

---

# On the Robustness of Neural Collapse and the Neural Collapse of Robustness

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Neural Collapse refers to the curious phenomenon in the end of training of a neural  
2 network, where feature vectors and classification weights converge to a very simple  
3 geometrical arrangement (a simplex). While it has been observed empirically in  
4 various cases and has been theoretically motivated, its connection with crucial  
5 properties of neural networks, like their generalization and robustness, remains  
6 unclear. In this work, we study the stability properties of these simplices. We  
7 find that the simplex structure disappears under small adversarial attacks, and  
8 that perturbed examples "leap" between simplex vertices. We further analyze  
9 the geometry of networks that are optimized to be robust against adversarial  
10 perturbations of the input, and find that Neural Collapse is a pervasive phenomenon  
11 in these cases as well, with clean and perturbed representations forming aligned  
12 simplices, and giving rise to a robust simple nearest-neighbor classifier. By studying  
13 the propagation of the amount of collapse inside the network, we identify novel  
14 properties of both robust and non-robust machine learning models, and show that  
15 earlier, unlike later layers maintain reliable simplices on perturbed data.

## 16 1 Introduction

17 Reinforcing arguments about the simplicity of neural networks found by stochastic gradient descent  
18 in classification settings, Pappayan et al. [2020] made the surprising empirical observation that both  
19 the feature representations in the penultimate layer (grouped by their corresponding class) and the  
20 weights of the final layer form a *simplex equiangular tight frame* (ETF) with  $C$  vertices, where  
21  $C$  is the number of classes. Curiously, such a geometric arrangement becomes more pronounced  
22 well-beyond the point of (effectively) zero loss on the training data, motivating the common tendency  
23 of practitioners to optimize a network for as long as the computational budget allows. The collection  
24 of these empirical phenomena was termed *Neural Collapse*.

25 While the results of [Pappayan et al., 2020] fueled much research in the field, many questions remain  
26 regarding the connection of Neural Collapse with properties like generalization and robustness of  
27 Neural Networks. In particular with regards to *adversarial robustness*, the ability of a model to  
28 withstand adversarial modifications of the input without effective drops in performance, it has been  
29 originally claimed that the instantiation of Neural Collapse has positive effect on defending against  
30 adversarial attacks [Pappayan et al., 2020, Han et al., 2022]. However, this seems to at least superficially  
31 contradict the fact that neural networks are not a priori adversarially robust [Szegedy et al., 2014,  
32 Carlini and Wagner, 2017].

33 In this paper, we thoroughly study the stability properties of the simplices under adversarial attacks  
34 and then investigate whether Neural Collapse happens and whether it is necessary for adversarially  
35 robust models. In particular, our contributions and findings, partially illustrated in Figure 1, are:

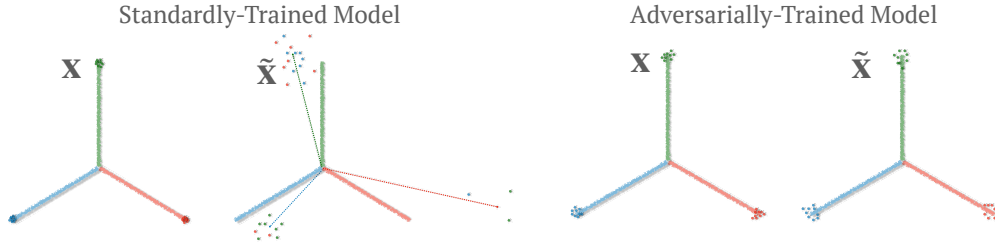


Figure 1: Visualisation of our findings. Sticks represent class-means. Small dots correspond to the representation of an individual datum, and the color represents the ground-truth label. **Left to Right:** clean representations with standardly-trained (ST) networks; perturbed representations with ST networks; clean representations with adversarially-trained (AT) networks; perturbed representations with AT networks. With ST nets, the adversarial perturbations push the representation to “leap” towards another cluster with slight angular deviation. AT makes the simplex resilient to such adversarial attacks, with higher and intra-class variance.

- 36 • **Is NC robust?** We initiate the study of the neural collapse phenomenon in the context of adversarial  
37 robustness, both for standardly trained networks under adversarial attacks and for adversarially  
38 trained robust networks to investigate the stability and prevalence of the NC phenomenon. Our  
39 work exposes considerable additional fundamental, and we think, surprising, geometrical structure:
- 40 • **No!** For standardly trained networks we find that small, imperceptible adversarial perturbations  
41 of the training data remove any simplicial structure at the representation layer: neither variance  
42 collapse nor simplex representations appear under standard metrics. Further analysis through class-  
43 targeted attacks that preserve class-balance shows a “cluster-leaping” phenomenon: representations  
44 of adversarially perturbed data jump to the (angular) vicinity of the original class means.
- 45 • **Yes for AT networks! Two identical simplices emerge.** Adversarially trained, robust, networks  
46 exhibit a simplex structure both on original clean and adversarially perturbed data, albeit of higher  
47 variance. These two simplices turn out to be the same. We find that the simple nearest-neighbor  
48 classifiers extracted from such models also exhibit robustness.
- 49 • **Early layers are more robust.** Analyzing NC metrics in the representations of the inner layers,  
50 we observe that initial layers exhibit a higher degree of collapse on adversarial data. The resulting  
51 simplices, when used for Nearest Neighbour clustering, give surprisingly robust classifiers. This  
52 phenomenon disappears in later layers.

## 53 2 Background

54 Papayan et al. [2020] demonstrated the prevalence of NC on networks optimized by SGD, by tracing  
55 the following quantities (please refer to Appendix B.2 for exact definitions):

56 **(NC1) Variability collapse:** For all classes, the within-class variation of the last layer representations  
57 collapses to zero.

