
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 Optimizing a CrossViT for the Brain-Score Competition

20 In this short paper, we try to solve an interesting question that was one of the motivations of this work:
21 “*Are Vision Transformers good models of the human ventral stream?*” Our approach to answering this
22 question will rely on using the Brain-Score platform (Schrimpf et al., 2020a) and participating in
23 their first yearly competition with a Transformer-based model. This platform quantifies the similarity
24 via bounded [0,1] scores of responses between a computer model and a set of non-human primates.
25 Here the ground truth is collected via neurophysiological recordings and/or behavioral outputs when
26 primates are performing psychophysical tasks, and the scores are computed by some derivation of
27 Representational Similarity Analysis (Kriegeskorte et al., 2008) when pitted against artificial neural
28 network activations of modern computer vision models.

29 We discuss an interesting finding, where amidst the constant debate of the biological plausibility of
30 Vision Transformers – which have been deemed less biologically plausible than convolutional neural
31 networks (as discussed in: URL_1 URL_2, though also see Conwell et al. (2021)) –, we find that
32 when these Transformers are optimized under certain conditions, they may achieve high explainable
33 variance with regards to many areas in primate vision, and surprisingly the highest score to date at
34 the time of the competition for explainable variance in area V4, that still remains a mystery in visual

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).

35 neuroscience (see Pasupathy et al. (2020) for a review). Our final model and highest scoring model
36 was based on several insights:

37 **Adversarial-Training:** Work by Santurkar et al. (2019); Engstrom et al. (2019b); Dapello et al.
38 (2020), has shown that convolutional neural networks trained adversarially¹ yield human perceptually-
39 aligned distortions when attacked. This is an interesting finding, that perhaps extends to vision
40 transformers, but has never been qualitatively tested before though recent works – including this
41 one (See Figure 2) – have started to investigate in this direction (Tuli et al., 2021; Caro et al., 2020).
42 Thus we projected that once we picked a specific vision transformer architecture, we would train it
43 adversarially.

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

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

53 **Localized texture-based computation:** Despite the emergence of a *global* texture-bias in object
54 recognition when training Deep Neural Networks (Geirhos et al., 2019) – object recognition is a
55 compositional process (Brendel & Bethge, 2019; Deza et al., 2020). Recently, works in neuroscience
56 have also suggested that *local* texture computation is perhaps pivotal for object recognition to either
57 create an ideal basis set from which to represent objects (Long et al., 2018; Jagadeesh & Gardner,
58 2022) and/or encode robust representations (Harrington & Deza, 2022).

59 After searching for several models in the com-
60 puter vision literature that resemble a Transformer
61 model that ticks all the boxes above, we opted for a
62 CrossViT-18† (that includes multi-resolution + local
63 texture-based computation) that was trained with
64 rotation-based augmentations and also adversarial
65 training (See Appendix A.3 for exact training de-
66 tails, our *best* model also used $p = 0.25$ grayscale
67 augmentation, though this contribution to model
68 Brain-Score is minimal).

Table 2: Selected Layers of CrossViT-18†

Benchmark	Layer
V1,V2,V4	blocks.1.blocks.1.0.norm1
IT	blocks.1.blocks.1.4.norm2
Behavior	blocks.2.revert_projs.1.2

69 **Results:** Our best performing model #991 achieved
70 2nd place in the overall Brain-Score 2022 competition (Schrimpf et al., 2020b) as shown in Table 1.
71 At the time of submission, it holds the first place for the highest explainable variance of area V4
72 and the second highest score in the IT area. Our model also currently ranks 6th across all Brain-
73 Score submitted models as shown on the main brain-score website (including those outside the
74 competition and since the start of the platform’s conception, totaling 216). A general schematic of
75 how Brain-Scores are calculated can be seen in Figure 1.

¹Adversarial 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 3: 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.

76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91

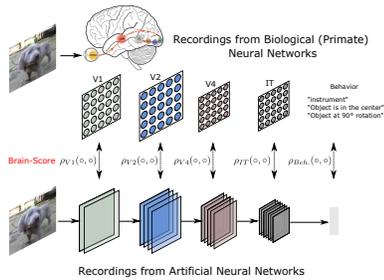


Figure 1: A schematic of how brain-score is calculated as similarity metrics obtained from neural responses and model activations.

92
93
94
95
96
97
98
99
100
101

We also investigated the differential effects of rotation invariance and adversarial training used on top of a pretrained CrossViT-18 \dagger as shown in Table 3. We observed that each step independently helps to improve the overall Brain-Score, quite ironically at the expense of ImageNet Validation accuracy (Zhang et al., 2019). Interestingly, when both methods are combined (Adversarial training and rotation invariance), the model outperforms the baseline behavioral score by a large margin (+0.062), the IT score by (+0.047), the V4 score by (+0.036), the V2 score by (+0.068), and the V1 score by (+0.020). Finally, our best model also retains a great standard accuracy at ImageNet from its pretrained version albeit a 10% drop, yet the performance on ImageNet Validation Accuracy of our model (73.53%) is still greater than a more biologically principled model such as the adversarially trained VOneResNet-50 (71.7%) (Dapello et al., 2020).

102

2 Assessment of CrossViT-18 \dagger -based models

103
104
105
106
107
108

As we have seen that the *optimization* procedure heavily influences the brain-score of each CrossViT-18 \dagger model, and thus its alignment to human vision (at a coarse level accepting the premise of the Brain-Score competition). We will now explore how different variations of such CrossViT's change as a function of their training procedure, and thus their learned representations via a suite of experiments that are more classical in computer vision. Additional experiments with CrossViT-18 \dagger -based models can be seen at Appendix B.

