 roads lead to Rome*” argument enforced via a joint optimization rule even for non
18 biologically-motivated models of vision such as Vision Transformers.

19 1 Introduction

20 Research and design of modern deep learning and computer vision systems such as the NeoCogni-
21 tron (Fukushima & Miyake, 1982), H-Max Model (Serre et al., 2005) and classical CNNs (LeCun
22 et al., 2015) have often stemmed from breakthroughs in visual neuroscience dating from Kuffler
23 (1953) and Hubel & Wiesel (1962). Today, research in neuroscience passes through a phase of
24 symbiotic development where several models of artificial visual computation (mainly deep neural
25 networks), may inform visual neuroscience (Richards et al., 2019) shedding light on puzzles of
26 development (Lindsey et al., 2019), physiology (Dapello et al., 2020), representation (Jagadeesh &
27 Gardner, 2022) and perception (Harrington & Deza, 2022).

28 Of particular recent interest is the development of Vision Transformers (Dosovitskiy et al., 2021). A
29 model that originally generated several great breakthroughs in natural language processing (Vaswani
30 et al., 2017), and that has now slowly begun to dominate the field of machine visual computation.
31 However, in computer vision, we still do not understand why Vision Transformers perform so well
32 when adapted to the visual domain (Bhojanapalli et al., 2021). Is this new excel in performance
33 due to their self-attention mechanism; a relaxation of their weight-sharing constraint? Their greater
34 number of parameters? Their optimization procedure? Or perhaps a combination of all these factors?
35 Naturally, given the uncertainty of the models’ *explainability*, their use has been carefully limited as
36 a model of visual computation in biological (human) vision.

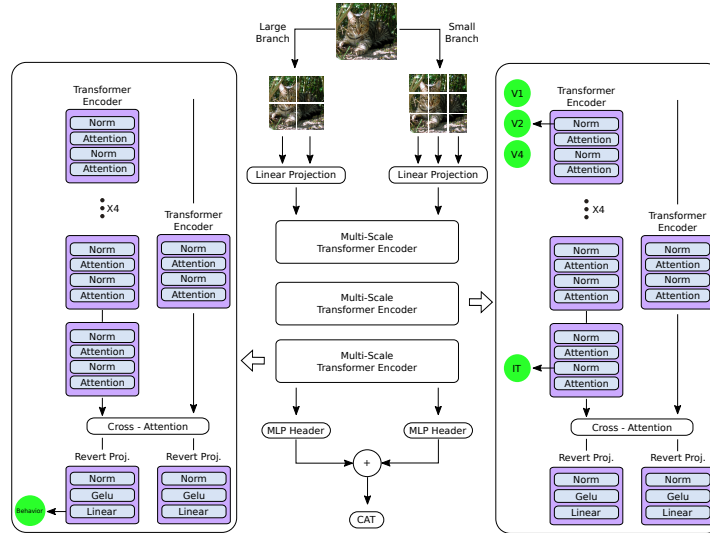


Figure 1: Diagram of CrossViT-18[†] (Chen et al., 2021) architecture & specification of selected layers for the V1, V2, V4, IT brain areas and the behavioral benchmark. Our Brain-Score 2022 competition entry was a variation of this model where the architecture is cloned, and the network is adversarially trained with hard data-augmentation rotations starting from a pre-trained ImageNet model.

37 This is a double-edged sword: On one hand, perceptual psychologists still rely heavily on relatively
 38 low-scoring ImageNet-based accuracy models such as AlexNet, ResNet & VGG despite their *limited*
 39 degree of biological plausibility (though some operations are preserved, *eg.* local filtering, half-wave
 40 rectification, pooling). On the other hand, a new breed of models such as Vision Transformers has
 41 surged, but their somewhat non-biologically inspired computations have no straightforward mapping
 42 to approximate the structure of the human ventral stream¹ – thus discarding them as serious models of
 43 the human visual system. Alas, even if computer vision scientists may want to remain on the sidelines
 44 of the usefulness of a biological/non-biological plausibility debate, the reality is that computer vision
 45 systems are still far from perfect. The existence of Adversarial examples, both artificial (Goodfellow
 46 et al., 2015; Szegedy et al., 2014) and natural (Hendrycks et al., 2021b), reflects that there is still
 47 a long way to go to close the human-machine perceptual alignment gap (Geirhos et al., 2021).
 48 Beyond the theoretical milestone of closing this gap, this will be beneficial for automated systems in
 49 radiology (Hosny et al., 2018), surveillance (Deza et al., 2019), driving (Huang & Chen, 2020), and
 50 art (Ramesh et al., 2022).