58 **(NC2, Equiangular):** Class-Means converge to equal, maximal pairwise angles.

59 **(NC2, Equinorm):** Class-Means converge to equal length.

60 **(NC3) Convergence to self-duality:** The linear classifier and the class-means converge to each other  
61 (after rescaling).

62 **(NC4) Simplification to Nearest Class Center (NCC) classifier:** The prediction of the network is  
63 equivalent to that of the NCC classifier formed by the (non-centered) class-means.

64 We adopt the natural extensions of the above metrics for adversarially trained models, and further  
65 study some new quantities, relevant to our analysis, which are defined and explained in Appendix  
66 B.2.

### 67 3 Experiments

68 In this section, we present our main experimental results measuring neural collapse in standardly  
 69 (ST) (with SGD) and adversarially trained (AT) models [Madry et al., 2018]. We consider image  
 70 classification tasks on CIFAR-10 and CIFAR-100 and we train two large convolutional networks, a  
 71 standard VGG and a Pre-Activation ResNet18, from random initializations. We launch 3 independent  
 72 runs and report the mean and standard deviation throughout our paper. Further results for varying  
 73 choices of hyperparameters can be found in the Appendix.

74 **Remark:** When collecting feature representations for adversarially perturbed data, we always  
 75 compute the *current epoch's* perturbations.

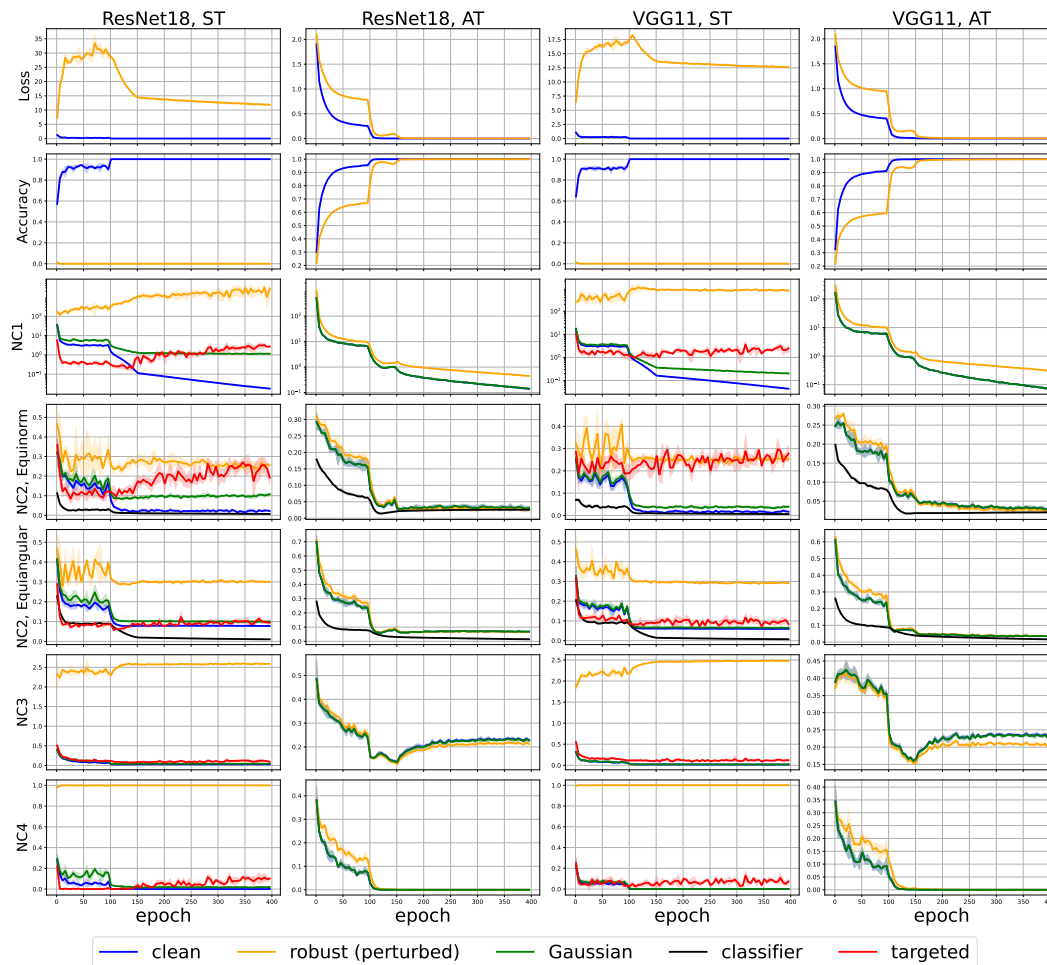


Figure 2: Accuracy, Loss, and NC evolution for standardly (ST) and adversarially (AT) trained VGG and ResNet. For AT models, clean and Gaussian curves coincide. Setting: CIFAR-10,  $\ell_\infty$  adversary.

#### 76 3.1 Standardly trained neural nets

77 The first and third column of Figure 2 show the evolution of the NC quantities as described in Section  
 78 2 for standardly trained models. We use both adversarially perturbed and Gaussian reference data to  
 79 study the stability of the original simplices. As expected, NC metrics converge on the clean training  
 80 data. Neural Collapse is slightly attenuated on Gaussian reference data, but disappears strikingly  
 81 for adversarially perturbed data, suggesting that the simplex formed by clean training data is robust  
 82 against random perturbations, but fragile to adversarial attacks. The results certainly corroborate the  
 83 conclusion that the representation class-means of perturbed points with ground-truth label  $c$  do not  
 84 form any geometrically-meaningful structure at all.

