
Track 1:

Ensemble everything everywhere: Multi-scale aggregation for adversarial robustness

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Adversarial examples pose a significant challenge to the robustness, reliability and
2 alignment of deep neural networks. We propose a novel, easy-to-use approach to
3 achieving high-quality representations that lead to adversarial robustness through
4 the use of multi-resolution input representations and dynamic self-ensembling of
5 intermediate layer predictions. We demonstrate that intermediate layer predictions
6 exhibit inherent robustness to adversarial attacks crafted to fool the full classifier,
7 and propose a robust aggregation mechanism based on Vickrey auction that we
8 call *CrossMax* to dynamically ensemble them. By combining multi-resolution
9 inputs and robust ensembling, we achieve significant adversarial robustness on
10 CIFAR-10 and CIFAR-100 datasets without any adversarial training or extra data,
11 reaching an adversarial accuracy of $\approx 72\%$ (CIFAR-10) and $\approx 48\%$ (CIFAR-100)
12 on the RobustBench AutoAttack suite ($L_\infty = 8/255$) with a finetuned ImageNet-
13 pre-trained ResNet152. This represents a result comparable with the top three
14 models on CIFAR-10 and a +5 % gain compared to the best current dedicated
15 approach on CIFAR-100. Adding simple adversarial training on top, we get
16 $\approx 78\%$ on CIFAR-10 and $\approx 51\%$ on CIFAR-100, improving SOTA by 5 % and
17 9 % respectively and seeing greater gains on the harder dataset. We validate our
18 approach through extensive experiments and provide insights into the interplay
19 between adversarial robustness, and the hierarchical nature of deep representations.
20 We show that simple gradient-based attacks against our model lead to human-
21 interpretable images of the target classes as well as interpretable image changes.
22 As a byproduct, using our multi-resolution prior, we turn pre-trained classifiers and
23 CLIP models into controllable image generators and develop successful transferable
24 attacks on large vision language models.

25 1 Introduction

26 Adversarial examples in the domain of image classification are small, typically human-imperceptible
27 perturbations P to an image X that nonetheless cause a classifier, $f : X \rightarrow y$, to misclassify the
28 perturbed image $X + P$ as a target class t chosen by the attacker, rather than its correct, ground truth
29 class. This is despite the perturbed image $X + P$ still looking clearly like the ground truth class to a
30 human, highlighting a striking and consistent difference between machine and human vision (first
31 described by Szegedy et al. [2013]). Adversarial vulnerability is ubiquitous in image classification,
32 from small models and datasets [Szegedy et al., 2013] to modern large models such CLIP [Radford
33 et al., 2021], and successful attacks transfer between models and architectures to a surprising degree
34 [Goodfellow et al., 2015] without comparable transfer to humans.

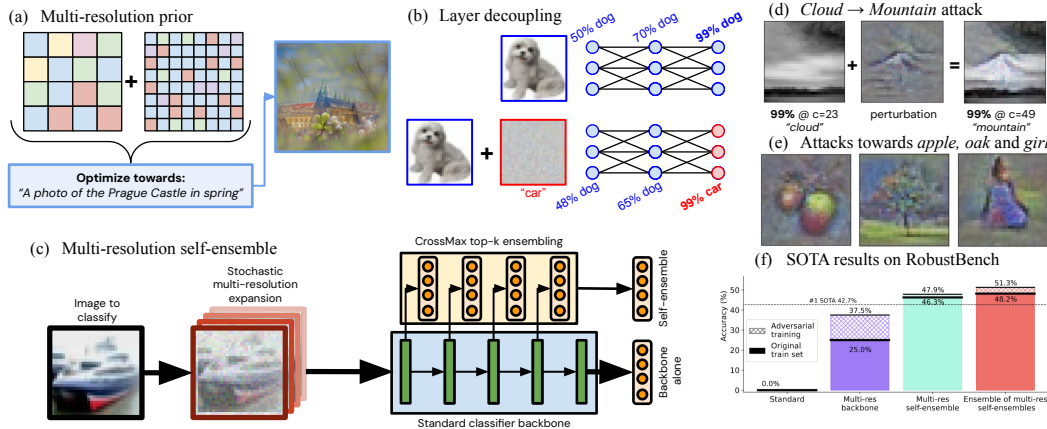


Figure 1: We use a multi-resolution decomposition (a) of an input image and a partial decorrelation of predictions of intermediate layers (b) to build a classifier (c) that has, by default, adversarial robustness comparable or exceeding state-of-the-art (f), even without any adversarial training. Optimizing inputs against it leads to interpretable changes (d) and images generated from scratch (e).

35 We hypothesize that the existence of adversarial attacks is due to the significant yet subtle mismatch
 36 between what humans do when they classify objects and how they learn such a classification in
 37 the first place (the *implicit* classification function in their brains), and what is conveyed to a neural
 38 network classifier explicitly during training by associating fixed pixel arrays with discrete labels (the
 39 learned machine classification function).

40 In this paper, we take a step towards aligning the implicit human and explicit machine classification
 41 functions, and consequently observe very significant gains in adversarial robustness against standard
 42 attacks as a result of a few, simple, well-motivated changes, and without any explicit adversarial
 43 training. While, historically, the bulk of improvement on robustness metrics came from adversarial
 44 training [Chakraborty et al., 2018], comparably little attention has been dedicated to improving the
 45 model backbone, and even less to rethinking the training paradigm itself. Our method can also
 46 be easily combined with adversarial training, further increasing the model’s robustness cheaply.
 47 Beyond benchmark measures of robustness, we show that if we optimize an image against our models
 48 directly, the resulting changes are human interpretable, suggesting at least much-harder-to-find
 49 instances of noise-like superstimuli that we usually find by attacking a model. This suggests an
 50 overall higher-quality, natural representations being learned by the model.

51 We operate under what we call the **Interpretability-Robustness Hypothesis: A model whose**
 52 **adversarial attacks typically look human-interpretable will also be adversarially robust.** The aim
 53 of this paper is to support this hypothesis and to construct first versions of such robust classifiers,
 54 without necessarily reaching their peak performance via extensive hyperparameter tuning.

55 Firstly, inspired by biology, we design an active adversarial defense by constructing and training a
 56 classifier whose input, a standard $H \times W \times 3$ image, is stochastically turned into a $H \times W \times (3N)$
 57 channel-wise stack of multiple downsampled and noisy versions of the same image. The classifier
 58 itself learns to make a decision about these N versions *at once*, mimicking the effect of microsaccades
 59 in the human (and mammal) vision systems. Secondly, we show experimentally that hidden layer
 60 features of a neural classifier show significant decorrelation between their representations under
 61 adversarial attacks – an attack fooling a network to see a *dog* as a *car* does not fool the intermediate
 62 representations, which still see a *dog*. We aggregate intermediate layer predictions into a self-
 63 ensemble dynamically, using a novel ensembling technique that we call a *CrossMax* ensemble.
 64 Thirdly, we show that our Vickrey-auction-inspired *CrossMax* ensembling yields significant gains
 65 in adversarial robustness when ensembling predictors as varied as 1) independent brittle models, 2)
 66 predictions of intermediate layers of the same model, 3) predictions from several checkpoints of
 67 the same model, and 4) predictions from several self-ensemble models. We use the last option to
 68 gain $\approx 5\%$ in adversarial accuracy at the $L_\infty = 8/255$ RobustBench’s AutoAttack on top of the best
 69 models on CIFAR-100. When we add light adversarial training on top, we outperform current best
 70 models by $\approx 5\%$ on CIFAR-10, and by $\approx 9\%$ on CIFAR-100, showing a promising trend where the
 71 harder the dataset, the more useful our approach compared to brute force adversarial training (see
 72 Figure 4).

73 **2 Key Observations and Techniques**

74 In this section we will describe the three key methods that we use in this paper. In Section 2.1
 75 we introduce the idea of multi-resolution inputs, in Section 2.2 we introduce our robust *CrossMax*
 76 ensembling method, and in Section 2.3 we showcase the de-correlation between adversarially induced
 mistakes at different layers of the network and how to use it as an active defense.

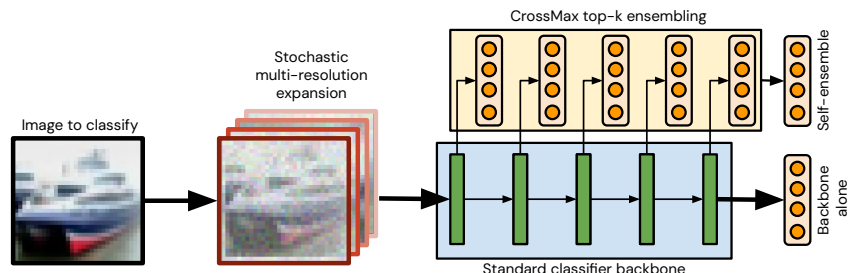


Figure 2: Combining channel-wise stacked augmented and down-sampled versions of the input image with robust intermediate layer class predictions via *CrossMax* self-ensemble. The resulting model gains a considerable adversarial robustness without any adversarial training or extra data.

77

78 **2.1 The multi-resolution prior**

79 Drawing inspiration from biology, we use multiple versions of the same image at once, down-sampled
 80 to lower resolutions and augmented with stochastic jitter and noise. We train a model to classify this
 81 channel-wise stack of images simultaneously. We show that this by default yields gains in adversarial
 82 robustness without any explicit adversarial training.

83 We turn an input image X of full resolution $R \times R$ and 3 channels (RGB) into its N variations of
 84 different resolutions $r \times r$ for $r \in \rho$. For CIFAR-10 and CIFAR-100, we are (arbitrarily) choosing res-
 85 olutions $\rho = \{32, 16, 8, 4\}$ and concatenating the resulting image variations $\text{rescale}_R(\text{rescale}_r(X))$
 86 channel-wise to a $R \times R \times (3|\rho|)$ augmented image \tilde{X} . This is shown in Figure 7. Similar approaches
 87 have historically been used to represent images, such as Gaussian pyramids introduced in [Burt and
 88 Adelson, 1983]. To each variant we add 1) random noise both when downsampled and at the full
 89 resolution $R \times R$ (in our experiments of strength 0.1 out of 1.0), 2) a random jitter in the $x - y$
 90 plane (± 3 in our experiments), 3) a small, random change in contrast, and 4) a small, random
 91 color-grayscale shift. This can also be seen as an effective reduction of the input space dimension
 92 available to the attacker, as discussed in [Fort, 2023].

93 **2.2 CrossMax robust ensembling**

94 The standard way of ensembling predictions of multiple networks is to either take the mean of their
 95 logits, or the mean of their probabilities. This increases both the accuracy as well as predictive
 96 uncertainty estimates of the ensemble [Lakshminarayanan et al., 2017, Ovadia et al., 2019]. Such
 97 aggregation methods are, however, susceptible to being swayed by an outlier prediction by a single
 98 member of the ensemble or its small subset. This produces a single point of failure. The pitfalls of
 99 uncertainty estimation and ensembling have been highlighted in [Ashukha et al., 2021], while the
 100 effect of ensembling on the learned classification function was studied by Fort et al. [2022].