51 These two lines of thought bring us to an interesting question that was one of the motivations of this
 52 paper: “*Are Vision Transformers good models of the human ventral stream?*” Our approach to answer
 53 this question will rely on using the Brain-Score platform (Schrimpf et al., 2020a; BrainScore-Org,
 54 2022) and participating in their first yearly competition with a Transformer-based model. This
 55 platform quantifies the similarity via bounded [0,1] scores of responses between a computer model
 56 and a set of non-human primates. Here the ground truth is collected via neurophysiological recordings
 57 and/or behavioral outputs when primates are performing psychophysical tasks, and the scores are
 58 computed by some derivation of Representational Similarity Analysis (Kriegeskorte et al., 2008)
 59 when pitted against artificial neural network activations of modern computer vision models.

60 Altogether, if we find that a specific model yields high Brain-Scores, this may suggest that such flavor
 61 of Vision Transformers-based models obey a necessary but not sufficient condition of biological
 62 plausibility – or at least relatively so with respect to their Convolutional Neural Network (CNN)
 63 counter-parts. As it turns out, we will find out that the answer to the previously posed question
 64 is complex, and depends heavily on how the artificial model is optimized (trained). Thus the
 65 main contribution of this paper is to understand *why* this particular Transformer-based model when
 66 optimized under certain conditions performs vastly better in the Brain-Score competition achieving

¹Even at their start, the patch embedding operation is not obviously mappable to retinal, LGN, or V1-like primate computation.

Rank	Model ID #	Description	Brain-Score						ρ -Hierarchy
			Avg	V1	V2	V4	IT	Behavior	
1	1033	Bag of Tricks (Riedel, 2022) [New SOTA]	0.515	0.568	0.360	0.481	0.514	0.652	-0.2
2	991	CrossViT-18† (Adv + Rot) [Ours]	0.488	0.493	0.342	0.514	0.531	0.562	+0.8
3	1044	Gated Recurrence (Azeglio et al., 2022)	0.463	0.509	0.303	0.482	0.467	0.554	-0.4
4	896	N/A	0.456	0.538	0.336	0.485	0.459	0.461	-0.4
5	1031	N/A	0.453	0.539	0.332	0.475	0.510	0.410	-0.2

Table 1: Ranking of all entries in the Brain-Score 2022 competition as of February 28th, 2022. Scores in **blue** indicate **world record** (highest of all models at the time of the competition), while scores in **bold** display the highest scores of **competing entries**. Column ρ -Hierarchy indicates the Spearman rank correlation between per-Area Brain-Score and Depth of Visual Area (V1 \rightarrow IT).

67 SOTA in such benchmark, and *not* to develop another competitive/SOTA model for ImageNet (which
68 has shown to not be a good target Beyer et al. (2020)). The authors firmly believe that the former goal
69 tackled in the paper is much under-explored compared to the latter, and is also of great importance to
70 the intersection of the visual neuroscience and machine learning communities.

71 2 Optimizing a CrossViT for the Brain-Score Competition

72 Now, we discuss an interesting finding, where amidst the constant debate of the biological plausibility
73 of Vision Transformers – which have been deemed less biologically plausible than convolutional
74 neural networks², though also see Conwell et al. (2021)) –, we find that when these Transformers are
75 optimized under certain conditions, they may achieve high explainable variance with regards to many
76 areas in primate vision, and surprisingly the highest score to date at the time of the competition for
77 explainable variance in area V4, that still remains a mystery in visual neuroscience (see Pasupathy
78 et al. (2020) for a review). Our final model and highest scoring model was based on several insights:

79 **Adversarial-Training:** Work by Santurkar et al. (2019); Engstrom et al. (2019b); Dapello et al.
80 (2020), has shown that convolutional neural networks trained adversarially³ yield human perceptually-
81 aligned distortions when attacked. This is an interesting finding, that perhaps extends to vision
82 transformers, but has never been qualitatively tested before though recent works – including this
83 one (See Figure 4) – have started to investigate in this direction (Tuli et al., 2021; Caro et al., 2020).
84 Thus we projected that once we picked a specific vision transformer architecture, we would train it
85 adversarially.