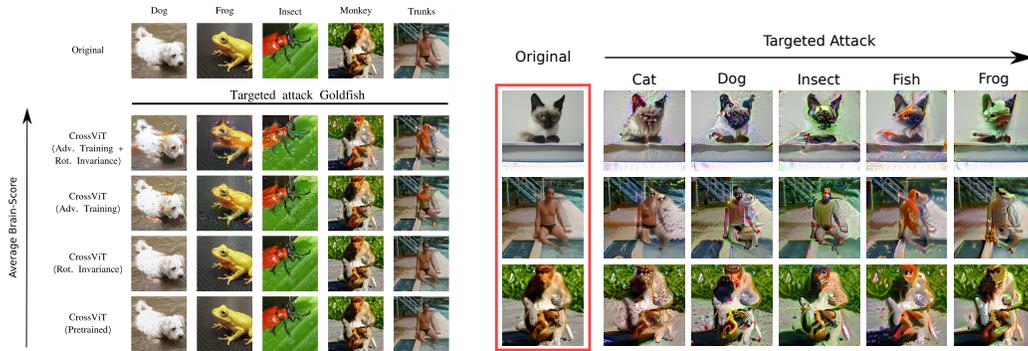
109

2.1 Adversarial Attacks

110
111
112
113

One of our most interesting qualitative results is that the *direction* of the adversarial attack made on our highest performing model resembles a distortion class that seems to fool a human observer too (Figure 2). Alas, while the adversarial attack can be conceived as a type of *eigendistortion* as in Berardino et al. (2017) we *find* that the Brain-Score optimized Transformer models are more

² 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]



(a) 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. (b) 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.

Figure 2: Exploring Human-Machine Perceptual Alignment via Adversarial Attacks.

114 perceptually aligned to human observers when judging distorted stimuli. Similar results were
 115 previously found by Santurkar et al. (2019) with ResNets, though there has not been any rigorous &
 116 unlimited time verification of this phenomena in humans similar to the work of Elsayed et al. (2018).

117 2.2 Feature Inversion

118 The last assessment we provided was inspired by feature inversion models that are a window to the
 119 representational soul of each model (Mahendran & Vedaldi, 2015). Oftentimes, models that are
 120 aligned with human visual perception in terms of their inductive biases and priors will show renderings
 121 that are very similar to the original image even when initialized from a noise image (Feather et al.,
 122 2019). We use the list of stimuli from Harrington & Deza (2022) to compare how several of these
 123 stimuli look like when they are rendered from the penultimate layer of a pretrained and our winning
 124 entry CrossViT-based model. A collection of synthesized images can be seen in Figure 3.

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

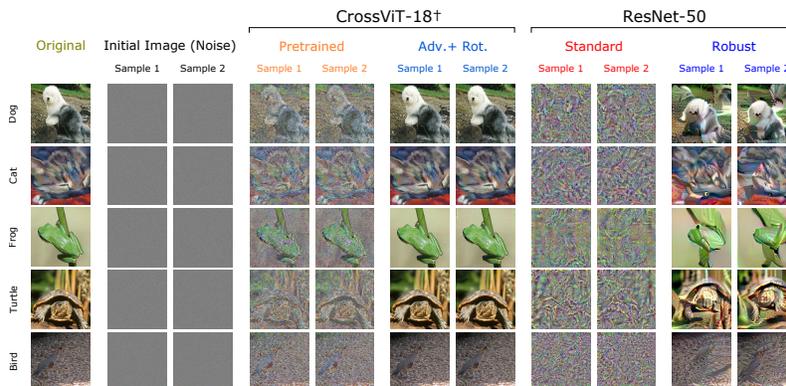


Figure 3: A summary of Feature Inversion models when applied on two different randomly samples noise images from a subset of the stimuli used in Harrington & Deza (2022). Standard and Pretrained models poorly invert the original stimuli leaving high spatial frequency artifacts. Adversarial training improves image inversion models, and this is even more evident for Transformer models.

131 **3 Discussion**

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

136 One possibility is that the cross-attention mechanism of the CrossViT-18† is a proxy for Gramian-like
 137 operations that encode local texture computation (vs global *a la* Geirhos et al. (2019)) which have
 138 been shown to be pivotal for object representation in humans (Long et al., 2018; Jagadeesh & Gardner,
 139 2022; Harrington & Deza, 2022). This initial conjecture is corroborated by our image inversion
 140 effects (Section 2.2) where we find that CrossViT’s preserves the structure stronger than Residual
 141 Networks (ResNets), while vanilla ViT’s shows strong grid-like artifacts (See Figures 12, 13 in the
 142 supplementary material).

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

149 On the other hand, the contributions to visual neuroscience from this paper are non-obvious. Tra-
 150 ditionally, work in vision science has started from investigating phenomena in biological systems
 151 via psychophysical experiments and/or neural recordings of highly controlled stimuli in animals, to
 152 later verify their use or emergence when engineered in artificial perceptual systems. We are now in
 153 a situation where we have “by accident” stumbled upon a perceptual system that can successfully
 154 model (with half the full explained variance) visual processing in human area V4 – a region of which
 155 its functional goal still remains a mystery to neuroscientists (Vacher et al., 2020; Bashivan et al.,
 156 2019) –, giving us the chance to reverse engineer and dissect the contributions of the optimization
 157 procedure to a fixed architecture. We have done our best to pin-point a causal root to this phenomena,
 158 but we can only make an educated guess that a system with a cross-attention mechanism can *even*
 159 *under regular training* achieve high V4 Brain-Scores, and these are maximized when optimized with
 160 our joint adversarial training and rotation invariance procedure.

161 Ultimately, does this mean that Vision Trans-
 162 formers are good models of the Human Ventral
 163 Stream? We think that an answer to this ques-
 164 tion is a response to the nursery rhyme: “*It looks*
 165 *like a duck, and walks like a duck, but it’s not*
 166 *a duck!*” One may be tempted to affirm that it
 167 is a duck if we are only to examine the family
 168 of in-distribution images from ImageNet at infer-
 169 ence; but when out of distribution stimuli are
 170 shown to both machine and human perceptual
 171 systems we will have a chance to accurately as-
 172 sess their degree of perceptual similarity³. We
 173 can tentatively expand this argument further by
 174 studying adversarial images for both perceptual
 175 systems (See also Figure 4). Future images used
 176 in the Brain-Score competition that will better
 177 assess human-machine representational similar-
 178 ity should use these adversarial-like images to
 179 test if the family of mistakes that machines make
 180 are similar in nature than to the ones made by hu-
 181 mans (See For example Golan et al. (2020)). If
 182 that is to be the case, then we are one step closer
 183 to building machines that can *see* like humans.