101 We draw our intuition from *Vickrey auctions* [Wilson, 1977] which are designed to incentivize truthful
 102 bidding. Viewing members of ensembles as individual bidders, we can limit the effect of wrong,
 103 yet overconfident predictions by using the 2nd highest, or generally k^{th} highest prediction per class.
 104 This also produces a cat-and-mouse-like setup for the attacker, since *which* classifier produces the
 105 k^{th} highest prediction for a particular class changes dynamically as the attacker tries to increase that
 106 prediction. A similar mechanism is used in balanced allocation [Azar et al., 1999] and specifically
 107 in the *k random choices* algorithm for load balancing [Mitzenmacher et al., 2001]. Our *CrossMax*
 108 aggregation is shown in Algorithm 1.

109 **2.3 Only partial overlap between the adversarial susceptibility of intermediate layers**

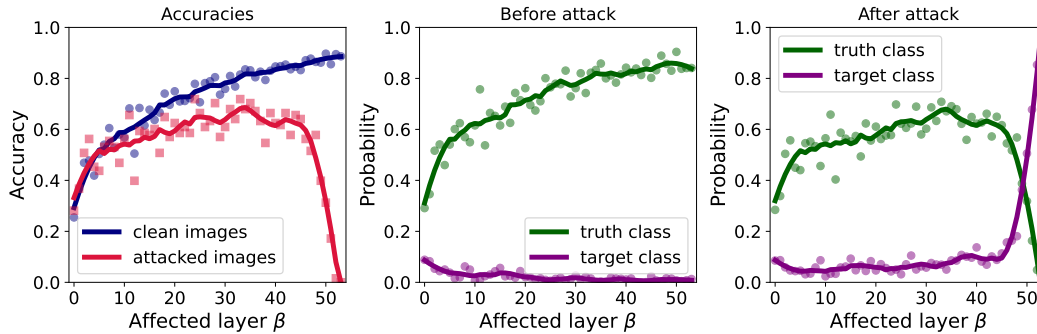


Figure 3: The impact of adversarial attacks ($L_\infty = 8/255$, 128 attacks) against the full classifier on the accuracy and probabilities at all intermediate layers for an ImageNet-1k pretrained ResNet152 finetuned on CIFAR-10 via trained linear probes.

110 A key question of both scientific and immediately practical interest is whether an adversarially
 111 modified image X' that looks like the target class t to a classifier $f : X \rightarrow y$ also has intermediate
 112 layer representations that look like that target class. In [Olah et al., 2017], it is shown via feature
 113 visualization that neural networks build up their understanding of an image hierarchically starting
 114 from edges, moving to textures, simple patterns, all the way to parts of objects and full objects
 115 themselves. This is further explored by Carter et al. [2019]. Does an image of a *car* that has been
 116 adversarially modified to look like a *tortoise* to the final layer classifier carry the intermediate features
 117 of the target class *tortoise* (e.g. the patterns on the shell, the legs, a tortoise head), of the original
 118 class *car* (e.g. wheels, doors), or something else entirely? We answer this question empirically.

119 In Figure 3 we showcase this effect using an ImageNet-pretrained ResNet152 [He et al., 2015]
 120 finetuned on CIFAR-10. Images attacked to look like some other class than their ground truth (to
 121 the final layer classification) do not look like that to intermediate layers, as shown by the target class
 122 probability only rising in the very last layers (see Figure 3). We can therefore confirm that indeed the
 123 activations of attacked images do not look like the target class in the intermediate layers, which offers
 124 two immediate use cases: 1) as a warning flag that the image has been tampered with and 2) as an
 125 active defense, which is strictly harder.

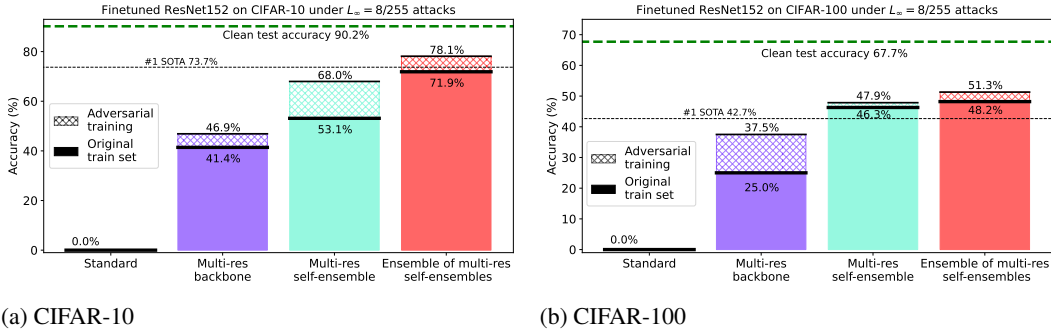
126 **3 Training and Experimental Results**

127 In this section we present in detail how we combine the previously described methods and techniques
 128 into a robust classifier on CIFAR-10 and CIFAR-100. We start both with a pretrained model and
 129 finetune it, as well as with a freshly initialized model. It turns out that finetuning a pre-existing model
 130 for robustness is technically easier and faster, therefore we predominantly focus on this approach.
 131 However, to demonstrate that the success of our technique does not simply come from massive
 132 pretraining, we also train a model from scratch. The concrete details of the model and training can be
 133 found in Appendix A.

134 **3.1 Adversarial vulnerability evaluation**

135 To make sure we are using as strong an attack suite as possible to measure our networks' robustness
 136 and to be able to compare our results to other approaches, we use the RobustBench [Croce et al.,
 137 2020] library and its AutoAttack method, which runs a suite of four strong, consecutive adversarial
 138 attacks on a model in a sequence and estimates its adversarial accuracy (e.g. if the attacked images
 139 were fed back to the network, what would be the classification accuracy with respect to their ground
 140 truth classes). To evaluate our models using the hardest method possible, we ran the AutoAttack
 141 with the rand flag that is tailored against models using randomness. The results without adversarial
 142 training are shown in Table 1 and with adversarial training at Table 2. The visual representation of
 143 the results is presented in Figure 4.

144 **3.2 Multi-resolution finetuning of a pretrained model**



(a) CIFAR-10 (b) CIFAR-100

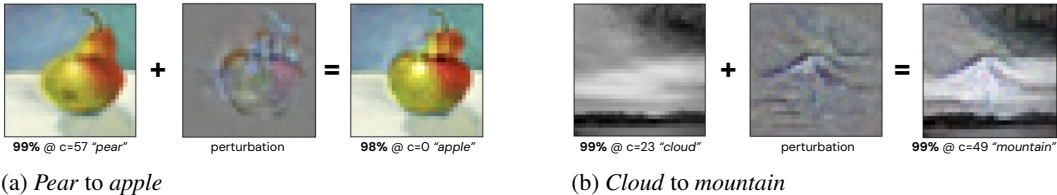
Figure 4: Adversarial robustness evaluation for finetuned ResNet152 models under $L_\infty = 8/255$ attacks of RobustBench AutoAttack (*rand* version, which is stronger against our models). On CIFAR-10, a CrossMax 3-ensemble of our self-ensemble multi-resolution models reaches #3 on the leaderboard, while on CIFAR-100 a 3-ensemble of our multi-resolution models is #1, leading by $\approx +5\%$ in adversarial accuracy. When we add light adversarial training, our models surpass SOTA on CIFAR-10 by $\approx +5\%$ and on CIFAR-100 by a strong $\approx +9\%$.

145 We demonstrate that this quickly leads to very significant adversarial robustness that matches and in
 146 some cases (CIFAR-100) significantly improves upon current best, dedicated approaches, without
 147 using any extra data or adversarial training. We see stronger gains on CIFAR-100 rather than CIFAR-
 148 10, suggesting that its edge might lie at harder datasets, which is a very favourable scaling compared
 149 to brute force adversarial training.

150 The steps we take are as follows: 1) Take a pretrained model (in our case ResNet18 and ResNet152
 151 pretrained on ImageNet) 2) Replace the first layer with a fresh initialization that can take in $3N$
 152 instead of 3 channels 3) Replace the final layer with a fresh initialization to project to 10 (for CIFAR-
 153 10) or 100 (for CIFAR-100) classes 4) Train the full network with a *small* (this is key) learning rate
 154 for a few epochs.

155 We find that using a small learning rate is key, which could be connected to the effects described for
 156 example in Thilak et al. [2022] and Fort et al. [2020]. While the network might reach a good clean
 157 test accuracy for high learning rates as well, only for small learning rates will it also get significantly
 158 robust against adversarial attacks, as shown in Figure 9. In Table 1 we present our results of finetuning
 159 an ImageNet pretrained ResNet152 on CIFAR-10 and CIFAR-100 for 10 epochs at the constant
 160 learning rate of 3.3×10^{-5} with Adam followed by 3 epochs at 3.3×10^{-6} . The details of our light
 161 adversarial finetuning are discussed in Appendix B.

162 **3.3 Visualizing attacks against multi-resolution models**



(a) Pear to apple (b) Cloud to mountain

Figure 5: Examples of an adversarial attack on an image towards a target label. We use simple gradient steps with respect to our multi-resolution ResNet152 finetuned on CIFAR-100. The resulting attacks use the underlying features of the original image and make semantically meaningful, human-interpretable changes to it. Additional examples available in Figure 21.

163 We wanted to visualize the attacks against our multi-resolution models. In Figure 5 we start with a
 164 test set image of CIFAR-100 (a pear, cloud, camel and elephant) and over 400 steps with SGD and
 165 $\eta = 1$ minimize the loss with respect to a target class (apple, mountain, rabbit and dinosaur). We

166 allow for large perturbations, up to $L_\infty = 128/255$, to showcase the alignment between our model
 167 and the implicit human visual system classification function. In case of the *pear*, the perturbation
 168 uses the underlying structure of the fruit to divide it into 2 apples by adding a well-placed edge. The
 169 resulting image is very obviously an apple to a human as well as the model itself. In case of the cloud,
 170 its white color is repurposed by the attack to form the snow of a mountain, which is drawn in by a
 171 dark sharp contour. In case of the elephant, it is turned into a dinosaur by being recolored to green
 172 and made spikier – all changes that are very easily interpretable to a human.

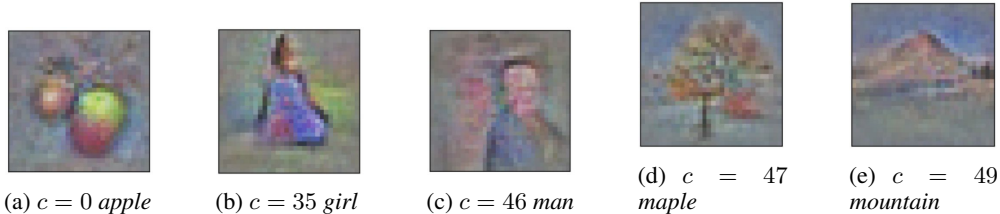


Figure 6: Examples of adversarial attacks on our multi-resolution ResNet152 finetuned on CIFAR-100. The attacks are generated by starting from a uniform image (128,128,128) and using gradient descent of the cross-entropy loss with SGD at $\eta = 1$ for 400 steps towards the target label.