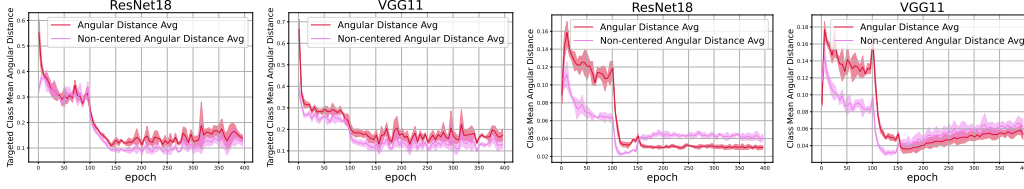


Figure 3: Angular distance. *Left and Inner Left*: Average between targeted attack class-means and clean class-means on **ST** network. *Inner Right and Right*: Average between perturbed class-means and clean class-means on **AT** network. Setting: CIFAR-10,  $\ell_\infty$  adversary.

85 Figure 3 (left) shows Simplex Similarity and non-centered Angular Distance of the simplices formed  
 86 by targeted adversarial examples and by clean examples as described in Appendix B.2. These results  
 87 give us a full glimpse of how standardly trained networks are non-robust and fail under adversarial  
 88 attacks: adversarial perturbations break the simplex ETF by “leaping” the representation from one  
 89 class-mean to another, forming a norm-imbalanced less concentrated structure around the original  
 90 simplex.

### 91 3.2 Neural Collapse during Adversarial Training

92 We train neural nets adversarially to full convergence with perfect clean and robust training accuracy  
 93 and measure NC metrics for clean and perturbed (epoch-wise) training data in Figure 2 (columns 2  
 94 and 4). Interestingly, we find that Neural Collapse *qualitatively* occurs in this setting as well, both for  
 95 clean and perturbed data, and two simplices emerge. Notice, however, that the extent of variability  
 96 collapse (NC1) on the perturbed points is smaller than on the “clean” data or the Gaussian noise  
 97 benchmark, indicating that clean examples are more concentrated around the vertices. To understand  
 98 the relative positioning of the two simplices, we investigate the Simplex Similarity and Angular  
 99 Distance between non-centered class-means in Figure 3 (right). The vanishing distance suggests  
 100 these two simplices are exactly the same. These results suggest that Adversarial Training nudges  
 101 the network to learn simple representational structures (namely, a simplex ETF) not only on clean  
 102 examples but also on perturbed examples to achieve robustness against adversarial perturbations.  
 103 Equivalently, the simplices induced by robust networks are *not fragile* anymore, but *resilient*. Note  
 104 also that NC4 results imply that there is a simple nearest-neighbor classifier that is robust against  
 105 adversarial perturbations generated from the network.

106 Curiously, this is not the case for all training algorithms that produce robust models. In particular,  
 107 a state-of-the-art algorithm that aims to balance clean and robust accuracy, TRADES [Zhang et al.,  
 108 2019], shows fundamentally different behavior (see Figure 4 in the Appendix). Even though both  
 109 terms of the loss (see Equation 4) are driven to zero, we do not observe Neural Collapse; the amount of  
 110 collapse is roughly one order of magnitude larger than for vanilla AT, and the feature representations  
 111 do not approach the ETF formation, even well past the onset of the terminal phase. *We view this as*  
 112 *evidence that the prevalence of Neural Collapse is not necessary for robust classification.*

### 113 3.3 Layerwise Analysis

114 Furthermore, the Appendix contains our detailed layerwise analysis on both ST and AT models.  
 115 We observe that initial layers exhibit a higher degree of collapse on adversarial data, while NCC  
 116 classifiers defined on intermediate layers show surprising robustness, even if the whole model fails to  
 117 do so.

## 118 4 Conclusion

119 Neural Collapse is an interesting phenomenon displayed by Neural Networks used in classification  
 120 tasks. We empirically studied and quantified the sensitivity of this geometric arrangement to input  
 121 perturbations, and, further, displayed that Neural Collapse can appear (but not always does!) in  
 122 Neural Networks trained to be robust. We conclude that Neural Collapse is prevalent in many deep  
 123 learning settings, including adversarially trained networks, though it does not seem to be necessary  
 124 for robustness.