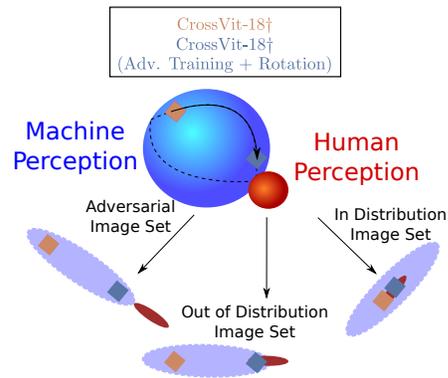


Figure 4: A cartoon inspired by Feather et al. (2019, 2021) depicting how our model changes its perceptual similarity depending on its optimization procedure. The arrows outside the spheres represent projections of such perceptual spaces that are observable by the images we show each system. While it 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.

³Consider for example, that some stimuli used in Brain-Score are a basis set of Gabor filters, which are never encountered in nature

184 References

- 185 Simone Azeglio, Simone Poetto, Marco Nurisso, and Luca Savant Aira. Improving neural predictivity
186 in the visual cortex with gated recurrent connections. In *Brain-Score Workshop*, 2022. URL
187 <https://openreview.net/forum?id=SKWlBbDWXWc>.
- 188 Pouya Bashivan, Kohitij Kar, and James J DiCarlo. Neural population control via deep image
189 synthesis. *Science*, 364(6439):eaav9436, 2019.
- 190 Alexander Bernardino, Valero Laparra, Johannes Ballé, and Eero Simoncelli. Eigen-distortions of
191 hierarchical representations. *Advances in neural information processing systems*, 30, 2017.
- 192 Wieland Brendel and Matthias Bethge. Approximating cnns with bag-of-local-features models works
193 surprisingly well on imagenet. *arXiv preprint arXiv:1904.00760*, 2019.
- 194 Peter J Burt and Edward H Adelson. The laplacian pyramid as a compact image code. In *Readings in*
195 *computer vision*, pp. 671–679. Elsevier, 1987.
- 196 Josue Ortega Caro, Yilong Ju, Ryan Pyle, Sourav Dey, Wieland Brendel, Fabio Anselmi, and Ankit
197 Patel. Local convolutions cause an implicit bias towards high frequency adversarial examples.
198 *arXiv preprint arXiv:2006.11440*, 2020.
- 199 Chun-Fu Richard Chen, Quanfu Fan, and Rameswar Panda. Crossvit: Cross-attention multi-scale
200 vision transformer for image classification. In *Proceedings of the IEEE/CVF International Confer-*
201 *ence on Computer Vision*, pp. 357–366, 2021.
- 202 Colin Conwell, Jacob S. Prince, George A. Alvarez, and Talia Konkle. What can 5.17 billion
203 regression fits tell us about artificial models of the human visual system? In *SVRHM 2021*
204 *Workshop @ NeurIPS*, 2021. URL https://openreview.net/forum?id=i_xiyGq6FNT.
- 205 Joel Dapello, Tiago Marques, Martin Schrimpf, Franziska Geiger, David Cox, and James J DiCarlo.
206 Simulating a primary visual cortex at the front of cnns improves robustness to image perturbations.
207 *Advances in Neural Information Processing Systems*, 33:13073–13087, 2020.
- 208 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale
209 hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,
210 pp. 248–255. Ieee, 2009.
- 211 Arturo Deza, Qianli Liao, Andrzej Banburski, and Tomaso Poggio. Hierarchically compositional
212 tasks and deep convolutional networks. *arXiv preprint arXiv:2006.13915*, 2020.
- 213 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
214 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit,
215 and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale.
216 In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=YicbFdNTTy>.
- 217
- 218 Alexander S. Ecker, Fabian H. Sinz, Emmanouil Froudarakis, Paul G. Fahey, Santiago A. Ca-
219 dena, Edgar Y. Walker, Erick Cobos, Jacob Reimer, Andreas S. Tolia, and Matthias Bethge. A
220 rotation-equivariant convolutional neural network model of primary visual cortex. In *International*
221 *Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=H1fU8iAqKX>.
- 222
- 223 Gamaleldin Elsayed, Shreya Shankar, Brian Cheung, Nicolas Papernot, Alexey Kurakin, Ian Good-
224 fellow, and Jascha Sohl-Dickstein. Adversarial examples that fool both computer vision and
225 time-limited humans. *Advances in neural information processing systems*, 31, 2018.
- 226 Logan Engstrom, Andrew Ilyas, Hadi Salman, Shibani Santurkar, and Dimitris Tsipras. Robustness
227 (python library), 2019a. URL <https://github.com/MadryLab/robustness>.
- 228 Logan Engstrom, Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Brandon Tran, and Alek-
229 sander Madry. Adversarial robustness as a prior for learned representations. *arXiv preprint*
230 *arXiv:1906.00945*, 2019b.