173 In Figure 6 we start with a uniform gray image of color (128, 128, 128) and by changing it to
 174 maximize the probability of a target class with respect to our model, we generate an image. The
 175 resulting images are very human-interpretable. We also generated 4 examples per CIFAR-100 class
 176 for all 100 classes in Figure 23 to showcase that we do not cherry-pick the images shown.

177 4 Discussion and Conclusion

178 Our work demonstrates that taking inspiration from biology and stochastically translating an input
 179 image into a multi-resolution stack of inputs that are classified *simultaneously* by a model leads to
 180 higher-quality, natural representations, significant adversarial robustness, and human-interpretable
 181 attacks. Combining this with a novel, robust ensembling method inspired by Vickrey auctions that
 182 we call *CrossMax*, we demonstrate that we can further improve the model’s adversarial robustness
 183 by combining its intermediate layer predictions into a *self-ensemble*. This is due to our empirical
 184 observation that intermediate layer representations are not fooled by attacks against the classifier as a
 185 whole, and that their induced errors are only partially correlated.

186 We are able to match the current state-of-the-art adversarial accuracy results on CIFAR-10 and surpass
 187 them by $\approx 5\%$ CIFAR-100 on a strong adversarial benchmark RobustBench without any extra data
 188 or dedicated adversarial training, that is usually needed to produce a robust model. When we add
 189 light adversarial training on top, we see that our methods are complementary to it and that we surpass
 190 the best models on CIFAR-10 by $\approx 5\%$ and by a very significant $\approx 9\%$ on CIFAR-100, taking it
 191 from $\approx 40\%$ to $\approx 50\%$ in a single step. Our methods seem to perform better on the harder dataset,
 192 suggesting a favourable scaling compared to the usual brute force adversarial training.

193 Our approach not only enhances robustness but also aligns the learned representations more closely
 194 with human visual processing, leading to more interpretable and reliable models. We demonstrate
 195 this by optimizing images against the outputs of our classifier directly and obtaining either human-
 196 interpretable changes, when applied to an existing image, or completely new, interpretable images,
 197 when starting from a uniform, empty image. This is in stark contrast to the usual result of such a
 198 procedure which would be a noise-like picture that would look very convincing to the network but
 199 would not resemble anything to a human.

200 References

- 201 Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow,
202 and Rob Fergus. Intriguing properties of neural networks, 2013.
- 203 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
204 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever.
205 Learning transferable visual models from natural language supervision, 2021.
- 206 Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial
207 examples, 2015. URL <https://arxiv.org/abs/1412.6572>.
- 208 Anirban Chakraborty, Manaar Alam, Vishal Dey, Anupam Chattopadhyay, and Debdeep Mukhopad-
209 hyay. Adversarial attacks and defences: A survey, 2018. URL <https://arxiv.org/abs/1810.00069>.
- 210
- 211 P. Burt and E. Adelson. The laplacian pyramid as a compact image code. *IEEE Transactions on*
212 *Communications*, 31(4):532–540, 1983. doi: 10.1109/TCOM.1983.1095851.
- 213 Stanislav Fort. Scaling laws for adversarial attacks on language model activations, 2023.
- 214 Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive
215 uncertainty estimation using deep ensembles, 2017.
- 216 Yaniv Ovadia, Emily Fertig, Jie Ren, Zachary Nado, D Sculley, Sebastian Nowozin, Joshua V. Dillon,
217 Balaji Lakshminarayanan, and Jasper Snoek. Can you trust your model’s uncertainty? evaluating
218 predictive uncertainty under dataset shift, 2019.
- 219 Arsenii Ashukha, Alexander Lyzhov, Dmitry Molchanov, and Dmitry Vetrov. Pitfalls of in-domain
220 uncertainty estimation and ensembling in deep learning, 2021.
- 221 Stanislav Fort, Ekin Dogus Cubuk, Surya Ganguli, and Samuel S. Schoenholz. What does a deep
222 neural network confidently perceive? the effective dimension of high certainty class manifolds and
223 their low confidence boundaries, 2022. URL <https://arxiv.org/abs/2210.05546>.
- 224 Robert B. Wilson. Counterspeculation, auctions, and competitive sealed tenders. *Journal of Finance*,
225 31(3):1106–1115, 1977.
- 226 Yossi Azar, Andrei Z Broder, Anna R Karlin, and Eli Upfal. Balanced allocations. *SIAM Journal on*
227 *Computing*, 29:180–200, 1999.
- 228 Michael Mitzenmacher, Andrea W. Richa, and Ramesh Sitaraman. The power of two random choices:
229 A survey of techniques and results. *Harvard University*, 2001.
- 230 Chris Olah, Alexander Mordvintsev, and Ludwig Schubert. Feature visualization. *Distill*, 2017. doi:
231 10.23915/distill.00007. <https://distill.pub/2017/feature-visualization>.
- 232 Shan Carter, Zan Armstrong, Ludwig Schubert, Ian Johnson, and Chris Olah. Activation atlas. *Distill*,
233 2019. doi: 10.23915/distill.00015. <https://distill.pub/2019/activation-atlas>.
- 234 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
235 recognition, 2015.
- 236 Francesco Croce, Maksym Andriushchenko, Vikash Sehwal, Edoardo Debenedetti, Nicolas Flam-
237 marion, Mung Chiang, Prateek Mittal, and Matthias Hein. Robustbench: a standardized adversarial
238 robustness benchmark, 2020.
- 239 Vimal Thilak, Etai Littwin, Shuangfei Zhai, Omid Saremi, Roni Paiss, and Joshua Susskind. The
240 slingshot mechanism: An empirical study of adaptive optimizers and the grokking phenomenon,
241 2022.
- 242 Stanislav Fort, Gintare Karolina Dziugaite, Mansheej Paul, Sepideh Kharaghani, Daniel M. Roy,
243 and Surya Ganguli. Deep learning versus kernel learning: an empirical study of loss landscape
244 geometry and the time evolution of the neural tangent kernel, 2020. URL <https://arxiv.org/abs/2010.15110>.
- 245

- 246 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale
247 hierarchical image database. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009.*
248 *IEEE Conference on*, pages 248–255. IEEE, 2009. URL [https://ieeexplore.ieee.org/
249 abstract/document/5206848/](https://ieeexplore.ieee.org/abstract/document/5206848/).
- 250 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Im-
251 age Recognition. In *Proceedings of 2016 IEEE Conference on Computer Vision and Pattern*
252 *Recognition*, CVPR '16, pages 770–778. IEEE, June 2016. doi: 10.1109/CVPR.2016.90. URL
253 <http://ieeexplore.ieee.org/document/7780459>.
- 254 Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by
255 reducing internal covariate shift, 2015.
- 256 Alex Krizhevsky. Learning multiple layers of features from tiny images. pages 32–33, 2009. URL
257 <https://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf>.
- 258 Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar-100 (canadian institute for advanced
259 research). URL <http://www.cs.toronto.edu/~kriz/cifar.html>.
- 260 ShengYun Peng, Weilin Xu, Cory Cornelius, Matthew Hull, Kevin Li, Rahul Duggal, Mansi Phute,
261 Jason Martin, and Duen Horng Chau. Robust principles: Architectural design principles for
262 adversarially robust cnns, 2023.
- 263 Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International*
264 *Conference on Learning Representations (ICLR)*, San Diego, CA, USA, 2015.
- 265 Jiequan Cui, Zhuotao Tian, Zhisheng Zhong, Xiaojuan Qi, Bei Yu, and Hanwang Zhang. Decoupled
266 kullback-leibler divergence loss, 2023. URL <https://arxiv.org/abs/2305.13948>.
- 267 Zekai Wang, Tianyu Pang, Chao Du, Min Lin, Weiwei Liu, and Shuicheng Yan. Better diffusion mod-
268 els further improve adversarial training, 2023. URL <https://arxiv.org/abs/2302.04638>.
- 269 Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond empirical
270 risk minimization, 2018.
- 271 Brian R. Bartoldson, James Diffenderfer, Konstantinos Parasyris, and Bhavya Kailkhura. Adversarial
272 robustness limits via scaling-law and human-alignment studies, 2024.
- 273 A van der Schaaf and J H van Hateren. Modelling the power spectra of natural images: Statistics and
274 information. *Vision Research*, 36(17):2759–2770, September 1996. ISSN 0042-6989. Relation:
275 <http://www.rug.nl/informatica/organisatie/overorganisatie/iwi> Rights: University of Groningen.
276 Research Institute for Mathematics and Computing Science (IWI).
- 277 Kurt Hornik, Maxwell Stinchcombe, and Halbert White. Multilayer feedforward networks are
278 universal approximators. *Neural Networks*, 2(5):359–366, 1989.
- 279 Herbert Robbins and Sutton Monro. A stochastic approximation method. *The Annals of Mathematical*
280 *Statistics*, 22(3):400–407, 1951.
- 281 Robert G Keys. Cubic convolution interpolation for digital image processing. *IEEE Transactions on*
282 *Acoustics, Speech, and Signal Processing*, 29(6):1153–1160, 1981.
- 283 Stanislav Fort. Pixels still beat text: Attacking the openai clip model with text
284 patches and adversarial pixel perturbations. URL https://stanislavfort.github.io/blog/OpenAI_CLIP_stickers_and_adversarial_examples, 2021a.
- 286 Stanislav Fort. Adversarial examples for the openai clip in its zero-shot classification
287 regime and their semantic generalization, jan 2021b. URL https://stanislavfort.github.io/2021/01/12/OpenAI_CLIP_adversarial_examples.html, 2021b.
- 289 Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori,
290 Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali
291 Farhadi, and Ludwig Schmidt. Openclip, July 2021. URL [https://doi.org/10.5281/zenodo.
292 5143773](https://doi.org/10.5281/zenodo.5143773). If you use this software, please cite it as below.

- 293 Mehdi Cherti, Romain Beaumont, Ross Wightman, Mitchell Wortsman, Gabriel Ilharco, Cade
294 Gordon, Christoph Schuhmann, Ludwig Schmidt, and Jenia Jitsev. Reproducible scaling laws for
295 contrastive language-image learning. In *Proceedings of the IEEE/CVF Conference on Computer
296 Vision and Pattern Recognition*, pages 2818–2829, 2023.
- 297 Rylan Schaeffer, Dan Valentine, Luke Bailey, James Chua, Cristóbal Eyzaguirre, Zane Durante, Joe
298 Benton, Brando Miranda, Henry Sleight, John Hughes, Rajashree Agrawal, Mrinank Sharma, Scott
299 Emmons, Sanmi Koyejo, and Ethan Perez. When do universal image jailbreaks transfer between
300 vision-language models?, 2024. URL <https://arxiv.org/abs/2407.15211>.
- 301 Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an ensemble
302 of diverse parameter-free attacks, 2020.

303 A Model and training details

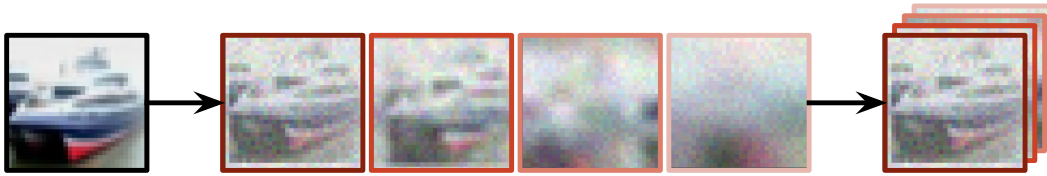


Figure 7: An image input being split into N progressively lower resolution versions that are then stacked channel-wise, forming a $3N$ -channel image input to a classifier.