125 **References**

- 126 Ido Ben-Shaul and Shai Dekel. Nearest class-center simplification through intermediate layers. *CoRR*,  
127 abs/2201.08924, 2022.
- 128 Nicholas Carlini and David A. Wagner. Towards evaluating the robustness of neural networks. In  
129 *2017 IEEE Symposium on Security and Privacy, SP 2017, San Jose, CA, USA, May 22-26, 2017*,  
130 pages 39–57. IEEE Computer Society, 2017.
- 131 Corinna Cortes and Vladimir Vapnik. Support-vector networks. *Mach. Learn.*, 20(3):273–297, 1995.
- 132 Francesco Croce, Maksym Andriushchenko, Vikash Sehwal, Edoardo Debenedetti, Nicolas Flam-  
133 marion, Mung Chiang, Prateek Mittal, and Matthias Hein. Robustbench: a standardized adversarial  
134 robustness benchmark. In *Thirty-fifth Conference on Neural Information Processing Systems*  
135 *Datasets and Benchmarks Track (Round 2)*, 2021.
- 136 Weinan E and Stephan Wojtowytsch. On the emergence of simplex symmetry in the final and penulti-  
137 mate layers of neural network classifiers. In Joan Bruna, Jan Hesthaven, and Lenka Zdeborova,  
138 editors, *Proceedings of the 2nd Mathematical and Scientific Machine Learning Conference*, volume  
139 145 of *Proceedings of Machine Learning Research*, pages 270–290. PMLR, 16–19 Aug 2022.
- 140 Cong Fang, Hangfeng He, Qi Long, and Weijie J. Su. Exploring deep neural networks via layer-peeled  
141 model: Minority collapse in imbalanced training. *Proceedings of the National Academy of Sciences*  
142 *of the United States of America*, 118, 2021.
- 143 Tomer Galanti, András György, and Marcus Hutter. On the role of neural collapse in transfer learning.  
144 *CoRR*, abs/2112.15121, 2021.
- 145 Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial  
146 examples. In Yoshua Bengio and Yann LeCun, editors, *3rd International Conference on Learning*  
147 *Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*,  
148 2015.
- 149 X.Y. Han, Vardan Papyan, and David L. Donoho. Neural collapse under MSE loss: Proximity to and  
150 dynamics on the central path. In *International Conference on Learning Representations*, 2022.  
151 URL [https://openreview.net/forum?id=w1UbvWH\\_R3](https://openreview.net/forum?id=w1UbvWH_R3).
- 152 Hangfeng He and Weijie J. Su. A law of data separation in deep learning. 2022.
- 153 Like Hui, Mikhail Belkin, and Preetum Nakkiran. Limitations of neural collapse for understanding  
154 generalization in deep learning. *CoRR*, abs/2202.08384, 2022.
- 155 Wenlong Ji, Yiping Lu, Yiliang Zhang, Zhun Deng, and Weijie J Su. An unconstrained layer-peeled  
156 perspective on neural collapse. In *International Conference on Learning Representations*, 2022.  
157 URL <https://openreview.net/forum?id=WZ3yjh8coDg>.
- 158 Alexey Kurakin, Ian J. Goodfellow, and Samy Bengio. Adversarial examples in the physical world.  
159 In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April*  
160 *24-26, 2017, Workshop Track Proceedings*, 2017.
- 161 Xiao Li, Sheng Liu, Jinxin Zhou, Xinyu Lu, Carlos Fernandez-Granda, Zhihui Zhu, and Qing Qu.  
162 Principled and efficient transfer learning of deep models via neural collapse. *arXiv preprint*  
163 *arXiv:2212.12206*, 2022.
- 164 Yan Li, Ethan X. Fang, Huan Xu, and Tuo Zhao. Implicit bias of gradient descent based adversarial  
165 training on separable data. In *8th International Conference on Learning Representations, ICLR*  
166 *2020, Addis Ababa, Ethiopia, April 26-30, 2020*.
- 167 Bochen Lv and Zhanxing Zhu. Implicit bias of adversarial training for deep neural networks. In  
168 *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April*  
169 *25-29, 2022*. OpenReview.net, 2022.
- 170 Kaifeng Lyu and Jian Li. Gradient descent maximizes the margin of homogeneous neural networks.  
171 In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia,*  
172 *April 26-30, 2020*. OpenReview.net, 2020.



- 173 Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.  
174 Towards Deep Learning Models Resistant to Adversarial Attacks. In *International Conference on*  
175 *Learning Representations*, 2018.
- 176 Dustin G. Mixon, Hans Parshall, and Jianzong Pi. Neural collapse with unconstrained features.  
177 *Sampling Theory, Signal Processing, and Data Analysis*, 20:11, 2022.
- 178 Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, Omar Fawzi, and Pascal Frossard. Universal  
179 adversarial perturbations. In *Proceedings of the IEEE conference on computer vision and pattern*  
180 *recognition*, pages 1765–1773, 2017.
- 181 Behnam Neyshabur, Ryota Tomioka, and Nathan Srebro. In search of the real inductive bias: On the  
182 role of implicit regularization in deep learning. In Yoshua Bengio and Yann LeCun, editors, *3rd*  
183 *International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9,*  
184 *2015, Workshop Track Proceedings*, 2015. URL <http://arxiv.org/abs/1412.6614>.
- 185 Nicolas Papernot, Patrick D. McDaniel, Ian J. Goodfellow, Somesh Jha, Z. Berkay Celik, and  
186 Ananthram Swami. Practical black-box attacks against machine learning. In Ramesh Karri, Ozgur  
187 Sinanoglu, Ahmad-Reza Sadeghi, and Xun Yi, editors, *Proceedings of the 2017 ACM on Asia*  
188 *Conference on Computer and Communications Security, AsiaCCS 2017, Abu Dhabi, United Arab*  
189 *Emirates, April 2-6, 2017*, pages 506–519. ACM, 2017.
- 190 Vardan Papyan, XY Han, and David L Donoho. Prevalence of neural collapse during the terminal  
191 phase of deep learning training. *Proceedings of the National Academy of Sciences*, 117(40):  
192 24652–24663, 2020.
- 193 Akshay Rangamani, Marius Lindegaard, Tomer Galanti, and Tomaso Poggio. Feature learning in  
194 deep classifiers through intermediate neural collapse. Technical report, Center for Brains, Minds  
195 and Machines (CBMM), 2023.
- 196 Leslie Rice, Eric Wong, and J. Zico Kolter. Overfitting in adversarially robust deep learning. In  
197 *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July*  
198 *2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages 8093–8104.  
199 PMLR, 2020.
- 200 Robert E. Schapire, Yoav Freund, Peter Barlett, and Wee Sun Lee. Boosting the margin: A new  
201 explanation for the effectiveness of voting methods. In Douglas H. Fisher, editor, *Proceedings of*  
202 *the Fourteenth International Conference on Machine Learning (ICML 1997), Nashville, Tennessee,*  
203 *USA, July 8-12, 1997*, pages 322–330. Morgan Kaufmann, 1997.
- 204 Ali Shafahi, Mahyar Najibi, Amin Ghiasi, Zheng Xu, John P. Dickerson, Christoph Studer, Larry S.  
205 Davis, Gavin Taylor, and Tom Goldstein. Adversarial training for free! In Hanna M. Wallach,  
206 Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett,  
207 editors, *Advances in Neural Information Processing Systems 32: Annual Conference on Neural*  
208 *Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC,*  
209 *Canada*, pages 3353–3364, 2019.
- 210 Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The implicit  
211 bias of gradient descent on separable data. *J. Mach. Learn. Res.*, 19:70:1–70:57, 2018.
- 212 Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow,  
213 and Rob Fergus. Intriguing properties of neural networks, 2014.
- 214 Tom Tirer and Joan Bruna. Extended unconstrained features model for exploring deep neural collapse.  
215 In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato,  
216 editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of  
217 *Proceedings of Machine Learning Research*, pages 21478–21505. PMLR, 17–23 Jul 2022.
- 218 Tom Tirer, Haoxiang Huang, and Jonathan Niles-Weed. Perturbation analysis of neural collapse.  
219 *arXiv preprint arXiv:2210.16658*, 2022.
- 220 Eric Wong, Leslie Rice, and J. Zico Kolter. Fast is better than free: Revisiting adversarial training. In  
221 *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia,*  
222 *April 26-30, 2020*. OpenReview.net, 2020.