86 **Multi-Resolution:** Pyramid approaches (Burt & Adelson, 1987; Simoncelli & Freeman, 1995; Heeger
87 & Bergen, 1995) have been shown to correlate highly with good models of Brain-Scores (Marques
88 et al., 2021). We devised that our Transformer had to incorporate this type of processing either
89 implicitly or explicitly in its architecture.

90 **Rotation Invariance:** Object identification is generally rotationally invariant (depending on the
91 category; *e.g.* not the case for faces (Kanwisher et al., 1998)). So we implicitly trained our model to
92 take in different rotated object samples via hard rotation-based data augmentation. This procedure is
93 different from pioneering work of Ecker et al. (2019) which explicitly added rotation equivariance to
94 a convolutional neural network.

95 **Localized texture-based computation:** Despite the emergence of a *global* texture-bias in object
96 recognition when training Deep Neural Networks (Geirhos et al., 2019) – object recognition is a
97 compositional process (Brendel & Bethge, 2019; Deza et al., 2020). Recently, works in neuroscience
98 have also suggested that *local* texture computation is perhaps pivotal for object recognition to either
99 create an ideal basis set from which to represent objects (Long et al., 2018; Jagadeesh & Gardner,
100 2022) and/or encode robust representations (Harrington & Deza, 2022).

101 After searching for several models in the computer vision literature that resemble a Transformer
102 model that ticks all the boxes above, we opted for a CrossViT-18† (that includes multi-resolution
103 + local texture-based computation) that was trained with rotation-based augmentations and also

²Discussed in: URL_1 URL_2

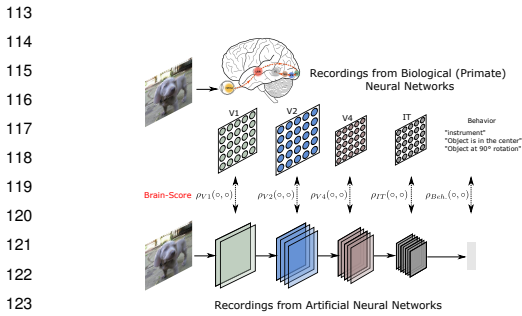
³Adversarial training is the process in which an image in the training distribution of a network is perturbed adversarially (*e.g.* via PGD); the perturbed image is re-labeled to its original non-perturbed class, and the network is optimized via Empirical Risk Minimization (Madry et al., 2018).

Model ID #	Description	ImageNet (\uparrow)	Brain-Score (\uparrow)					
		Validation Accuracy (%)	Avg	V1	V2	V4	IT	Behavior
N/A	Pixels (Baseline)	N/A	0.053	0.158	0.003	0.048	0.035	0.020
N/A	AlexNet (Baseline)	63.3	0.424	0.508	0.353	0.443	0.447	0.370
N/A	VOneResNet50-robust (SOTA)	71.7	0.492	0.531	0.391	0.471	0.522	0.545
991	CrossViT-18 \dagger (Adv + Rot)	73.53	0.488	0.493	0.342	0.514	0.531	0.562
1084	CrossViT-18 \dagger (Adv)	64.60	0.462	0.497	0.343	0.508	0.519	0.441
1095	CrossViT-18 \dagger (Rot)	79.22	0.458	0.458	0.288	0.495	0.503	0.547
1057	CrossViT-18 \dagger	83.05	0.442	0.473	0.274	0.478	0.484	0.500

Table 2: A list of different models submitted to the Brain-Score 2022 competition. Scores in **bold** indicate the highest performing model per column. Scores in **blue** indicate **world record** (highest of all models at the time of the competition). All CrossViT-18 \dagger entries in the table are ours.