304 The pretrained models we use are the ImageNet [Deng et al., 2009] trained ResNet18 and ResNet152
305 [He et al., 2016]. Our hyperparameter search was very minimal and we believe that additional gains
306 are to be had with a more involved search easily. The only architectural modification we make is to
307 change the number of input channels in the very first convolutional layer from 3 to $3N$, where N is
308 the number of channel-wise stacked down-sampled images we use as input. We also replaced the
309 final linear layer to map to the correct number of classes (10 for CIFAR-10 and 100 for CIFAR-100).
310 Both the new convolutional layer as well as the final linear layer are initialized at random. The batch
311 norm [Ioffe and Szegedy, 2015] is on for finetuning a pretrained model (although we did not find a
312 significant effect beyond the speed of training).

313 We focused on the CIFAR-* datasets [Krizhevsky, 2009, Krizhevsky et al.] that comprise 50,000
314 $32 \times 32 \times 3$ images. We arbitrarily chose $N = 4$ and the resolutions we used are 32×32 , 16×16 ,
315 8×8 , 4×4 (see Figure 7). We believe it is possible to choose better combinations, however, we
316 did not run an exhaustive hyperparameter search there. The ResNets we used expect 224×224
317 inputs. We therefore used a bicubic interpolation to upsample the input resolution for each of the
318 12 channels independently.

319 To each image (the $32 \times 32 \times 3$ block of RGB channels) we add a random jitter in the $x - y$ plane in
320 the ± 3 range. We also add a random noise of standard deviation 0.2 (out of 1.0). We believe that
321 the biological jitter and noise are key aspects of a successful robust classifier, and therefore want to
322 mimic their function here as well.

323 For training from scratch, we use a standard ResNet18 with the modifications above. We chose it
324 since we primarily wanted to show the effect of multi-resolution inputs and multi-layer prediction
325 aggregation rather than to find the maximum possible performance. We turn off batch normalization
326 [Ioffe and Szegedy, 2015] not to conflate the effects we are exploring. While it is possible that
327 additional architectural choices could lead to more robustness (as convincingly demonstrated in Peng
328 et al. [2023]), we wanted to show the effect of our multi-resolution and self-ensemble choices in
329 isolation.

330 All training is done using the Adam [Kingma and Ba, 2015] optimizer at a flat learning rate η that we
331 always specify. Optimization is applied to all trainable parameters and the batch norm is turned on in
332 case of finetuning, but turned off for training from scratch.

333 Linear probes producing predictions at each layer are just single linear layers that are trained on top
334 of the pre-trained and frozen backbone network, mapping from the number of hidden neurons in
335 that layer (flattened to a single dimension) to the number of classes (10 for CIFAR-10 and 100 for
336 CIFAR-100). We trained them using the same learning rate as the full network for 1 epoch each.

337 B Adversarial finetuning

338 Adversarial training, which adds attacked images with their correct, ground truth labels back to
339 the training set, is a standard brute force method for increasing models’ adversarial robustness.
340 [Chakraborty et al., 2018] It is ubiquitous among the winning submissions on the RobustBench leader
341 board, e.g. in Cui et al. [2023] and Wang et al. [2023]. To verify that our technique does not only
342 somehow replace the need for dedicated adversarial training, but rather that it can be productively
343 combined with it for even stronger adversarial robustness, we re-ran all our finetuning experiments
344 solely on adversarially modified batches of input images generated on the fly.

345 In all cases, we see an additive benefit of adversarial training on top of our techniques. In particular,
 346 for CIFAR-10 we outperform current SOTA by approximately 5 % while on CIFAR-100 and by
 347 approximately 9 % on CIFAR-100, which is a very large increase. The fact that our techniques benefit
 348 even from a very small amount of additional adversarial training (units of epochs of a single step
 349 attack) shows that our multi-resolution inputs and intermediate layer aggregation are a good prior for
 350 getting broadly robust networks.

351 B.1 Training from scratch

352 We train a ResNet18 from scratch on CIFAR-10 as a back-
 353 bone, and then train additional linear heads for all of its
 354 intermediate layers to form a CrossMax self-ensemble. We
 355 find that, during training, augmenting our input images X
 356 with an independently drawn images X' with a randomly
 357 chosen mixing proportion p as $(1 - p)X + pX'$ increases
 358 the robustness of the trained model. This simple augmen-
 359 tation technique is known as mixup and is described in
 360 Zhang et al. [2018]. We believe that this works well due
 361 to our multi-resolution inputs that are the correct prior for
 362 robustness, and show that without them such mixing does
 363 not increase robustness. For finetuning a pretrained model,
 364 however, this is not needed.

365 For our ResNet18 model trained from scratch on CIFAR-
 366 10, we keep the pairs of images that are mixed in mixup
 367 fixed for 20 epochs at a time, producing a characteristic
 368 pattern in the training accuracies. Every 5 epochs we
 369 re-draw the random mixing proportions in the $[0, 1/2]$
 370 range. We trained the model for 380 epochs with the
 371 Adam optimizer [Kingma and Ba, 2015] at learning rate
 372 10^{-3} and dropped it to 10^{-4} for another 120 epochs. The
 373 final checkpoint is the weight average of the last 3 epochs.
 374 The training batch size is 512. These choices are arbitrary
 375 and we did not run a hyperparameter search over them.

376 The results on the full RobustBench AutoAttack suite of attacks for CIFAR-10 are shown in Table 1
 377 for self-ensemble constructed on top of the multi-resolution ResNet18 backbone (the linear heads on
 378 top of each layer were trained for 2 epochs with Adam at 10^{-3} learning rate).

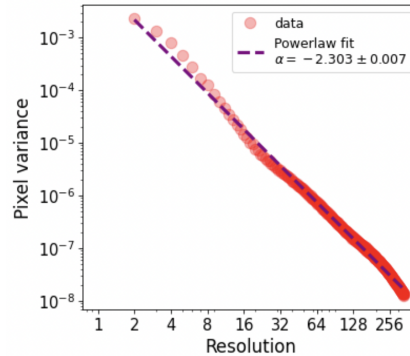


Figure 8: The image spectrum of generated multi-resolution attacks. The adversarial attacks generated over multiple resolutions at once end up showing very white-noise-like distribution of powers over frequencies (the slope for natural images is ≈ -2). This is in contrast with standard noise-like attacks.

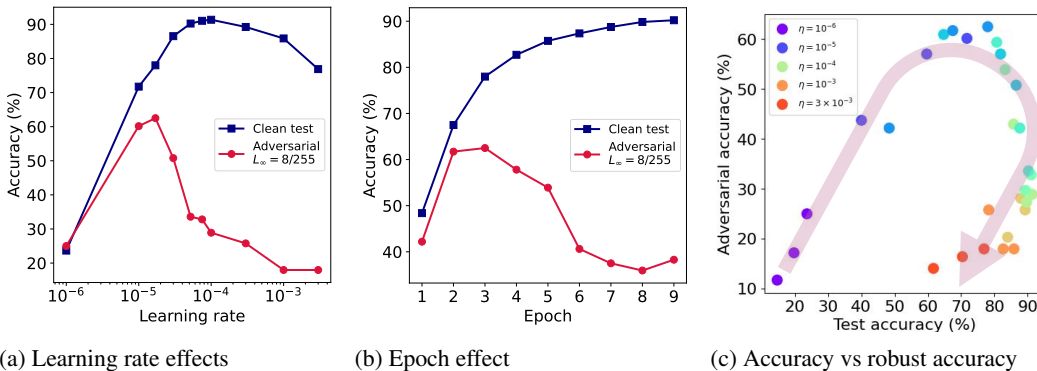


Figure 9: Finetuning a pretrained model with multi-resolution inputs. The left panel shows the test accuracy and adversarial accuracy after the first two attacks of RobustBench AutoAttack at $L_\infty = 8/255$ after 3 epochs of finetuning an ImageNet pretrained ResNet152. The middle panel shows the effect of training epoch for a single finetuning run at the learning rate $\eta = 1.7 \times 10^{-5}$. The right panel shows a hysteresis-like curve where high test accuracies are both compatible with low and high adversarial accuracies. The test accuracies are over the full 10,000 images while the adversarial accuracies are evaluated on 128 test images.

C Results tables

Dataset	Adv. train	Model	Method	#	Test acc	rand RobustBench AutoAttack $L_\infty = 8/255$ # samples (%)		
						Adv acc	APGD→ CE	APGD DLR
CIFAR-10	×	ResNet18*	Self-ensemble	1024	76.94	64.06	51.56	44.53
CIFAR-10	×	ResNet152	Multi-res backbone	128	89.17	41.44	32.81	21.88
CIFAR-10	×	ResNet152	3-ensemble	128	91.06	67.97	61.72	59.38
CIFAR-10	×	ResNet152	Self-ensemble	128	87.14	53.12	50.00	43.75
CIFAR-10	×	ResNet152	3-ensemble of self-ensembles	128	90.20	71.88	68.75	68.75
CIFAR-10	✓	[30]	SOTA #1			73.71		
CIFAR-100	×	ResNet152	Multi-res backbone	128	65.70	25.00	21.88	13.28
CIFAR-100	×	ResNet152	3-ensemble	128	66.63	47.66	39.06	37.50
CIFAR-100	×	ResNet152	Self-ensemble	512	65.71	46.29	34.77	30.08
						±2.36	±2.09	±2.13
CIFAR-100	×	ResNet152	3-ensemble of self-ensembles	512	67.71	48.16	40.63	37.32
						±2.65	±2.11	±1.98
CIFAR-100	✓	[28]	SOTA #1			42.67		

Table 1: Full *randomized* (=the strongest against our approach) RobustBench AutoAttack adversarial attack suite results for 128 test samples at the $L_\infty = 8/255$ strength. In this table we show the results of attacking our multi-resolution ResNet152 models finetuned on CIFAR-10 and CIFAR-100 from an ImageNet pretrained state without any adversarial training or extra data for 20 epochs with Adam at $\eta = 3.3 \times 10^{-5}$. We use our *CrossMax* ensembling on the model itself (self-ensemble), the final 3 epochs (3-ensemble), and on self-ensembles from 3 different runs (3-ensemble of self-ensembles). We also include results for a ResNet18 trained from *scratch* on CIFAR-10. Despite its simplicity, our method gets adversarial robustness of $\approx 72\%$ on CIFAR-10 (ranking #3 on RobustBench leaderboard) and $\approx 48\%$ on CIFAR-100, surpassing current best models by +5%. Unlike other approaches, we do not use any extra data or adversarial training and our models gain adversarial robustness by default. Additional adversarial training helps, as shown in Table 2.