- 223 Huan Xu, Constantine Caramanis, and Shie Mannor. Robustness and regularization of support vector  
224 machines. *J. Mach. Learn. Res.*, 10:1485–1510, 2009.
- 225 Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric P. Xing, Laurent El Ghaoui, and Michael I. Jordan.  
226 Theoretically principled trade-off between robustness and accuracy. In Kamalika Chaudhuri and  
227 Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine  
228 Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings  
229 of Machine Learning Research*, pages 7472–7482. PMLR, 2019.
- 230 Jinxin Zhou, Xiao Li, Tianyu Ding, Chong You, Qing Qu, and Zhihui Zhu. On the optimization  
231 landscape of neural collapse under MSE loss: Global optimality with unconstrained features. In  
232 Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato,  
233 editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of  
234 *Proceedings of Machine Learning Research*, pages 27179–27202. PMLR, 17–23 Jul 2022.
- 235 Zhihui Zhu, Tianyu Ding, Jinxin Zhou, Xiao Li, Chong You, Jeremias Sulam, and Qing Qu. A  
236 geometric analysis of neural collapse with unconstrained features. In A. Beygelzimer, Y. Dauphin,  
237 P. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*,  
238 2021.

## 239 Appendix

### 240 A Relevant Work

241 **Neural Collapse & Geometric properties of Optimization in Deep Learning.** The term Neural  
242 Collapse was coined by Pappayan et al. [2020] to describe phenomena about the feature representations  
243 of the last layer and the classification weights of a deep neural network at *convergence*. It collectively  
244 refers to the onset of variability collapse of within-class representations (NC1), the formation of  
245 two simplices (NC2) - one from the class-mean representations and another from the classification  
246 weights - that are actually dual (NC3), and, finally, the underlying simplicity of the prediction rule  
247 of the network, which becomes nothing but a simple nearest-neighbor classifier (NC4) (see Section  
248 B for formal definitions). [Pappayan et al., 2020], using ideas from Information Theory, showed that  
249 the formation of a simplex is optimal in the presence of vanishing within-class variability. Mixon  
250 et al. [2022] introduced the *unconstrained features* model (independently proposed by [Fang et al.,  
251 2021] as the Layer-Peeled model), a model where the feature representations are considered as free  
252 optimization variables, and showed that a global optimizer of this problem (for the MSE loss) exhibits  
253 Neural Collapse. Many derivative works have proven Neural Collapse modifying this model, by  
254 either considering other loss functions or trying to incorporate more deep learning elements into  
255 it [Fang et al., 2021, Zhu et al., 2021, Ji et al., 2022, E and Wojtowytsch, 2022, Zhou et al., 2022,  
256 Tirer and Bruna, 2022, Han et al., 2022]. The notion of maximum separability dates back to Support  
257 Vector Machines [Cortes and Vapnik, 1995], while the bias of gradient-based optimization algorithms  
258 towards such solutions has been used to explain the success of boosting methods [Schapire et al.,  
259 1997], and, more recently, to motivate the generalization properties of neural networks [Neyshabur  
260 et al., 2015, Soudry et al., 2018, Lyu and Li, 2020]. The connection between Neural Collapse and  
261 generalization of neural networks (on in-distribution and transfer-learning tasks) has been explored in  
262 [Galanti et al., 2021, Hui et al., 2022]. Finally, the propagation of Neural Collapse inside the network  
263 has been studied by [Ben-Shaul and Dekel, 2022, He and Su, 2022, Hui et al., 2022, Li et al., 2022,  
264 Tirer et al., 2022, Rangamani et al., 2023].

265 **Adversarial Examples & Robustness.** Neural Networks are famously susceptible to adversarial  
266 perturbations of their inputs, even of very small magnitude [Szegedy et al., 2014]. Most of the attacks  
267 that drive the performance of networks to zero are gradient-based [Goodfellow et al., 2015, Carlini  
268 and Wagner, 2017]. These perturbations are surprisingly consistent between different architectures  
269 and hyperparameters, they are in many cases transferable between models [Papernot et al., 2017],  
270 and they can also be made universal (one perturbation for all inputs) [Moosavi-Dezfooli et al., 2017].  
271 For training robust models, one can resort to algorithms from robust optimization [Xu et al., 2009,  
272 Goodfellow et al., 2015, Madry et al., 2018]. In particular, the most effective algorithm used in  
273 deep learning is called *Adversarial Training* [Madry et al., 2018]. During adversarial training one  
274 alternates steps of generating adversarial examples and training on this data instead of the original  
275 one. Several variations of this approach have been proposed in the literature (e.g. Zhang et al. [2019],  
276 Shafahi et al. [2019], Wong et al. [2020]), modifying either the attack used for data generation or the  
277 loss used to measure mistakes. However, models produced by this algorithm, despite being relatively  
278 robust, still fall behind in terms of absolute performance [Croce et al., 2021], while there are still  
279 many unresolved conceptual questions about adversarial training [Rice et al., 2020]. In terms of the  
280 geometrical properties of the solutions, [Li et al., Lv and Zhu, 2022] showed that in some cases (either  
281 in the presence of separable data or/and homogeneous networks) adversarial training converges to a  
282 solution that maximally separates the *adversarial* points.