104 adversarial training (See Appendix A.3 for exact training details, our *best* model also used $p = 0.25$
 105 grayscale augmentation, though this contribution to model Brain-Score is minimal).

106 **Results:** Our best performing model #991 achieved 2nd place in the overall Brain-Score 2022
 107 competition (Schrimpf et al., 2020b) as shown in Table 1. At the time of submission, it holds the
 108 first place for the highest explainable variance of area V4 and the second highest score in the IT area.
 109 Our model also currently ranks 6th across all Brain-Score submitted models as shown on the main
 110 brain-score website (including those outside the competition and since the start of the platform’s
 111 conception, totaling 216). A general schematic of how Brain-Scores are calculated can be seen in
 112 Figure 2.



124 Figure 2: A schematic of how brain-score is calculated as similarity metrics obtained from neural
 125 responses and model activations.
 126
 127
 128

Additionally, in comparison with the biologically-inspired model (VOneRes-
 Net50+ Adv. training), our model achieves greater scores in the IT, V4 and Behavioral
 benchmarks. Critically we notice that our best-performing model (#991) has a *positive*
 ρ -Hierarchy coefficient⁴ compared to the new state of the art model (#1033) and other
 remaining entries, where this coefficient is negative. This was an unexpected result that
 we found as most biologically-driven models obtain higher Brain-Scores at the initial stages
 of the visual hierarchy (V1) (Dapello et al., 2020), and these scores decrease as a function
 of hierarchy with generally worse Brain-Scores in the final stages (*e.g.* IT).

129 We also investigated the differential effects of rotation invariance and adversarial training used on
 130 top of a pretrained CrossViT-18 \dagger as shown in Table 2. We observed that each step independently
 131 helps to improve the overall Brain-Score, quite ironically at the expense of ImageNet Validation
 132 accuracy (Zhang et al., 2019). Interestingly, when both methods are combined (Adversarial training
 133 and rotation invariance), the model outperforms the baseline behavioral score by a large margin
 134 (+0.062), the IT score by (+0.047), the V4 score by (+0.036), the V2 score by (+0.068), and the V1
 135 score by (+0.020). Finally, our best model also retains a great standard accuracy at ImageNet from its
 136 pretrained version albeit a 10% drop, yet the performance on ImageNet Validation Accuracy of our
 137 model (73.53%) is still greater than a more biologically principled model such as the adversarially
 138 trained VOneResNet-50 (71.7%) (Dapello et al., 2020).

139 3 Assessment of CrossViT-18 \dagger -based models

140 As we have seen that the *optimization* procedure heavily influences the brain-score of each CrossViT-
 141 18 \dagger model, and thus its alignment to human vision (at a coarse level accepting the premise of the
 142 Brain-Score competition). We will now explore how different variations of such CrossViT’s change as
 143 a function of their training procedure, and thus their learned representations via a suite of experiments

⁴ ρ -Hierarchy coefficient: We define this as the Spearman rank correlation between the Brain-Scores of areas [V1,V2,V4,IT] with hierarchy: [1,2,3,4]

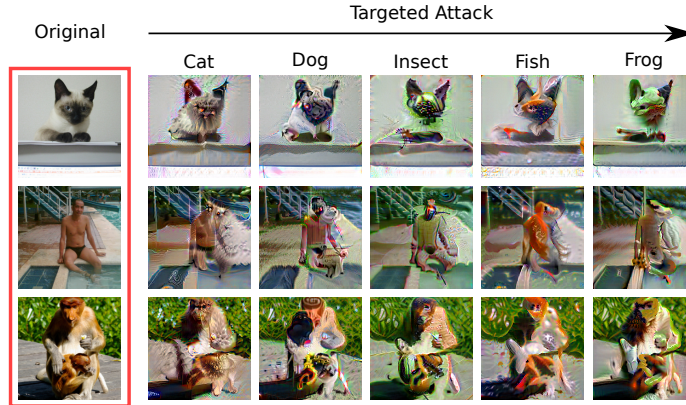


Figure 3: An extended demonstration of our winning model (CrossViT-18† [Adv. Training + Rot. invariance]) where a targeted attack is done for 3 images and the resulting stimuli is perceptually aligned with a human judgment of the fooled class. To our knowledge, this is the first time perceptually-aligned adversarial attacks have been shown to emerge in Transformer-based models.