Dataset	Adv. train	Model	Method	#	Test acc	rand RobustBench AutoAttack $L_\infty = 8/255$ # samples (%)		
						Adv acc	APGD→ CE	APGD DLR
CIFAR-10	✓	ResNet152	Multi-res backbone	128	87.19	46.88	34.38	32.03
CIFAR-10	✓	ResNet152	Self-ensemble	128	84.58	67.94	64.06	54.69
CIFAR-10	✓	ResNet152	3-ensemble of self-ensembles	128	87.00	78.13	73.44	72.65
CIFAR-10	✓	[30]	SOTA #1			73.71		
CIFAR-100	✓	ResNet152	Multi-res backbone	128	62.72	37.50	32.03	22.66
CIFAR-100	✓	ResNet152	Self-ensemble	512	58.93	47.85	36.72	33.98
						±2.66	±3.01	±2.72
CIFAR-100	✓	ResNet152	3-ensemble of self-ensembles	512	61.17	51.28	44.60	43.04
						±1.95	±2.00	±1.97
CIFAR-100	✓	[28]	SOTA #1			42.67		

Table 2: Full *randomized* (=the strongest against our approach) RobustBench AutoAttack adversarial attack suite results for 128 test samples at the $L_\infty = 8/255$ strength. In this table we show the results of attacking our multi-resolution ResNet152 models finetuned on CIFAR-10 and CIFAR-100 from an ImageNet pretrained state **with** light adversarial training.

380 D Additional Insights and Applications

381 We want to support our multi-resolution input choice as an active defense by demonstrating that
 382 by reversing it and representing an adversarial perturbation *explicitly* as a sum of perturbations at
 383 different resolutions, we get human-interpretable perturbations by default.

384 E Single-resolution adversarial attacks

385 Natural images contain information expressed on all frequencies, with an empirically observed
 386 power-law scaling. The higher the frequency, the lower the spectral power, as $\propto f^{-2}$ [van der Schaaf
 387 and van Hateren, 1996].

388 While having a single perturbation P of the full resolution $R \times R$ theoretically suffices to express
 389 anything, we find that this choice induces a specific kind of high frequency prior. Even simple neural
 390 networks can theoretically express any function [Hornik et al., 1989], yet the specific architecture
 391 matters for what kind of a solution we obtain given our data, optimization, and other practical choices.
 392 Similarly, we find that an alternative formulation of the perturbation P leads to more natural looking
 393 and human interpretable perturbations despite the attacker having access to the highest-resolution
 394 perturbation as well and could in principle just use that.

395 F Multi-resolution attacks

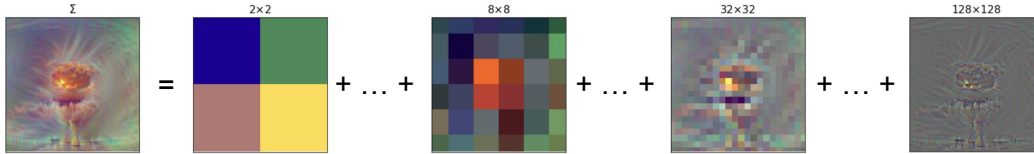


Figure 10: The result of expressing an image as a set of resolutions and optimizing it towards the CLIP embedding of the text 'a photo of a nuclear explosion'. The plot shows the resulting sum of resolutions (left panel, marked with ρ) and selected individual perturbations P_r of resolutions 2×2 , 8×8 , 32×32 and 128×128 . The intensity of each is shifted and rescaled to fit between 0 and 1 to be recognizable visually, however, the pixel values in the real P_r fall of approximately as r^{-1} .

396 We express the single, high resolution perturbation P as a sum
 397 of perturbations $P = \sum_{r \in \rho} \text{rescale}_R(P_r)$, where P_r is of the
 398 resolution $r \times r$ specified by a set of resolutions ρ , and the
 399 rescale_R function rescales and interpolates an image to the full
 400 resolution $R \times R$. When we jointly optimize the set of pertur-
 401 bations $\{P_r\}_{r \in \rho}$, we find that: a) the resulting attacked image
 402 $X + \sum_{r \in \rho} \text{rescale}_R(P_r)$ is much more human-interpretable,
 403 b) the attack follows a power distribution of natural images.

404 When attacking a classifier, we choose a target label t and
 405 optimize the cross-entropy loss of the predictions stemming
 406 from the perturbed image as if that class t were ground truth. To
 407 add to the robustness and therefore interpretability of the attack
 408 (as hypothesized in our *Interpretability-Robustness Hypothesis*),
 409 we add random jitter in the x-y plane and random pixel noise,
 410 and design the attack to work on a set of models.

411 An example of the multi-resolution sum is show in Figure 12.

412 There we use a simple Stochastic Gradient Descent [Robbins and Monro, 1951] optimization with the
 413 learning rate of 5×10^{-3} and a cosine decay schedule over 50 steps. We add a random pixel noise
 414 of 0.6 (out of 1), jitter in the x-y plane in the ± 5 range and a set of all perturbations from 1×1 to
 415 224×224 interpolated using bicubic interpolation [Keys, 1981]. In Figure 12 we see that despite
 416 the very limited expressiveness of the final layer class label, we can still recover images that look like
 417 the target class to a human. We also tested them using Gemini Advanced and GPT-4, asking what



Figure 11: An attack on vision language models. GPT-4 sees *Rick Astley from his famous "Never Gonna Give You Up" music video* tree. See Table 3 and 4 for details.

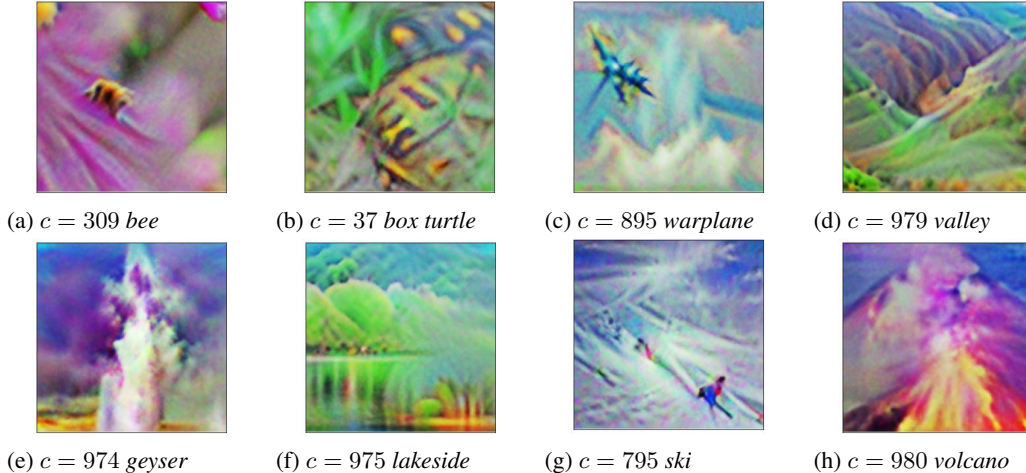


Figure 12: Examples of images generated as attacks on ImageNet-trained classifiers. These images were generated by minimizing the cross-entropy loss of seven pretrained classifiers with respect to the target ImageNet class. Spatial jitter in the ± 5 pixel range and pixel noise of standard deviation 0.6 were applied during SGD optimization with learning rate 5×10^{-3} over 50 steps with a cosine schedule. The perturbation was expressed as a sum of perturbations at all resolutions from 1×1 to 224×224 that were optimized at once.

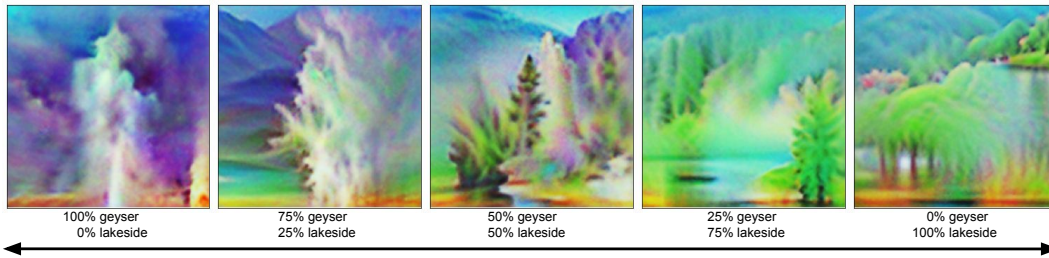


Figure 13: Optimizing towards a probability vector with a sliding scale between $c = 974$ *geyser* and $c = 975$ *lakeside*. Optimizing against pretrained classifiers generated semantically blended image of the two concepts.

418 the AI model sees in the picture, and got the right response in all 8 cases. To demonstrate that we
 419 can generate images beyond the original 1000 ImageNet classes, we experimented with setting the
 420 target label not as a one-hot vector, but rather with target probability p on class t_1 and $1 - p$ on t_2 .
 421 For classes $c = 974$ (*geyser*) and $c = 975$ (*lakeside*) we show, in Figure 13 that we get semantically
 422 meaningful combinations of the two concepts in the same image as we vary p from 0 to 1. $p = 1/2$
 423 gives us a *geyser* hiding beyond trees at a *lakeside*. This example demonstrates that in a limited way,
 424 classifiers can be used as controllable image generators.

425 G Multi-resolution attack on CLIP

426 The CLIP-style [Radford et al., 2021] models map an image I to an embedding vector $f_I : I \rightarrow v_I$
 427 and a text T to an embedding vector $f_T : T \rightarrow v_T$. The cosine between these two vectors corresponds
 428 to the semantic similarity of the image and the text, $\cos(v_I, v_T) = v_I \cdot v_T / (|v_I| |v_T|)$. This gives us
 429 score(I, T) that we can optimize.

430 Adversarial attacks on CLIP can be thought of as starting with a human-understandable image X_0 (or
 431 just a noise), and a target label text T^* , and optimizing for a perturbation P to the image that tries to
 432 increase the score($X_0 + P, T^*$) as much as possible. In general, finding such perturbations is easy,
 433 however, they end up looking very noise-like and non-interpretable. [Fort, 2021a,b].

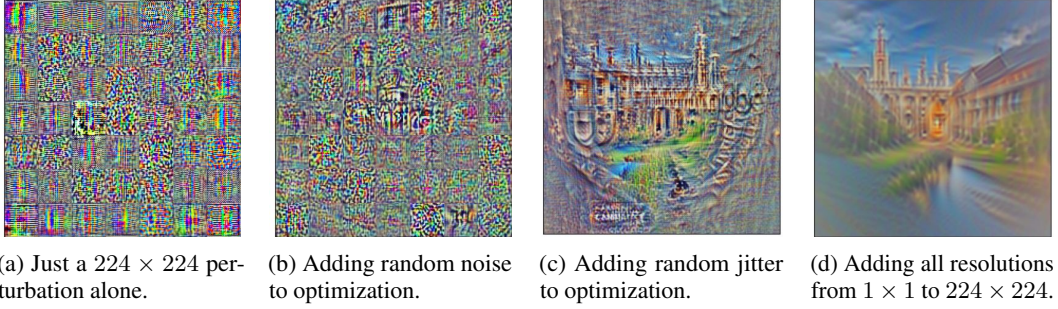


Figure 14: The effect of adding noise, jitter, and a full set of resolutions to an adversarial attack on CLIP towards the text '*a beautiful photo of the University of Cambridge, detailed*'. While using just a plain perturbation of the full resolution in Figure 14a, as is standard in the typical adversarial attack setup, we get a completely noise-like image. Adding random noise to the pixels during optimization leads to a glimpse of a structure, but still maintains a very noise-like pattern (Figure 14b). Adding random jitter in the x - y plane on top, we can already see interpretable shapes of *Cambridge* buildings in Figure 14c. Finally, adding perturbations of all resolutions, $1 \times 1, 2 \times 2, \dots, 224 \times 224$, we get a completely interpretable image as a result in Figure 14d.