### 283 B Methodology

284 In this section, we proceed with formal definitions of Neural Collapse (NC), adversarial attacks, and  
285 Adversarial Training (AT), together with the variants we study in this paper.

#### 286 B.1 Notation

287 Let  $\mathcal{X}$  be an input space, and  $\mathcal{Y}$  be an output space, with  $|\mathcal{Y}| = C$ . Denote by  $\mathcal{S}$  a given class-balanced  
288 dataset that consists of  $C$  classes and  $n$  data points per class. Let  $f : \mathcal{X} \rightarrow \mathcal{Y}$  be a neural network,  
289 with its final linear layer denoted as  $\mathbf{W}$ . For each class  $c$ , the corresponding classifier is denoted as  $\mathbf{w}_c$ ,



290 and the bias is called  $b_c$ . Denote the representation of the  $i$ -th sample within class  $c$  as  $\mathbf{h}_{i,c} \in \mathbb{R}^p$ , and  
 291 the union of such representations  $H(\mathcal{S})$ . We define the global-mean vector  $\boldsymbol{\mu}_G \in \mathbb{R}^p$ , and class-mean  
 292 vector  $\boldsymbol{\mu}_c \in \mathbb{R}^p$  associated with  $\mathcal{S}$  as  $\boldsymbol{\mu}_G \triangleq \frac{1}{nC} \sum_{i,c} \mathbf{h}_{i,c}$  and  $\boldsymbol{\mu}_c \triangleq \frac{1}{n} \sum_i \mathbf{h}_{i,c}$ ,  $c = 1, \dots, C$ . For  
 293 brevity, we refer in the text to the globally-centered class-means,  $\{\boldsymbol{\mu}_c - \boldsymbol{\mu}_G\}_{c=1}^C$ , as just *class-means*,  
 294 since these vectors are constituents of the simplex. We denote  $\tilde{\boldsymbol{\mu}}_c = (\boldsymbol{\mu}_c - \boldsymbol{\mu}_G) / \|\boldsymbol{\mu}_c - \boldsymbol{\mu}_G\|_2$   
 295 the normalized class-means. Unless otherwise specified, the term ‘‘representation’’ refers to the  
 296 penultimate layer of the network.

## 297 B.2 Neural Collapse Concepts

298 [Papayan et al., 2020] demonstrate the prevalence of NC on networks optimized by SGD, by tracing  
 299 the following quantities<sup>1</sup>. Throughout our paper, we closely follow Papayan et al. [2020] and Han  
 300 et al. [2022] on formalization of NC1-NC4. Before proceeding to the NC concepts, we introduce

301 **Simplex ETF:** A *standard simplex ETF* composed of  $C$  points is a set of points in  $\mathbb{R}^C$ , each point  
 302 belonging to a column of

$$\sqrt{\frac{C}{C-1}} (\mathbf{I} - \frac{1}{C} \mathbf{1}_C \mathbf{1}_C^\top),$$

303 where  $\mathbf{I} \in \mathbb{R}^{C \times C}$  is the identity matrix and  $\mathbf{1}_C = [1 \ \dots \ 1]^\top \in \mathbb{R}^C$  is the all-ones vector. In our  
 304 discussion, a *simplex* can be thought of as a standard simplex ETF up to partial rotations, reflections,  
 305 and rescaling.

306 **Between-class and within-class covariance:** Using terminology developed in Section B, we define  
 307 between-class covariance  $\Sigma_B \in \mathbb{R}^{p \times p}$  as

$$\Sigma_B \triangleq \text{AVG}_c (\boldsymbol{\mu}_c - \boldsymbol{\mu}_G) (\boldsymbol{\mu}_c - \boldsymbol{\mu}_G)^\top$$

308 and  $\Sigma_W \in \mathbb{R}^{p \times p}$  as

$$\Sigma_W \triangleq \text{AVG}_{i,c} (\mathbf{h}_{i,c} - \boldsymbol{\mu}_c) (\mathbf{h}_{i,c} - \boldsymbol{\mu}_c)^\top.$$

309 Next, we start the introduction of NC concepts. We use the following exact definitions proposed by  
 310 Han et al. [2022]:

311 **(NC1) Variability Collapse:**

$$\Sigma_B^\dagger \Sigma_W \rightarrow \mathbf{0},$$

312 where  $\dagger$  denotes the Moore-Penrose inverse. The NC1 curve corresponds to  $\text{Tr}(\Sigma_B^\dagger \Sigma_W)$ .

313 **(NC2) Convergence to Simplex ETF:**

$$\langle \tilde{\boldsymbol{\mu}}_c, \tilde{\boldsymbol{\mu}}_{c'} \rangle \rightarrow -\frac{1}{C-1} \quad \forall c \neq c'$$

314

$$\left| \|\boldsymbol{\mu}_c - \boldsymbol{\mu}_G\|_2 - \|\boldsymbol{\mu}_{c'} - \boldsymbol{\mu}_G\|_2 \right| \rightarrow 0 \quad \forall c, c'$$

315 The NC2 Equinorm curve corresponds to the variation of  $\|\boldsymbol{\mu}_c - \boldsymbol{\mu}_G\|_2$  across all labels  $c$ , the standard  
 316 deviation of these  $c$  quantities:  $\text{std}(\|\boldsymbol{\mu}_c - \boldsymbol{\mu}_G\|_2)$ . The NC2 Equiangular curve corresponds to  
 317  $\text{AVG}_{c \neq c'} \text{abs}(\langle \tilde{\boldsymbol{\mu}}_c, \tilde{\boldsymbol{\mu}}_{c'} \rangle + \frac{1}{C-1})$ , where *abs* is the absolute value operator.