144 that are more classical in computer vision. Additional experiments with CrossViT-18†-based models
 145 can be seen at Appendix B.

146 One of our most interesting qualitative results is
 147 that the *direction* of the adversarial attack made
 148 on our highest performing model resembles a
 149 distortion class that seems to fool a human ob-
 150 server too (Figures 4, 3). Alas, while the ad-
 151 versarial attack can be conceived as a type of
 152 *eigendistortion* as in Berardino et al. (2017) we
 153 *find* that the Brain-Score optimized Transformer
 154 models are more perceptually aligned to human
 155 observers when judging distorted stimuli. Simi-
 156 lar results were previously found by Santurkar
 157 et al. (2019) with ResNets, though there has not
 158 been any rigorous & unlimited time verification
 159 of this phenomena in humans similar to the work
 160 of Elsayed et al. (2018). Experimental details
 161 can be found in Appendix C

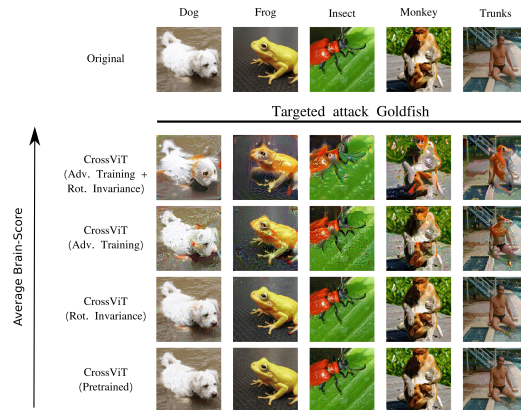


Figure 4: A qualitative demonstration of the human-machine perceptual alignment of the CrossViT-18† via the effects of adversarial perturbations. As the average Brain-Score increases in our system, the distortions seem to fool a human as well.

162 We also applied PGD attacks on our
 163 winning entry model (Adversarial Training
 164 + Rot. Invariance) on range $\epsilon \in$
 165 $\{1/255, 2/255, 4/255, 6/255, 8/255, 10/255\}$
 166 and step-size = $\frac{2.5}{\#PGD_{iterations}}$ as in the
 167 robustness Python library (Engstrom et al.,
 168 2019a) , in addition to three other controls:
 169 Adv. Training, Rotational Invariance, and a pretrained CrossViT, to evaluate how their adversarial
 170 robustness would change as a function of this particular distortion class. When doing this evaluation
 171 we observe in Figure 5 that Adversarially trained models are more robust to PGD attacks (three-step
 172 size flavors: 1 (FGSM), 10 & 20). One may be tempted to say that this is “expected” as the
 173 adversarially trained networks would be more robust, but the type of adversarial attack on which
 174 they are trained is different (FGSM as part of FAT (Wong et al., 2020) during training; and PGD
 175 at testing). Even if FGSM can be interpreted as a 1 step PGD attack, it is not obvious that this
 176 type of generalization would occur. In fact, it is of particular interest that the Adversarially trained
 177 CrossViT-18† with “fast adversarial training” (FAT) shows greater robustness to PGD 1 step attacks
 178 when the epsilon value used at testing time is very close to the values used at training (See Figure 5a).
 179 Naturally, for PGD-based attacks where the step size is greater (10 and 20; Figs. 5b,5c), our winning
 180 entry model achieves greater robustness against all other trained CrossViT’s independent of the ϵ
 181 values.