434 If we again express $P = \text{rescale}_{224}(P_1) + \text{rescale}_{224}(P_2) + \dots + P_{224}$, where P_r is a resolution $r \times r$
 435 image perturbation, and optimize $\text{score}(X_0 + \text{rescale}_{224}(P_1) + \text{rescale}_{224}(P_2) + \dots + P_{224}, T^*)$
 436 by simultaneously updating all $\{P_r\}_r$, the resulting image $X_0 + \sum_{r \in [1, 224]} \text{rescale}_R(P_r)$ looks like
 437 the target text T^* to a human rather than being just a noisy pattern. Even though the optimizer could
 438 choose to act only on the full resolution perturbation P_{224} , it ends up optimizing all of them jointly
 439 instead, leading to a more natural looking image. To further help with natural-looking attacks, we
 440 introduce pixel noise and the x - y plane jitter, the effect of which is shown in Figure 14.

441 We use SGD at the learning rate of 5×10^{-3} for 300 steps with a cosine decay schedule to maximize
 442 the cosine between the text description and our perturbed image. We use the OpenCLIP models
 443 [Ilharco et al., 2021, Cherti et al., 2023] (an open-source replication of the CLIP model [Radford
 444 et al., 2021]). Examples of the resulting "adversarial attacks", starting with a blank image with 0.5
 445 in its RGB channels, and optimizing towards the embedding of specific texts such as "*a photo of*
 446 *Cambridge UK, detailed*", and "*a photo of a sailing boat on a rough sea*" are shown in Figure 16. The
 447 image spectra are shown in Figure 8, displaying a very natural-image-like distribution of powers. The
 448 resulting images look very human-interpretable.

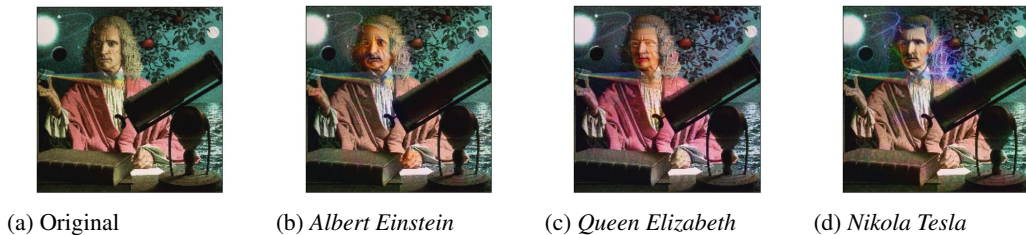


Figure 15: Starting with an image of Isaac Newton and optimizing a multi-resolution perturbation towards text embeddings of *Albert Einstein*, *Queen Elizabeth* and *Nikola Tesla* leads to a change in the face of the person depicted. This demonstrates how semantically well-targeted such multi-resolution attacks are. All 4 images are recognizable as the target person to humans as well as GPT-4o and Gemini Advanced.

449 Starting from a painting of Isaac Newton and optimizing towards the embeddings of "*Albert Ein-*
 450 *stein*", "*Queen Elizabeth*" and "*Nikola Tesla*", we show that the attack is very semantically targeted,
 451 effectively just changing the facial features of Isaac Newton towards the desired person. This is
 452 shown in Figure 15. This is exactly what we would ideally like adversarial attacks to be – when
 453 changing the content of what the model sees, the same change should apply to a human. We use a
 454 similar method to craft transferable attacks (see Figure 11 for an example) against commercial, closed
 455 source vision language models (GPT-4, Gemini Advanced, Claude 3 and Bing AI) in Table 3, in

456 which a *turtle* turns into a *cannon*, and in Table 4, where *Stephen Hawking* turns into the music video
 457 *Never Gonna Give You Up* by *Rick Astley*. The attacks also transfer to Google Lens, demonstrating
 458 that the multi-resolution prior also serves as a good *transfer* prior and forms an early version of a
 459 transferable image vision language model jailbreak. This is a constructive proof to the contrary of the
 460 non-transferability results in [Schaeffer et al. \[2024\]](#).

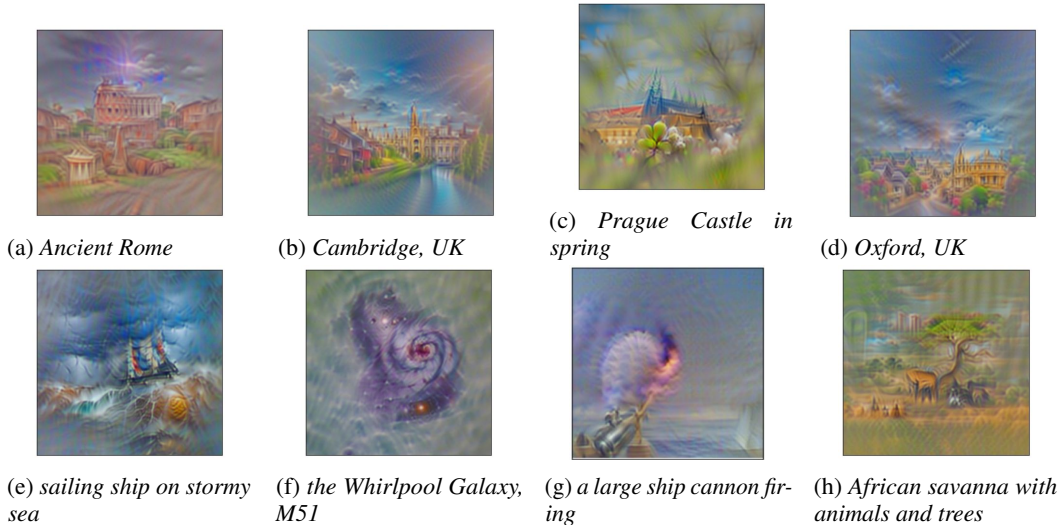


Figure 16: Examples of images generated with the multi-resolution prior, jitter and noise with the OpenCLIP models. The text whose embedding the image optimizes to approach is of the form ‘A beautiful photo of [X], detailed’ for different values of [X].

461 H Additional details on attack transfer between layers

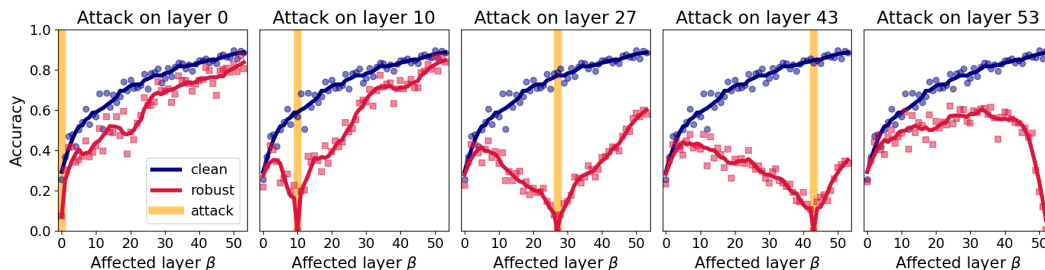


Figure 17: Transfer of adversarial attacks ($L_\infty = 8/255$, 512 attacks) against the activations of layer α on the accuracy of layer β for $\alpha = 0, 10, 27, 43, 53$ on ImageNet-1k pretrained ResNet152 finetuned on CIFAR-10 via trained linear probes. Each panel shows the effect of designing a pixel-level attack to confuse the linear probe at a particular layer. The blue curve is the test accuracy on the unperturbed data, and the red line shows the accuracy on the attacked images. The accuracy drops to 0 at the layer that is directly attacked (marked in orange), showing a successful attack. The effect is localized: attacking early layers mainly affects early layer predictions, middle layer attacks primarily affect middle layers, and likewise attacks on the final layers (the standard regime) primarily influence late layer performance. For more details, see Figure 19.

462 **I Transfer to massive commercial models**

463 In Table 3 we show the results of asking "What do you see in this photo?" and adding the relevant
 464 picture to four different, publicly available commercial AI models: GPT-4¹, Bing Copilot², Claude
 465 3 Opus³ and Gemini Advanced⁴. We find that, with an exception of Gemini Advanced, even a
 466 $L_\infty = 30/255$ attack generated in approximately 1 minute on a single A100 GPU (implying a cost
 467 at most in cents) fools these large models into seeing a *cannon* instead of a *turtle*. The attack also
 468 transfers to Google Lens.

	Original	$L_\infty = 20/255$	$L_\infty = 30/255$	$L_\infty = 40/255$	$L_\infty = 70/255$	$L_\infty = 100/255$
						
GPT-4	sea turtle swimming	turtle swimming in water	cannon mounted on stone base, firing	cannon with a notably ornate and rusted appearance	cannon mounted on a brick platform	stylized or artistically rendered depiction of a cannon
Bing Copilot	sea turtle gracefully swimming	sea turtle gracefully swimming	a cannon mounted on a stone base	cannon with a wheel, mounted on a stone base	old cannon mounted on a brick platform	color-saturated cannon mounted on wheels
Claude 3 Opus	sea turtle swimming in clear, turquoise water	sea turtle swimming underwater	old cannon submerged underwater	old decorative cannon sitting on a stone or concrete platform	old naval cannon set on a stone or brick platform	artistic painting or illustration of an old cannon
Gemini Advanced	sea turtle swimming underwater	sea turtle swimming underwater	sea turtle swimming	sea turtle swimming in a pool	cannon being fired by a turtle wearing a red bucket	artistic interpretation of a cannon firing

Table 3: Multi-resolution adversarial attacks of increasing L_∞ using OpenCLIP on an image of a sea turtle towards the text "a cannon" tested on GPT-4, Bing Copilot (Balanced), Claude 3 Sonnet and Gemini Advanced. All models we tested the images on were publicly available. The conversation included a single message "What do you see in this photo?" and an image. We chose the most relevant parts of the response.

469 Figure 18 compares attacks on robust and brittle models.

470 **J Attack transfer between layers**

471 This setup also allows us not only to investigate what the intermediate classification decision would
 472 be for an adversarially modified image X' that confuses the network's final layer classifier, but also to
 473 generally ask what the effect of confusing the classifier at layer α would do to the logits at a layer β .
 474 The results are shown in Figure 17 for 6 selected layers to attack, and the full attack layer \times read-out
 475 layer is show in Figure 19.