318 **(NC3) Convergence to self-duality:**

$$\frac{\mathbf{w}_c}{\|\mathbf{w}_c\|_2} - \frac{\boldsymbol{\mu}_c - \boldsymbol{\mu}_G}{\|\boldsymbol{\mu}_c - \boldsymbol{\mu}_G\|_2} \rightarrow \mathbf{0} \quad \forall c.$$

The NC3 curve corresponds to

$$\sqrt{\sum_c \left\| \frac{\mathbf{w}_c}{\|\mathbf{w}_c\|_2} - \frac{\boldsymbol{\mu}_c - \boldsymbol{\mu}_G}{\|\boldsymbol{\mu}_c - \boldsymbol{\mu}_G\|_2} \right\|_2^2}.$$

319 **(NC4) Simplification to NCC classifier:**

$$\arg \max_{c'} \langle \mathbf{w}_{c'}, \mathbf{h} \rangle + b_{c'} \rightarrow \arg \min_{c'} \|\mathbf{h} - \boldsymbol{\mu}_{c'}\|_2 \quad \forall \mathbf{h} \in H(\mathcal{S}).$$

320 The NC4 curve corresponds to the mismatch ratio of these two quantities.

<sup>1</sup>In particular, Neural Collapse becomes more evident in the so-called *Terminal Phase of Training*, the phase beyond the point of (effectively) zero training loss.

321 In our experiments, we calculate the NC statistics with the code provided by Han et al. [2022]<sup>2</sup>.  
 322 Furthermore, in this work, we compare representations of original and perturbed data, which imposes  
 323 ambiguity on which class-mean vectors  $\mu_c, \mu_G$  to use (from  $\mathcal{S}$  or  $\mathcal{S}'$ ). In the spirit of the original  
 324 definitions, for NC1-4 we will use the class means induced by the dataset  $\mathcal{S}'$ , even if different from  
 325 the training set. NC4 studies the predictive power of the NCC classifier on  $\mathcal{S}'$  by comparing it to  
 326 the network classification output, which at TPT is equivalent to the ground truth label. For study of  
 327 reference data  $\mathcal{S}'$  outside the TPT, we introduce two quantities, also applicable to any intermediate  
 328 layer:

329 **NCC-Network Matching Rate:** measures the rate at which the NCC classifier defined in NC4  
 330 trained on  $\mathcal{S}$  coincides with the output of the network on dataset  $\mathcal{S}'$ . Note that we use  $\mu_c$  calculated  
 331 by  $\mathcal{S}$ .

$$\arg \max_{c'} \langle \mathbf{w}_{c'}, \mathbf{h} \rangle + b_{c'} \stackrel{?}{=} \arg \min_{c'} \|\mathbf{h} - \mu_{c'}\|_2, h \in H(\mathcal{S}').$$

332 **NCC Accuracy:** measures the accuracy on dataset  $\mathcal{S}'$  of the NCC classifier defined in NC4 trained  
 333 on  $\mathcal{S}$ . Note that we use  $\mu_c$  calculated by  $\mathcal{S}$ .  $c_h$  denotes the ground-truth label of the input.

$$c_h \stackrel{?}{=} \arg \min_{c'} \|\mathbf{h} - \mu_{c'}\|_2, h \in H(\mathcal{S}').$$

334 Note that when  $\mathcal{S} = \mathcal{S}'$ , both NCC-Network Matching Rate and NCC Accuracy stem from (NC4).  
 335 We also introduce the following measures to quantify the proximity of two simplices over  $C$ -classes:

**Simplex Similarity:** We define the similarity measure between two  $C$ -class simplices with normal-  
 ized class means  $\tilde{\mu}_c, \tilde{\mu}'_c$  as

$$\text{AVG}_c \arccos \langle \tilde{\mu}_c, \tilde{\mu}'_c \rangle.$$

**Non-centered Angular Distance:** Similarly, given two simplices, without taking the global mean  $\mu_G$   
 and  $\mu'_G$  into account, we can calculate the angular distance with non-centered class-means directly:

$$\text{AVG}_c \arccos \left\langle \frac{\mu_c}{\|\mu_c\|_2}, \frac{\mu'_c}{\|\mu'_c\|_2} \right\rangle.$$

336 Note that the similarity and angular distance between a simplex and itself is zero.

### 337 B.3 Gradient-Based Adversarial Attack, Adversarial Training (AT), and TRADES

338 Given a deep neural network  $f$  with parameters  $\theta$ , a clean example  $(\mathbf{x}, y)$  and cross-entropy loss  
 339  $\mathcal{L}(\cdot, \cdot)$ , the *untargeted* adversarial perturbation is crafted by running multiple steps of projected  
 340 gradient descent (PGD) to maximize the CE loss [Kurakin et al., 2017, Madry et al., 2018] (in what  
 341 follows, we focus on  $\ell_\infty$  adversary with  $\ell_2$  deferred to the appendix):

$$\mathbf{x}^{k+1} = \Pi_{\mathcal{B}_{\mathbf{x}^0}^\epsilon} (\mathbf{x}^k + \alpha \cdot \text{sign}(\nabla_{\mathbf{x}^k} \mathcal{L}(f(\mathbf{x}^k), y))), \quad (1)$$