- 231 Jenelle Feather, Alex Durango, Ray Gonzalez, and Josh McDermott. Metamers of neural networks
232 reveal divergence from human perceptual systems. *Advances in Neural Information Processing*
233 *Systems*, 32, 2019.
- 234 Jenelle Feather, Alex Durango, Guillaume Leclerc, Aleksander Madry, and Josh McDermott. Adver-
235 sarial training aligns invariances between artificial neural networks and biological sensory systems.
236 *Cosyne Meeting Abstract*, 2021.
- 237 Jeremy Freeman and Eero P Simoncelli. Metamers of the ventral stream. *Nature neuroscience*, 14(9):
238 1195–1201, 2011.
- 239 Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A. Wichmann, and
240 Wieland Brendel. Imagenet-trained CNNs are biased towards texture; increasing shape bias
241 improves accuracy and robustness. In *International Conference on Learning Representations*,
242 2019. URL <https://openreview.net/forum?id=Bygh9j09KX>.
- 243 Tal Golan, Prashant C. Raju, and Nikolaus Kriegeskorte. Controversial stimuli: Pitting neural
244 networks against each other as models of human cognition. *Proceedings of the National Academy*
245 *of Sciences*, 117(47):29330–29337, 2020. doi: 10.1073/pnas.1912334117. URL <https://www.pnas.org/doi/abs/10.1073/pnas.1912334117>.
- 247 Anne Harrington and Arturo Deza. Finding biological plausibility for adversarially robust features
248 via metameric tasks. In *International Conference on Learning Representations*, 2022. URL
249 https://openreview.net/forum?id=yeP_zx9vqNm.
- 250 David J Heeger and James R Bergen. Pyramid-based texture analysis/synthesis. In *Proceedings of*
251 *the 22nd annual conference on Computer graphics and interactive techniques*, pp. 229–238, 1995.
- 252 Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common
253 corruptions and perturbations. In *International Conference on Learning Representations*, 2019.
254 URL <https://openreview.net/forum?id=HJz6tiCqYm>.
- 255 Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul
256 Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, Dawn Song, Jacob Steinhardt, and Justin Gilmer.
257 The many faces of robustness: A critical analysis of out-of-distribution generalization. *ICCV*,
258 2021.
- 259 Akshay V. Jagadeesh and Justin L. Gardner. Texture-like representation of objects in human visual cor-
260 tex. *Proceedings of the National Academy of Sciences*, 119(17):e2115302119, 2022. doi: 10.1073/
261 [pnas.2115302119](https://www.pnas.org/doi/abs/10.1073/pnas.2115302119). URL <https://www.pnas.org/doi/abs/10.1073/pnas.2115302119>.
- 262 Ahmadreza Jeddi, Mohammad Javad Shafiee, and Alexander Wong. A simple fine-tuning is all you
263 need: Towards robust deep learning via adversarial fine-tuning, 2020.
- 264 Nancy Kanwisher, Frank Tong, and Ken Nakayama. The effect of face inversion on the human
265 fusiform face area. *Cognition*, 68(1):B1–B11, 1998.
- 266 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
267 *arXiv:1412.6980*, 2014.
- 268 Klim Kireev, Maksym Andriushchenko, and Nicolas Flammarion. On the effectiveness of adversarial
269 training against common corruptions. *arXiv preprint arXiv:2103.02325*, 2021.
- 270 Nikolaus Kriegeskorte, Marieke Mur, and Peter A Bandettini. Representational similarity analysis-
271 connecting the branches of systems neuroscience. *Frontiers in systems neuroscience*, 2:4, 2008.
- 272 Alfred Laugros, Alice Caplier, and Matthieu Ospici. Are adversarial robustness and common
273 perturbation robustness independant attributes? In *Proceedings of the IEEE/CVF International*
274 *Conference on Computer Vision Workshops*, pp. 0–0, 2019.
- 275 Zhe Li, Josue Ortega Caro, Evgenia Rusak, Wieland Brendel, Matthias Bethge, Fabio Anselmi,
276 Ankit B Patel, Andreas S Tolias, and Xaq Pitkow. Robust deep learning object recognition models
277 rely on low frequency information in natural images. *bioRxiv*, 2022.

- 278 Bria Long, Chen-Ping Yu, and Talia Konkle. Mid-level visual features underlie the high-level
279 categorical organization of the ventral stream. *Proceedings of the National Academy of Sciences*,
280 115(38):E9015–E9024, 2018.
- 281 Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.
282 Towards deep learning models resistant to adversarial attacks. In *International Conference on*
283 *Learning Representations*, 2018. URL <https://openreview.net/forum?id=rJzIBfZAb>.
- 284 Aravindh Mahendran and Andrea Vedaldi. Understanding deep image representations by inverting
285 them. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.
286 5188–5196, 2015.
- 287 Tiago Marques, Martin Schrimpf, and James J DiCarlo. Multi-scale hierarchical neural network
288 models that bridge from single neurons in the primate primary visual cortex to object recognition
289 behavior. *bioRxiv*, 2021.
- 290 Anitha Pasupathy, Dina V Popovkina, and Taekjun Kim. Visual functions of primate area v4. *Annual*
291 *review of vision science*, 6:363–385, 2020.
- 292 Rishi Rajalingham, Elias B. Issa, Pouya Bashivan, Kohitij Kar, Kailyn Schmidt, and James J. DiCarlo.
293 Large-scale, high-resolution comparison of the core visual object recognition behavior of humans,
294 monkeys, and state-of-the-art deep artificial neural networks. *bioRxiv*, 2018. doi: 10.1101/240614.
295 URL <https://www.biorxiv.org/content/early/2018/02/12/240614>.
- 296 Alexander Riedel. Bag of tricks for training brain-like deep neural networks. In *Brain-Score*
297 *Workshop*, 2022. URL <https://openreview.net/forum?id=SudzH-vWQ-c>.
- 298 Ruth Rosenholtz, Jie Huang, Alvin Raj, Benjamin J Balas, and Livia Ilie. A summary statistic
299 representation in peripheral vision explains visual search. *Journal of vision*, 12(4):14–14, 2012.
- 300 Shibani Santurkar, Andrew Ilyas, Dimitris Tsipras, Logan Engstrom, Brandon Tran, and Aleksander
301 Madry. Image synthesis with a single (robust) classifier. *Advances in Neural Information Processing*
302 *Systems*, 32, 2019.
- 303 Martin Schrimpf, Jonas Kubilius, Ha Hong, Najib J Majaj, Rishi Rajalingham, Elias B Issa, Kohitij
304 Kar, Pouya Bashivan, Jonathan Prescott-Roy, Franziska Geiger, et al. Brain-score: Which artificial
305 neural network for object recognition is most brain-like? *BioRxiv*, pp. 407007, 2020a.
- 306 Martin Schrimpf, Jonas Kubilius, Michael J. Lee, N. Apurva Ratan Murty, Robert Ajemian, and
307 James J. DiCarlo. Integrative benchmarking to advance neurally mechanistic models of hu-
308 man intelligence. *Neuron*, 108(3):413–423, 2020b. ISSN 0896-6273. doi: [https://doi.org/](https://doi.org/10.1016/j.neuron.2020.07.040)
309 [10.1016/j.neuron.2020.07.040](https://doi.org/10.1016/j.neuron.2020.07.040). URL [https://www.sciencedirect.com/science/article/](https://www.sciencedirect.com/science/article/pii/S089662732030605X)
310 [pii/S089662732030605X](https://www.sciencedirect.com/science/article/pii/S089662732030605X).
- 311 Eero P Simoncelli and William T Freeman. The steerable pyramid: A flexible architecture for multi-
312 scale derivative computation. In *Proceedings., International Conference on Image Processing*,
313 volume 3, pp. 444–447. IEEE, 1995.
- 314 Shikhar Tuli, Ishita Dasgupta, Erin Grant, and Thomas L Griffiths. Are convolutional neural networks
315 or transformers more like human vision? *arXiv preprint arXiv:2105.07197*, 2021.
- 316 Jonathan Vacher, Aida Davila, Adam Kohn, and Ruben Coen-Cagli. Texture interpolation for probing
317 visual perception. *Advances in Neural Information Processing Systems*, 33:22146–22157, 2020.
- 318 Eric Wong, Leslie Rice, and J. Zico Kolter. Fast is better than free: Revisiting adversarial training,
319 2020.
- 320 Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric Xing, Laurent El Ghaoui, and Michael Jordan.
321 Theoretically principled trade-off between robustness and accuracy. In *International conference on*
322 *machine learning*, pp. 7472–7482. PMLR, 2019.