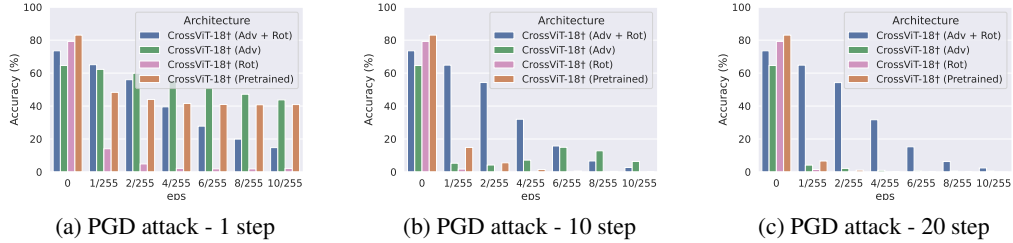


Figure 5: A suite of multiple steps [1,10,20] PGD-based adversarial attacks on clones of CrossViT-18† models that were optimized differently. Here we see that our winning entry (Adversarial training + Rotation Invariance) shows greater robustness (adversarial accuracy) than all other models as the number of steps of PGD-based attacks increases only for big step sizes of 10 & 20.

182 3.1 Feature Inversion

183 The last assessment we provided was inspired by feature inversion models that are a window to the
 184 representational soul of each model (Mahendran & Vedaldi, 2015). Oftentimes, models that are
 185 aligned with human visual perception in terms of their inductive biases and priors will show renderings
 186 that are very similar to the original image even when initialized from a noise image (Feather et al.,
 187 2019). We use the list of stimuli from Harrington & Deza (2022) to compare how several of these
 188 stimuli look like when they are rendered from the penultimate layer of a pretrained and our winning
 189 entry CrossViT-based model. A collection of synthesized images can be seen in Figure 6.

190 Even when these images are rendered starting from different noise images, Transformer-based models
 191 are remarkably good at recovering the structure of these images. This hints at a coherence with the
 192 results of Tuli et al. (2021) who have argued that Transformer-based models have a stronger shape
 193 bias than most CNN’s (Geirhos et al., 2019). We think this is due to their initial patch-embedding
 194 stage that preserves the visual organization of the image, though further investigation is necessary to
 195 validate this conjecture.

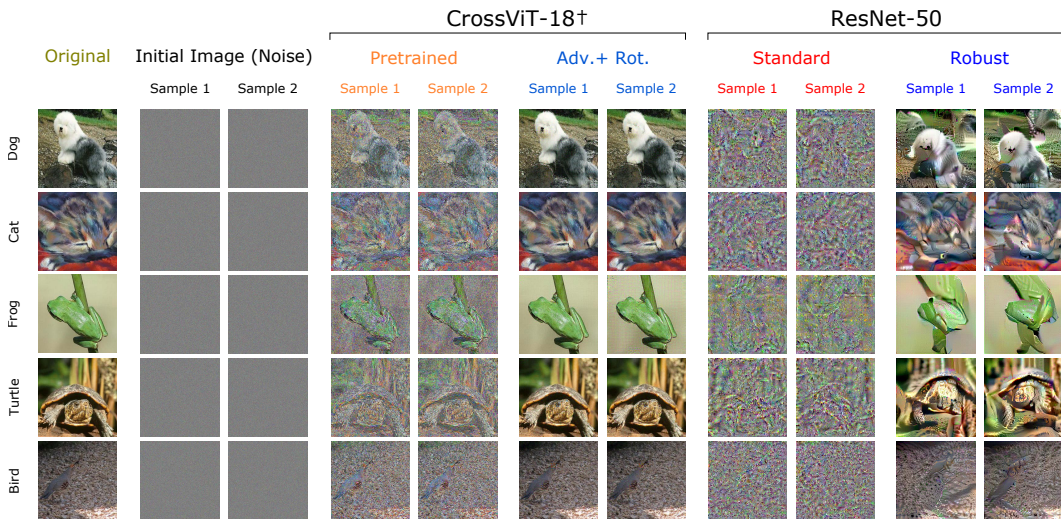


Figure 6: A summary of Feature Inversion models when applied on two different randomly samples noise images from a subset of the stimuli used in Harrington & Deza (2022). Standard and Pretrained models poorly invert the original stimuli leaving high spatial frequency artifacts.

196 **4 Discussion**

197 A question from this work that motivated the writing of this paper beyond the achievement of a high
198 score in the Brain-Score competition is: How does a CrossViT-18† perform so well at explaining
199 variance in primate area V4 without many iterations of hyper-parameter engineering? In this paper,
200 we have only scratched the surface of this question, but some clues have emerged.