¹chatgpt.com
²bing.com/chat
³claude.ai/
⁴gemini.google.com

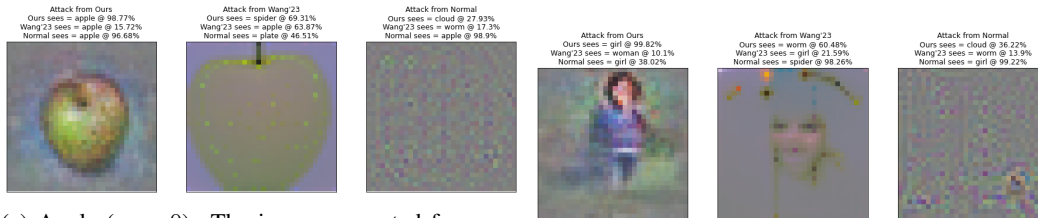
	Original	$L_\infty = 20/255$	$L_\infty = 30/255$	$L_\infty = 40/255$	$L_\infty = 70/255$	$L_\infty = 100/255$
						
GPT-4	Stephen Hawking	Stephen Hawking	Never Gonna Give You Up	Never Gonna Give You Up	Never Gonna Give You Up	singer or performer, possibly Rick Astley
Bing Copilot	individual sitting in a wheelchair	individual sitting on a bench	individual sitting down, holding a microphone, singing	person seated, holding a musical instrument	two individuals in an indoor setting	person in front of a microphone, singing
Claude 3 Opus	elderly man in a wheelchair	man in a wheelchair, smiling	young man with blonde hair, vintage-style microphone, singing	young man with blond hair, 1980s pop music	music video, 1980s, singer	music video, 1980s fashion
Gemini Advanced	Refused to answer.	Refused to answer.	Refused to answer.	Refused to answer.	Refused to answer.	Refused to answer.

Table 4: Multi-resolution adversarial attacks of increasing L_∞ using OpenCLIP on an image of *Stephen Hawking* towards the embedding of an image from the famous *Rick Astley's* song *Never Gonna Give You Up* from the 1980s tested on GPT-4, Bing Copilot (Balanced), Claude 3 Sonnet and Gemini Advanced. All models we tested the images on were publicly available. The conversation included a single message "What do you see in this photo?" and an image. We chose the most relevant part of the response. Unfortunately, Gemini refused to answer, likely due to the presence of a human face in the photo.

476 We find that attacks designed to confuse early layers of a network do not confuse its middle and
477 late layers. Attacks designed to fool middle layers do not fool early nor late layers, and attacks
478 designed to fool late layers do not confuse early or middle layers. In short, there seems to be roughly
479 a 3-way split: early layers, middle layers, and late layers. Attacks designed to affect one of these do
480 not generically generalize to others. We call this effect the *adversarial layer de-correlation*. This
481 de-correlation allows us to create a *self-ensemble* from a single model, aggregating the predictions
482 resulting from intermediate layer activations. To make sure that the ensemble is robust, we use the
483 *CrossMax* method described in Section 2.2 and Algorithm 1. While ensembling multiple equivalent
484 models, we did not have to care about their different quality, however, here early layers are typically
485 less accurate than late layers, as shown in Figure 3.

486 In Figure 24 we show the self-ensemble robustness under adversarial attacks of different strength
487 for an ImageNet pretrained ResNet152 and ViT-B/16, with linear heads at each layer separately
488 finetuned on CIFAR-10. The aggregation method in Algorithm 1 provides non-zero robust accuracy
489 for attacks of even $L_\infty = 5/255$, while standard ensembling using mean logits as well as just the
490 last layer prediction loses robust accuracy around 3/255. This is an early indication that CrossMax
491 self-ensembling can actively use the decorrelation of intermediate layer adversarial susceptibilities
492 for an active, white-box defense.

493 K Visualizing attacks on multi-resolution models



(a) Apple ($c = 0$): The image generated from our model looks like an *apple* to itself, the Wang et al. [2023] robust model, and a brittle ResNet152 alike. The attacks against Wang et al. [2023] and standard ResNet152, on the other hand, convince only themselves.

(b) Girl ($c = 35$): The image generated from our model looks like a *girl* to itself, a brittle ResNet152 alike, and as a *woman* to the Wang et al. [2023] robust model. The attacks against them, on the other hand, convince only themselves.

Figure 18: Examples of adversarial attacks on our multi-resolution ResNet152 finetuned on CIFAR-100 (left), the previous best model on CIFAR-100 $L_\infty = 8/255$ on RubustBench from Wang et al. [2023] (middle), and standard ResNet152 finetuned on CIFAR-100. The attacks are generated by starting from a uniform image (128,128,128) and using gradient descent of the cross-entropy loss with SGD at $\eta = 1$ for 400 steps towards the target label. The prediction results for each of the models are shown above the images.

494 Figure 22 shows 6 examples of successfully attacked CIFAR-100 test set images for an ensemble of 3 self-ensemble models – our most adversarially robust model. When looking at the misclassifications caused, we can easily see human-plausible ways in which the attacked image can be misconstrued as the most probable target class. For example, a crab with a body resembling a mushroom cap gets a foot of a mushroom added by the attack, causing a misclassification as 40% mushroom from a 90% crab. A blurry picture of a sting ray gets 3D-like shading added by the attack, making it look mouse-like and being classified as 30% shrew from a 90% ray. Overall, we see that the changes that are induced by the attacker seem to have a human-understandable explanation. Figure 20 shows an example of a successful $L_\infty = 64/255$ (much larger than the standard 8/255 perturbations) RobustBench AutoAttack on a test image of a *bicycle* converting it, in a human-interpretable way, to a *snake* by re-purposing parts of the bicycle frame as the snake body.

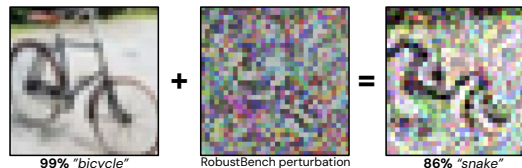


Figure 20: An example of a $L_\infty = 64/255$ RobustBench AutoAttack on our model, changing a *bicycle* into a *snake* in an interpretable way.

511 We are also very interested in the existence of adversarial attacks on the human visual system and we believe that our work should be an update against their likelihood. We use biologically inspired methods (multiple resolutions, jitter, noise) that work as a defense against a white-box attacker. When flipped around, the same ideas generate human-interpretable images. The intermediate layer representations could also be viewed as using shallower circuits in the brain, and their partial robustness might suggest the same in humans. Given that moving closer (in a very rudimentary way) to the human visual system in these regards gave us both a practical defense and an image generator, we believe that we should update against adversarial vulnerability of humans.

519 L Additional experiments for CrossMax

520 To demonstrate experimentally different characteristics of prediction aggregation among several classifiers, we trained 10 ResNet18 models, starting from an ImageNet pretrained model, changing their final linear layer to output 10 classes of CIFAR-10. We then used the first 2 attacks of the RobustBench AutoAttack suite (APGD-T and APGD-CE; introduced by Croce and Hein [2020] as particularly strong attack methods) and evaluated the robustness of our ensemble of 10 models under adversarial attacks of different L_∞ strength. The results are shown in Figure 25.

526 The aggregation methods we show are 1) our CrossMax (Algorithm 1) (using *median* since the 10 models are expected to be equally good), 2) a standard logit mean over models, 3) median over

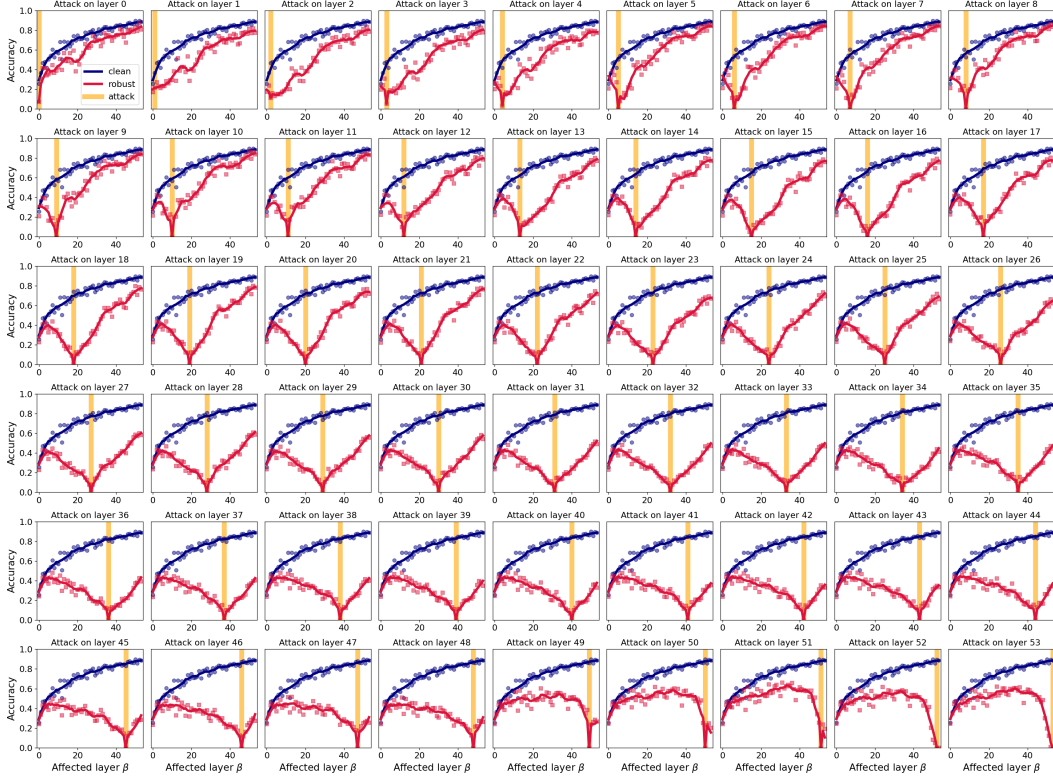


Figure 19: Attack transfer between layers of the ResNet154 model pre-trained on ImageNet-1k. The individual linear heads were finetuned on CIFAR-10 on top of the frozen model.

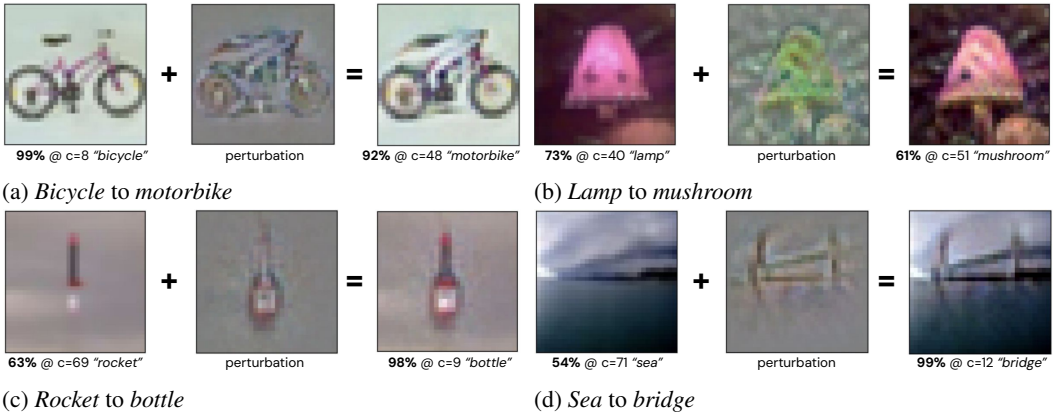


Figure 21: Additional examples of an adversarial attack on an image towards a target label. We use simple gradient steps with respect to our multi-resolution ResNet152 finetuned on CIFAR-100. The resulting attacks use the underlying features of the original image and make semantically meaningful, human-interpretable changes to it. Additional examples available in Figure 5.