323 **Checklist**

- 324 1. For all authors...
- 325 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
326 contributions and scope? [Yes]
- 327 (b) Did you describe the limitations of your work? [Yes] These are mentioned in the
328 Discussion
- 329 (c) Did you discuss any potential negative societal impacts of your work? [N/A] We don't
330 anticipate negative social impacts from this work.
- 331 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
332 them? [Yes] Yes, we have read the guidelines and ensured our paper conforms to them.
- 333 2. If you are including theoretical results...
- 334 (a) Did you state the full set of assumptions of all theoretical results? [N/A] We do not
335 have theoretical results.
- 336 (b) Did you include complete proofs of all theoretical results? [N/A] We do not have
337 theoretical results.
- 338 3. If you ran experiments...
- 339 (a) Did you include the code, data, and instructions needed to reproduce the main ex-
340 perimental results (either in the supplemental material or as a URL)? [Yes] We have
341 included a hyperlink from the publicly available Brain-Score 2022 competition. If ac-
342 cepted we will provide de-anonymized links to our entry model and all other competing
343 models from Table 1
- 344 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
345 were chosen)? [Yes] These were all specified in the Appendix and main body when
346 relevant
- 347 (c) Did you report error bars (e.g., with respect to the random seed after running exper-
348 iments multiple times)? [No] Our experiments did not include error bars as running
349 them were expensive, and internal pilot trials showed that variation was minimal (pilot
350 models converged to nearly identical behavior when run twice). However, models were
351 initialized uniformly from the same PreTrained Model/seed when applicable to analyze
352 the contribution of each training/fine-tuning regime.
- 353 (d) Did you include the total amount of compute and the type of resources used (e.g., type
354 of GPUs, internal cluster, or cloud provider)? [Yes] Please See Appendix A.3
- 355 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 356 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 357 (b) Did you mention the license of the assets? [Yes] Assets were all OpenSource
- 358 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
359 These are provided through-out the paper either cited in the main body, as footnotes, or
360 highlighted in the Appendix.
- 361 (d) Did you discuss whether and how consent was obtained from people whose data you're
362 using/curating? [N/A] We did not use any human-based data
- 363 (e) Did you discuss whether the data you are using/curating contains personally identifiable
364 information or offensive content? [N/A] We did not use human data.
- 365 5. If you used crowdsourcing or conducted research with human subjects...
- 366 (a) Did you include the full text of instructions given to participants and screenshots, if
367 applicable? [N/A] We did not collect human data.
- 368 (b) Did you describe any potential participant risks, with links to Institutional Review
369 Board (IRB) approvals, if applicable? [N/A] We did not collect human data.
- 370 (c) Did you include the estimated hourly wage paid to participants and the total amount
371 spent on participant compensation? [N/A] We did not collect human data.

372 A Experimental Setup

373 A.1 Dataset

374 We used the ImageNet 1k (Deng et al., 2009) dataset for training. ImageNet1K contains 1,000 classes
375 and the number of training and validation images are 1.28 million and 50,000, respectively. We
376 validate the effectiveness of our models in the different datasets proposed in the Brain-Score (Schrimpf
377 et al., 2020a) competition.

378 A.2 Custom Scheduler

379 The proposed learning rate scheduler is based
380 on Jeddi et al. (2020) and is formulated as
381 $LR = 0.00012 \times e - 0.0004$ for $e = 1$ and
382 $LR = \frac{0.00002}{2^{e-2}}$ for $1 < e \leq 6$. As shown
383 in Figure 5, we start with a small learning rate
384 and then it is smoothly increased for one epoch.
385 We empirically found that fine-tuning the trans-
386 former for more than 1 epoch resulted in an
387 under-fitting behavior of the adversarial robust-
388 ness. After this first epoch, the learning rate
389 is reduced very fast so that model performance
390 converges to a steady state, without having too
391 much time to overfit on the training data.

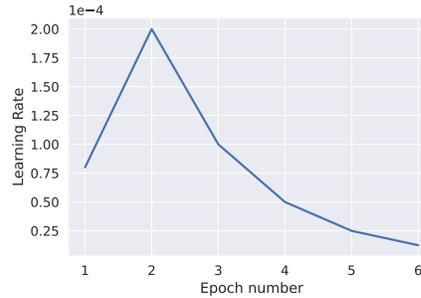


Figure 5: Custom scheduler used for training the Vision Transformer.

392 A.3 Training Setup

393 We used a pretrained CrossViT-18† (Chen et al.,
394 2021) downloaded from the timm library that
395 is adversarially trained via a fast gradient sign
396 method (FGSM) attack and random initializa-
397 tion (Wong et al., 2020). We opted for this strat-
398 egy, known as "Fast Adversarial Training" as it
399 allows a faster iteration in comparison with other
400 common approaches (e.g. adversarial training
401 with the PGD attack). In particular, all experi-
402 ments used $\epsilon = 2/255$ and step size $\alpha = 1.25\epsilon$
403 as proposed originally in (Wong et al., 2020).
404 However, in contrast to the previous method, we
405 follow a 5 epoch fine-tuning approach with a cus-
406 tom learning rate scheduler in order to avoid un-
407 derfitting. We optimize our networks with Adap-
408 tive Moment Estimation (Adam *a la* Kingma
409 & Ba (2014)) and employed mixed precision
410 for faster training. All input images were pre-
411 processed with resizing to 256×256 followed by
412 standard random cropping and horizontal mirroring.
413 In the case of our best performing model (#991), we additionally incorporated a random grayscale
414 transformation ($p = 0.25$) and a set of hard rotation transformations of $(0^\circ, 90^\circ, 180^\circ, 270^\circ)$ –
415 implicitly aiding for rotational invariance – due to the characteristics of images appearing in the
416 behavioral benchmark of Rajalingham et al. (2018). All our experiments were ran locally on a
GPU-Tesla V-100. Each adversarial training of a vision transformer took around 48 hours.