201 One possibility is that the cross-attention mechanism of the CrossViT-18† is a proxy for Gramian-like
202 operations that encode local texture computation (vs global *a la* Geirhos et al. (2019)) which have
203 been shown to be pivotal for object representation in humans (Long et al., 2018; Jagadeesh & Gardner,
204 2022; Harrington & Deza, 2022). This initial conjecture is corroborated by our image inversion
205 effects (Section 3.1) where we find that CrossViT’s preserves the structure stronger than Residual
206 Networks (ResNets), while vanilla ViT’s shows strong grid-like artifacts.

207 Equally relevant throughout this paper has been the critical finding of the role of the optimization
208 procedure and the influence it has on achieving high Brain-Scores – even for non-biologically plausible
209 architectures (Riedel, 2022). Indeed, the simple combination of adding rotation invariance as an
210 implicit inductive bias through data-augmentation, and adding “worst-case scenario” (adversarial)
211 images in the training regime seems to create a perceptually-aligned representation for neural
212 networks (Santurkar et al., 2019).

213 On the other hand, the contributions to visual neuroscience from this paper are non-obvious. Tra-
214 ditionally, work in vision science has started from investigating phenomena in biological systems
215 via psychophysical experiments and/or neural recordings of highly controlled stimuli in animals, to
216 later verify their use or emergence when engineered in artificial perceptual systems. We are now in
217 a situation where we have “by accident” stumbled upon a perceptual system that can successfully
218 model (with half the full explained variance) visual processing in human area V4 – a region of which
219 its functional goal still remains a mystery to neuroscientists (Vacher et al., 2020; Bashivan et al.,
220 2019) –, giving us the chance to reverse engineer and dissect the contributions of the optimization
221 procedure to a fixed architecture. We have done our best to pin-point a causal root to this phenomena,
222 but we can only make an educated guess that a system with a cross-attention mechanism can *even*
223 *under regular training* achieve high V4 Brain-Scores, and these are maximized when optimized with
224 our joint adversarial training and rotation invariance procedure.

225 Ultimately, does this mean that Vision Trans-
226 formers are good models of the Human Ventral
227 Stream? We think that an answer to this ques-
228 tion is a response to the nursery rhyme: “*It looks*
229 *like a duck, and walks like a duck, but it’s not*
230 *a duck!*” One may be tempted to affirm that it
231 is a duck if we are only to examine the family
232 of in-distribution images from ImageNet at in-
233 ference; but when out of distribution stimuli are
234 shown to both machine and human perceptual
235 systems we will have a chance to accurately as-
236 sess their degree of perceptual similarity⁵. We
237 can tentatively expand this argument further by
238 studying adversarial images for both perceptual
239 systems (See also Figure 7). Future images used
240 in the Brain-Score competition that will better
241 assess human-machine representational similar-
242 ity should use these adversarial-like images to
243 test if the family of mistakes that machines make
244 are similar in nature than to the ones made by hu-
245 mans (See For example Golan et al. (2020)). If
246 that is to be the case, then we are one step closer
247 to building machines that can *see* like humans.

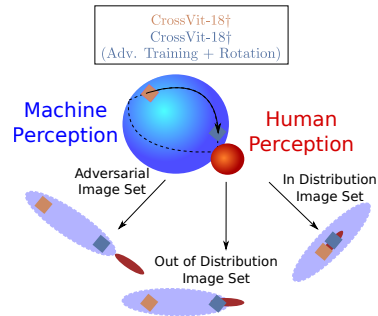


Figure 7: A cartoon inspired by Feather et al. (2019, 2021) depicting how our model changes its perceptual similarity depending on its optimization procedure. The arrows outside the spheres represent projections of such perceptual spaces that are observable by the images we show each system. While it may look like our winning model is “nearly human” it has still a long way to go, as the adversarial conditions have never been physiologically tested.

⁵Consider for example, that some stimuli used in Brain-Score are a basis set of Gabor filters, which are never encountered in nature