528 models, and 4) the performance of the individual models themselves. While an ensemble of 10
 529 models, either aggregated with a mean or median, is more robust than individual models at all
 530 attack strengths, it nonetheless loses robust accuracy very fast with the attack strength L_∞ and at the
 531 standard level of $L_\infty = 8/255$ it drops to $\approx 0\%$. Our *CrossMax* in Algorithm 1 provides > 0 robust
 532 accuracy even to $10/255$ attack strengths, and for $8/255$ gives a 17-fold higher robust accuracy than
 533 just plain mean or median. We use this aggregation for intermediate layer predictions (changing
 534 *median* to *top₃*) as well and see similar, transferable gains. We call this setup a *self-ensemble*.



Figure 22: Examples of successfully attacked CIFAR-100 images for an ensemble of self-ensembles – our most robust model. We can see human-plausible ways in which the attack changes the perceived class. For example, the skyscraper has a texture added to it to make it look tree-like.

Algorithm 1 CrossMax = An Ensembling Algorithm with Improved Adversarial Robustness

Require: Logits Z of shape $[B, N, C]$, where B is the batch size, N the number of predictors, and C the number of classes

Ensure: Aggregated logits

- 1: $\hat{Z} \leftarrow Z - \max(Z, \text{axis} = 2)$ {Subtract the max per-predictor over classes to prevent any predictor from dominating}
 - 2: $\tilde{Z} \leftarrow \hat{Z} - \max(\hat{Z}, \text{axis} = 1)$ {Subtract the per-class max over predictors to prevent any class from dominating}
 - 3: $Y \leftarrow \text{median}(\tilde{Z}, \text{axis} = 1)$ {Choose the median (or k^{th} highest for self-ensemble) logit per class}
 - 4: **return** Y
-

535 As an ablation, we tested variants of the *CrossMax* method. There are two normalization steps: A)
 536 subtracting the per-predictor max, and B) subtracting the per-class max. We exhaustively experiment
 537 with all combinations, meaning $\{_, A, B, AB, BA\}$, (robust accuracies at 4/255 are $\{4, 4, 0, 22, 0\}\%$)
 538 and find that performing *A* and then *B*, as in Algorithm 1, is by far the most robust method. We
 539 perform a similar ablation for a robust, multi-resolution self-ensemble model in Table 5 and reach
 540 the same verdict, in addition to confirming that the algorithm is very likely not accidentally masking
 541 gradients.

Aggregation fn	topk ₂					mean				
	_	A	B	BA	AB	_	A	B	BA	AB
Test acc	57.08	59.86	0.82	1.27	58.92	60.31	59.89	1.1	1.05	57.23
Adv acc	46.88	46.88	1.56	0.00	57.81	40.62	48.44	0.00	0.00	39.06

Table 5: CrossMax algorithm ablation. The Algorithm 1 contains two subtraction steps: A = the per-predictor max subtraction, and B = the per-class max subtraction. This Table shows the robust accuracies of a self-ensemble model on CIFAR-100 trained with light adversarial training, whose intermediate layer predictions were aggregated using different combinations and orders of the two steps. We also look at the effect of using the final topk₂ aggregation vs just using a standard mean. The best result is obtained by the Algorithm 1, however, we see that not using the topk does not lead to a critical loss of robustness as might be expected if there were accidental gradient masking happening.

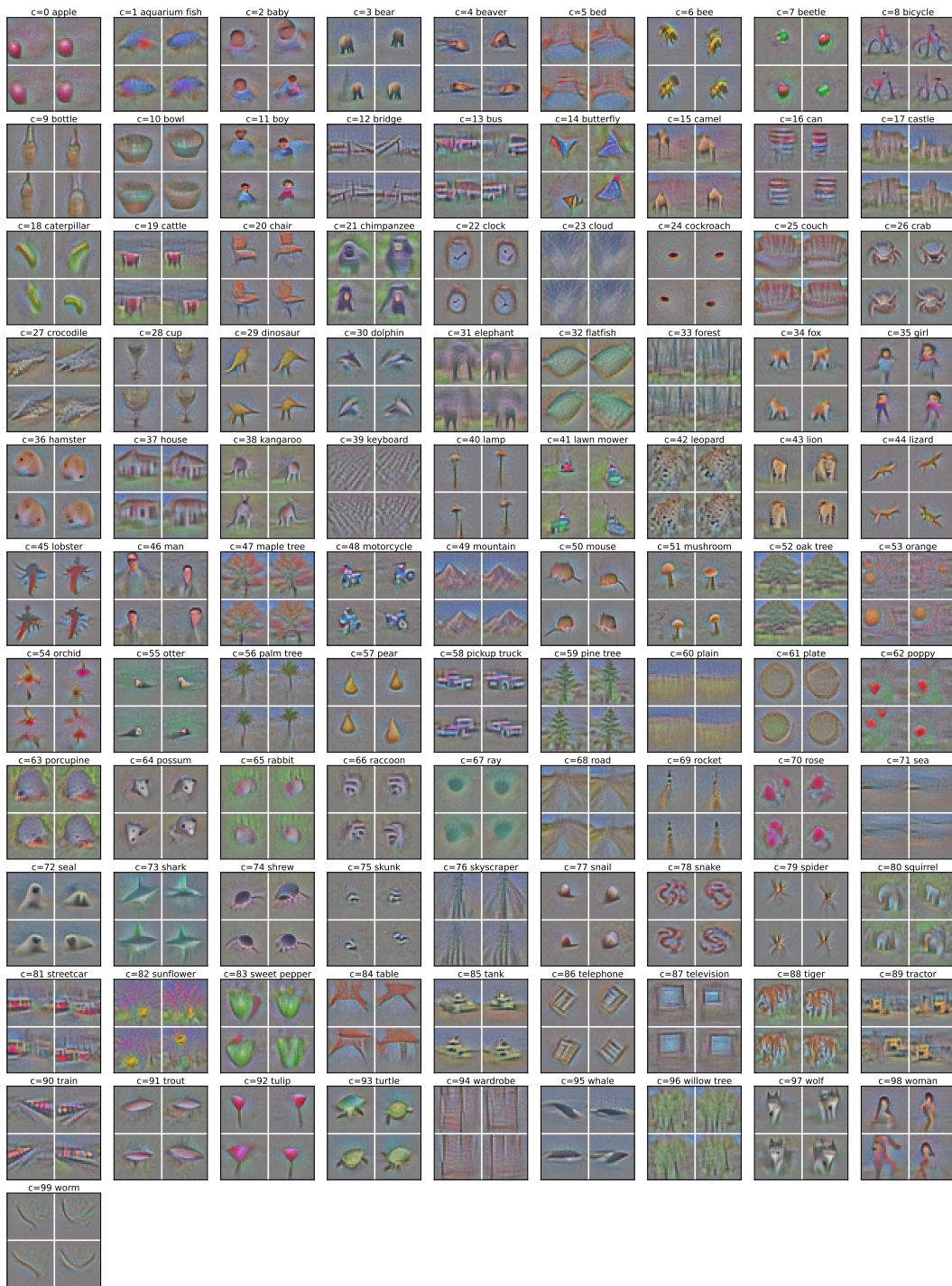
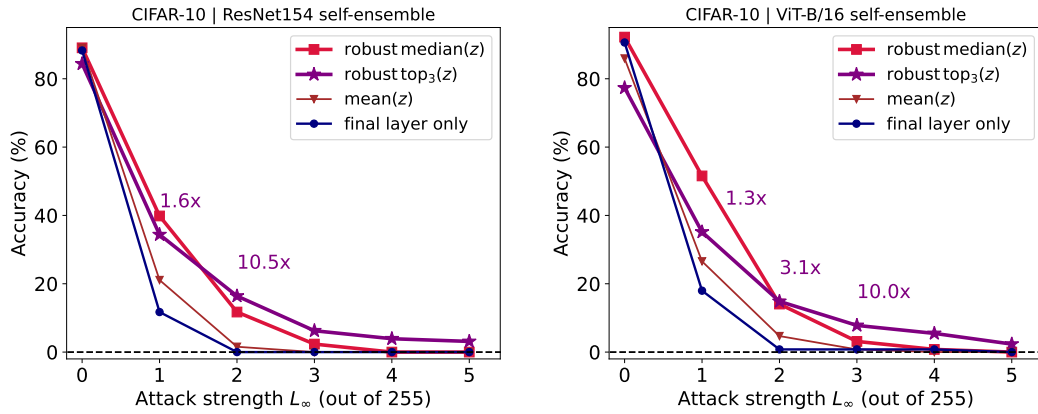


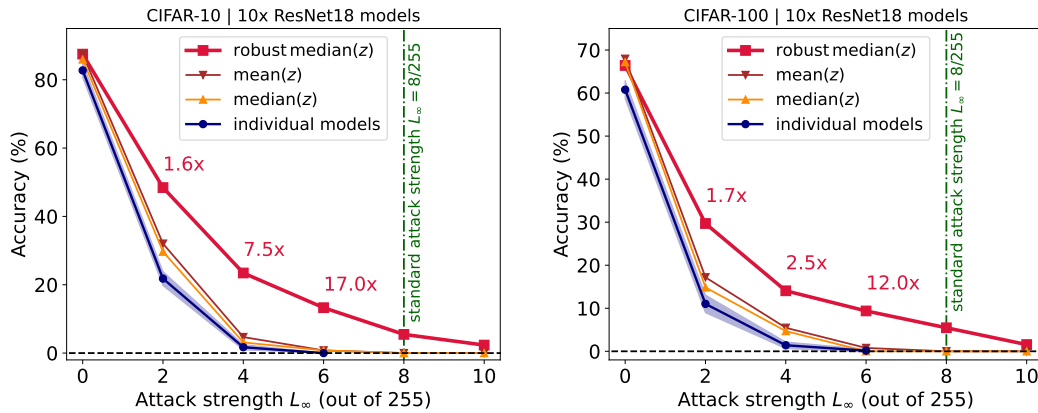
Figure 23: Examples of optimizing towards all 100 CIFAR-100 classes against our multi-resolution ResNet152 model, 4 examples for each. We use 400 simple gradient steps at learning rate $\eta = 1$ with SGD with respect to the model, starting from all grey pixels (128,128,128). The resulting attacks are easily recognizable as the target class to a human.



(a) ResNet154 self-ensemble on CIFAR-10

(b) ViT-B/16 self-ensemble on CIFAR-10

Figure 24: The robust accuracy of different types of self-ensembles of ResNet152 and ViT-B/16 with linear heads finetuned on CIFAR-10 under increasing L_∞ attack strength.



(a) CIFAR-10

(b) CIFAR-100

Figure 25: The robust accuracy of different types of ensembles of 10 ResNet18 models under increasing L_∞ attack strength. Our robust median ensemble, *CrossMax*, gives very non-trivial adversarial accuracy gains to ensembles of individually brittle models. For $L_\infty = 6/255$, its CIFAR-10 robust accuracy is 17-fold larger than standard ensembling, and for CIFAR-100 the factor is 12.