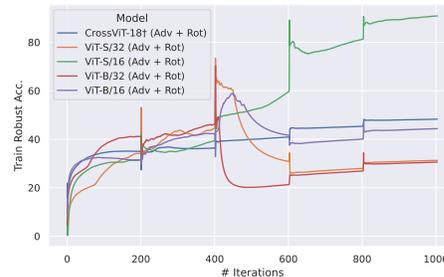


Figure 6: Training robust acc. of each Vision Transformer model (Adv + Rot). We clearly observed that ViT-S/16 has over-fitted during training.

417 Optionally include extra information (complete proofs, additional experiments and plots) in the
418 appendix. This section will often be part of the supplemental material.

419 B Additional Assessment of CrossViT-18†-based models

420 B.1 Robustness against adversarial attacks

421 We also applied PGD attacks on our winning entry model (Adversarial Training + Rot. Invariance) on
422 range $\epsilon \in \{1/255, 2/255, 4/255, 6/255, 8/255, 10/255\}$ and step-size = $\frac{2.5}{\#PGD_{iterations}}$ as in the
423 robustness Python library (Engstrom et al., 2019a), in addition to three other controls: Adv. Training,
424 Rotational Invariance, and a pretrained CrossViT, to evaluate how their adversarial robustness would
425 change as a function of this particular distortion class. When doing this evaluation we observe in
426 Figure 7 that Adversarially trained models are more robust to PGD attacks (three-step size flavors: 1
427 (FGSM), 10 & 20). One may be tempted to say that this is “expected” as the adversarially trained
428 networks would be more robust, but the type of adversarial attack on which they are trained is different
429 (FGSM as part of FAT (Wong et al., 2020) during training; and PGD at testing). Even if FGSM can
430 be interpreted as a 1 step PGD attack, it is not obvious that this type of generalization would occur.
431 In fact, it is of particular interest that the Adversarially trained CrossViT-18† with “fast adversarial
432 training” (FAT) shows greater robustness to PGD 1 step attacks when the epsilon value used at testing
433 time is very close to the values used at training (See Figure 7a). Naturally, for PGD-based attacks
434 where the step size is greater (10 and 20; Figs. 7b,7c), our winning entry model achieves greater
435 robustness against all other trained CrossViT’s independent of the ϵ values.

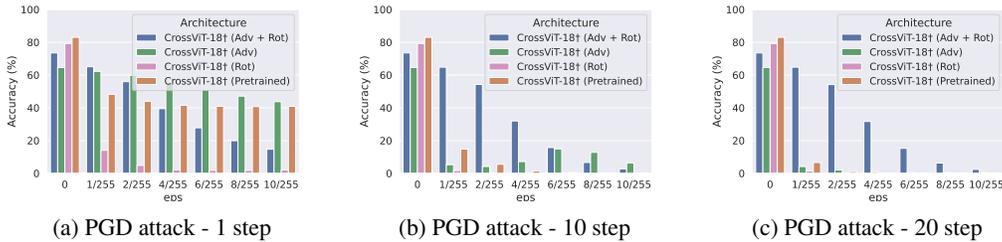


Figure 7: A suite of multiple steps [1,10,20] PGD-based adversarial attacks on clones of CrossViT-18† models that were optimized differently. Here we see that our winning entry (Adversarial training + Rotation Invariance) shows greater robustness (adversarial accuracy) than all other models as the number of steps of PGD-based attacks increases only for big step sizes of 10 & 20.

436 B.2 Mid-Ventral Stimuli Interpretation

437 In addition to the previous experiments, we wondered how well the two models: CrossViT-18†
438 (PreTrained) and CrossViT-18† (Adv. Training + Rot. Invariance) could linearly separate a small
439 subset of 2-class stimuli across their visual hierarchy. For this experiment, we used both the
440 original and texform stimuli (100 images per class) from Harrington & Deza (2022), where the
441 texform stimuli can be used to test the mechanisms of human peripheral computation (Rosenholtz
442 et al., 2012; Freeman & Simoncelli, 2011) or mid-ventral human computation (Long et al., 2018;
443 Jagadeesh & Gardner, 2022). Roughly speaking these texforms are very similar to their original
444 counter-part, where they match in global structure (*i.e.* form), but are locally distorted through a
445 texture-matching operation (*i.e.* texture) as seen in Figure 8 (Inset 0.). In this analysis, we will use a
446 t-SNE projection with a fixed random seed across both models and stimuli to evaluate the qualitative
447 similarity/differences of their 2D clustering patterns.

448 Here we are interested in exposing our models to this distortion class because recent work has used
449 these types of stimuli to show that human peripheral computation may act as a biological proxy for an
450 adversarially robust processing system (Harrington & Deza, 2022), and that humans may in-fact use
451 strong texture-like cues to perform object recognition (in IT) without the specific need for a strong
452 structural cue (Jagadeesh & Gardner, 2022).

453 We find that Pretrained CrossViT-18† models have trouble in early visual cortex read-out sections
454 to cluster both classes. In fact, several images are considered “visual outliers” for both original and
455 texform images. These differences are slowly resolved only for the original images as we go higher
456 in depth in the Transformer model until we get to the Behavior read-out layer. This is not the case for
457 the texforms, where the PreTrained CrossViT-18† can not tease apart the primate and insect classes

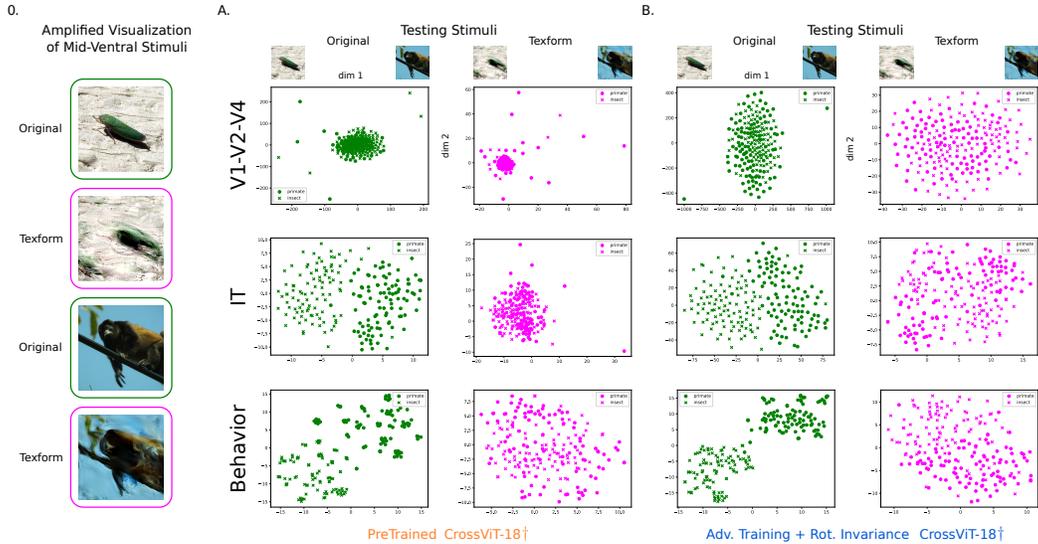


Figure 8: A comparison of how two CrossViT-18† models manage to classify original and texform stimuli. In (0.) we see a magnification of a texform, and in (A.,B.) we see how our winning Model Adv. + Rot. manages to create tighter vicinities across the visual stimuli, and ultimately – at the behavioral level – can separate both original and texform stimuli, while pretrained transformers seem to struggle with texform linear separability at the behavioral stage.

458 at such simulated behavioral stage. This story was to our surprise very different and more coherent
 459 with human visual processing for the Adv + Rot CrossViT-18† where outliers no longer exist – as
 460 there are none in the small dataset –, and the degree of linear separability for the original and texform
 461 stimuli increases to near perfect separation for both stimuli at the behavioral stage.

462 B.3 Common Corruption Benchmarks

463 We also looked into how adversarial training would affect the performance of the different sets of
 464 neural networks to common corruptions that are *not* adversarial. To do this, we ran our models and
 465 benchmarked them to the ImageNet-C dataset (Hendrycks & Dietterich, 2019).

466 One would have expected Brain-Aligned models like our adversarially-trained + rotationally invariant
 467 CrossViT to also present strong robustness to common corruptions. To our surprise, this was not
 468 the case as seen in Table 5. This is a puzzling result, though there have been several bodies of
 469 work suggesting that adversarial robustness and common corruptions robustness are independent
 470 phenomena (Laugros et al., 2019), however, Kireev et al. (2021) have proved otherwise contingent on
 471 the l_∞ radius⁴ – but now see Li et al. (2022).

Network	Clean Accuracy (†)	mce (↓)	Gauss	Shot	Impulse	Defocus	Glass	Motion	Zoom	Snow	Frost	Fog	Bright	Contrast	Elastic	Pixel	JPEG
ResNet50-Augmix	77.53	67.1	65.5	65.1	66.4	67.7	81	63.9	65.5	71.6	70.9	66.5	57.8	60.2	76.9	59.5	68.5
CrossViT-18† (Adv + Rot)	73.53	79.5	80.7	81.6	83.2	90.2	78.7	82.4	80	77.6	74	107.9	65	100.4	74.2	57.4	58.7
CrossViT-18† (Adv)	64.60	88.8	85	85.7	86.7	96.7	88	92.1	91.3	85.8	83.6	109.3	82.2	104.9	90	70.3	80.9
CrossViT-18† (Rot)	79.22	73.1	75.4	76.7	75	75.7	85.3	72.3	79.2	68.8	70.9	64.3	54.7	67.6	78.4	75.4	76.4
CrossViT-18†	83.05	51	46.1	48.8	46.4	61.2	72.6	54.4	65	44.9	42.1	37.2	41.5	37	67.2	46.8	54.2

Table 4: A table showing the comparison of mean corruption errors (mce)’s across CrossViT models contingent on their training regime. A ResNet50-Augmix is shown as a reference of a particularly strong model to common corruptions. Here lower scores are indicative of better robustness to the different distortion types of Hendrycks & Dietterich (2019).

⁴Also see Li et al. (2022) that shows that generally robust models (robust to adversarial + common corruptions) have a preference for low-spatial frequency statistics.

472 **B.4 ImageNet-R**

473 We also looked into how adversarial training would affect the performance of generalization to
 474 various abstract visual renditions. To do this, we ran our models and benchmarked them on the
 475 ImageNet-Rendition (ImageNet-R) dataset (Hendrycks et al., 2021).

476 We observe that the accuracy on ImageNet-R decreases when the CrossViT is adversarially trained.
 477 However, when we combine the rotation invariance and adversarial training regimes, the accuracy on
 478 ImageNet-R becomes competitive with its pretrained version. In addition, we also appreciate that this
 479 combination does not affect the IID/OOD Gap with respect to the pretrained CrossViT.

Network	ImageNet-200 (↑)	ImageNet-R (↑)	Gap (↓)
CrossViT-18† (Adv + Rot)	90.75	41.14	49.61
CrossViT-18† (Adv)	85.52	35.73	49.79
CrossViT-18† (Rot)	93.89	37.35	56.54
CrossViT-18†	95.64	45.7	49.94

Table 5: A table showing the comparison of the accuracy on Imagenet-R dataset across CrossViT models contingent in their training regime.

480 **C Comparison of CrossViT vs vanilla Transformer (ViT) Models**

481 In this section, we investigated what is the role of the architecture in our results. Did we arrive at
 482 a high-scoring Brain-Score model by virtue of the general Transformer architecture, or was there
 483 something particular about the CrossViT (dual stream Transformer), that in tandem with our training
 484 pipeline allowed for a more ventral-stream like representation? We repeated our analysis and training
 485 procedures with a collection of vanilla Vision Transformers (ViT) where we manipulated the patch
 486 size and number of layers with the conventions of Dosovitskiy et al. (2021) as shown in Figure 9.

487 Here we see that the Brain-Score on V2, V4, superior processing IT, Behavior and Average *increase*
 488 independent of the type of Vision Transformer used for our suite of models (CrossViT-18†, and
 489 multiple ViT flavors) except for the particular case of ViT-S/16 due to over-fitting (See Figure 6) that
 490 heavily reflects on the behavior score. To our surprise, adversarial training in some cases helped V1
 491 score and in some not, potentially due to an interaction with both patch size and transformer depth
 492 that has not fully been understood. In addition, to our knowledge, this is also the first time that it has
 493 been shown that adversarial training coupled with rotational invariance homogeneously increases
 494 brain-scores across Transformer-like architectures, as previous work has shown that classical CNNs
 495 (*i.e.* ResNets) increase Brain-Scores with adversarial training (Dapello et al., 2020). Additionally to
 496 the experiments on CrossViT-18†, we also evaluate the brain-scores on vanilla Vision transformers
 497 that can be seen in Table 6.

Description	ImageNet(↑)	Brain-Score(↑)					
	Validation Acc. (%)	Avg	V1	V2	V4	IT	Behavior
ViT-S/16	81.40	0.445	0.527	0.295	0.454	0.449	0.498
ViT-S/32	75.99	0.415	0.531	0.271	0.422	0.423	0.426
ViT-B/16	84.53	0.451	0.522	0.317	0.398	0.487	0.529
ViT-B/32	80.72	0.440	0.553	0.311	0.413	0.418	0.505
ViT-S/16 (Adv + Rot)	50.44	0.443	0.506	0.332	0.470	0.496	0.409
ViT-S/32 (Adv + Rot)	55.20	0.457	0.512	0.347	0.433	0.485	0.508
ViT-B/16 (Adv + Rot)	67.25	0.486	0.536	0.332	0.470	0.496	0.598
ViT-B/32 (Adv + Rot)	53.01	0.457	0.524	0.357	0.417	0.472	0.515

Table 6: ImageNet accuracy, Brain-Scores of each brain area & Behavior benchmark evaluated on vanilla vision transformers

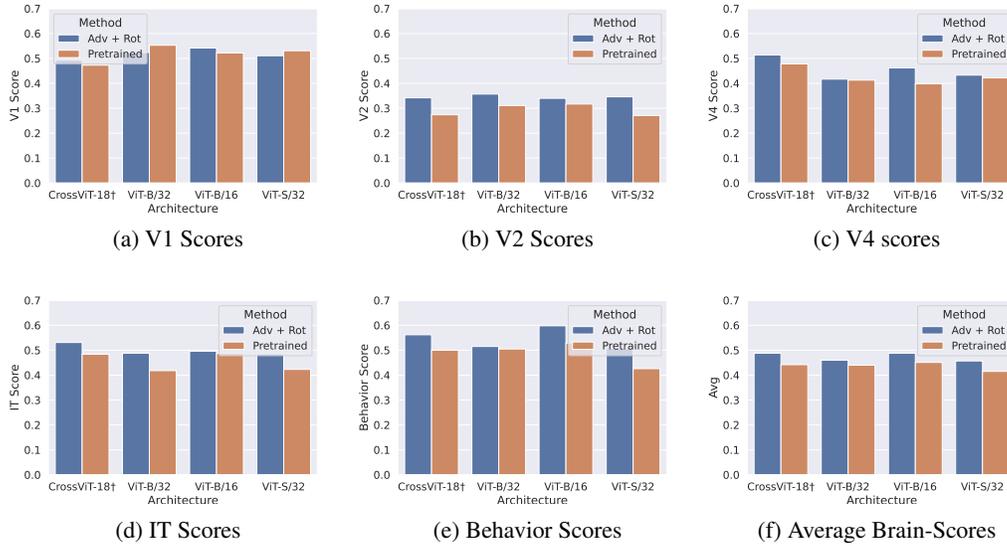


Figure 9: Similarity Brain-Score analysis on the different cortical areas of the ventral stream for vanilla transformers (ViT) and CrossViT. For nearly all Transformer variations, Adversarial Training with Joint Rotational Invariance increases per Area and Average Brain-Scores.

498 D Selection of the Best-BrainScore layers

499 Best performing layers on each vision transformer were selected by a brute-force approach. We
 500 evaluate each layer of the vision transformer models on each brain region and behavior dataset
 501 and select the layer that got the best score on the public benchmarks (in order to avoid overfitting)
 502 proportioned by Brain-Score organization. After this step, the "Adv + Rot" & pretrained versions
 503 of each transformer are submitted to the competition fixing best performing layers (See Table 7).
 504 We achieved our highest score at the time of our 4th submission, which was the lowest number of
 505 submissions in the competition (the winner of the competition performed nearly 60 submissions). All
 506 our results reflect the private scores obtained by each vision transformer model.

Model	V1	V2	V4	IT	Behavior
CrossViT-18f	blocks.1.blocks.1.0.norm1	blocks.1.blocks.1.0.norm1	blocks.1.blocks.1.0.norm1	blocks.1.blocks.1.4.norm2	blocks.2.revert_proj.1.2
ViT-S/16	blocks.1.mlp.act	blocks.3.attn.proj	blocks.3.norm2	blocks.9.norm1	pre_logits
ViT-S/16	blocks.1.mlp.act	blocks.3.attn.proj	blocks.3.norm2	blocks.9.norm1	pre_logits
ViT-S/32	blocks.1.mlp.act	blocks.10.norm1	blocks.2.mlp.act	blocks.10.norm1	pre_logits
ViT-B/16	blocks.1.mlp.act	blocks.6.norm2	blocks.2.mlp.act	blocks.8.norm1	pre_logits
ViT-B/32	blocks.1.mlp.act	blocks.6.norm2	blocks.2.mlp.act	blocks.11.norm1	pre_logits

Table 7: Layers selected for each brain region on each vision transformer.