342 where  $\mathbf{x}^0 = \mathbf{x}$  is the original example,  $\alpha$  is the step size,  $\tilde{\mathbf{x}} = \mathbf{x}^N$  is the final adversarial example,  
 343 and  $\Pi$  is the projection on the valid  $\epsilon$ -constraint set,  $\mathcal{B}_{\mathbf{x}^0}^\epsilon$ , of the data.  $\mathcal{B}_{\mathbf{x}^0}^\epsilon$  is usually taken as either an  
 344  $\ell_\infty$  or  $\ell_2$  ball centered in  $\mathbf{x}^0$ . Further, to control the predicted label of  $\tilde{\mathbf{x}}$ , a variant called *targeted*  
 345 *attack* minimizes the CE loss w.r.t. a target label  $y_t \neq y$ :

$$\mathbf{x}^{k+1} = \Pi_{\mathcal{B}_{\mathbf{x}^0}^\epsilon} (\mathbf{x}^k - \alpha \cdot \text{sign}(\nabla_{\mathbf{x}^k} \mathcal{L}(f(\mathbf{x}^k), y_t))). \quad (2)$$

346 With a standardly-trained network, both these methods can effectively reduce the accuracy to 0%. To  
 347 combat this phenomenon, robust optimization algorithms have been proposed. The most representa-  
 348 tive methodology, *adversarial training* [Madry et al., 2018], generates  $\tilde{\mathbf{x}}$  *on-the-fly* with Equation (1)  
 349 for each epoch from  $\mathbf{x}$ , and takes the model-gradient update on  $\tilde{\mathbf{x}}$  only.

350 An alternative robust training variant, TRADES [Zhang et al., 2019], is of particular interest as it  
 351 aims to address both robustness and clean accuracy. Thus the gradient steps of TRADES directly  
 352 involve both  $\mathbf{x}$  and  $\tilde{\mathbf{x}}$ , where  $\tilde{\mathbf{x}}$  is also obtained by PGD, but under the KL-divergence loss:

$$\mathbf{x}^{k+1} = \Pi_{\mathcal{B}_{\mathbf{x}^0}^\epsilon} (\mathbf{x}^k + \alpha \cdot \text{sign}(\nabla_{\mathbf{x}^k} \mathcal{L}_{KL}(f(\mathbf{x}), f(\mathbf{x}^k)))). \quad (3)$$

<sup>2</sup><https://colab.research.google.com/github/neuralcollapse/neuralcollapse/blob/main/neuralcollapse.ipynb>

353 The total TRADES loss is a summation of the CE loss on the clean data and a KL-divergence (KLD)  
354 loss between the predicted probability of  $\mathbf{x}$  and  $\bar{\mathbf{x}}$  with a regularization constant  $\beta$ :

$$\mathcal{L}_{CE}(f(\mathbf{x}), y) + \beta \cdot \mathcal{L}_{KL}(f(\mathbf{x}), f(\bar{\mathbf{x}})). \quad (4)$$

## 355 C Experimental Details

356 **Code.** For  $\ell_\infty$  and  $\ell_2$  PGD attacks with ST and AT, we used the code from Rice et al. [2020]<sup>3</sup>. For  
357 TRADES, we adopted the original implementation<sup>4</sup>. We have attached the code for reproducing NC  
358 results with ST, AT, and TRADES within a zip file.

359 **Plotting.** Throughout our paper, we plot all quantities per 5 epochs in all figures.

360 **Layerwise NC.** We study the layerwise NC1, NC2 and NC4 quantities for both PreActResNet18  
361 (ResNet18) and VGG11. With ResNet18, which consists of one convolutional layer, four residual  
362 blocks, and the final linear layer, we use the features after every block for the first five blocks (one  
363 convolutional layer and four residual blocks) as representations. With VGG, which consists of eight  
364 convolutional blocks (convolutional layer + batch-normalization + max-pooling) and the final linear  
365 layer, we use the features after each convolutional block as representations. We apply average-pooling  
366 subsampling on representations that are too large for feasible computation of NC1’s pseudo-inverse.

## 367 D TRADES Results on CIFAR-10 with $\ell_\infty$ adversary

368 For CIFAR-10’s results with TRADES, we have produced Figure 4, which depicts the evolution  
369 of loss, accuracy and all of the NC metrics under the standard  $\ell_\infty$  adversary. Note that we plot  
370 the KLD-loss here to showcase optimization convergence, to avoid the effect of the regularization  
371 constant  $\beta$ .

## 372 E Complementary Results on CIFAR-10, $\ell_2$ adversary.

373 Here we complement our main text with robust network experiments on CIFAR-10 for  $\ell_2$  adversarial  
374 perturbations.

375 Figure 5 illustrates NC results of Adversarial Training and TRADES training with the  $\ell_2$  adversary.  
376 All plots are consistent with our findings in the main text: Adversarial Training alters Neural Collapse  
377 such that the clean representation simplex overlaps with the perturbed representation simplex, whereas  
378 TRADES does not lead to any simplex ETF.

## 379 F Complementary Results on CIFAR-100

380 In this section, we reproduce our experiments on CIFAR-100. We illustrate results with  $(\ell_\infty, \ell_2)$   
381 adversaries and obtain the same conclusions as those on CIFAR-10. This suggests the universality of  
382 the intrinsic adversarial perturbation dynamics that we have detailed in the main text.

### 383 F.1 CIFAR-100 $\ell_\infty$ Standard and Adversarial Training Results

384 All results are summarized within Figure 6. Similar to the main text, we plot the untargeted attack  
385 illustration in Figure 7. Notably, on CIFAR-100 with ST, adversarial perturbations also push the  
386 representation to leap toward the predicted class’s simplex cluster with very small angular deviation.

### 387 F.2 CIFAR-100 $\ell_\infty$ TRADES Results

388 For CIFAR-100  $\ell_\infty$  trained with TRADES, Figure 8 depicts the results, and we observe that no  
389 simplex exists, consistent with previous results.

<sup>3</sup>[https://github.com/locuslab/robust\\_overfitting](https://github.com/locuslab/robust_overfitting)

<sup>4</sup><https://github.com/yaodongyu/TRADES>

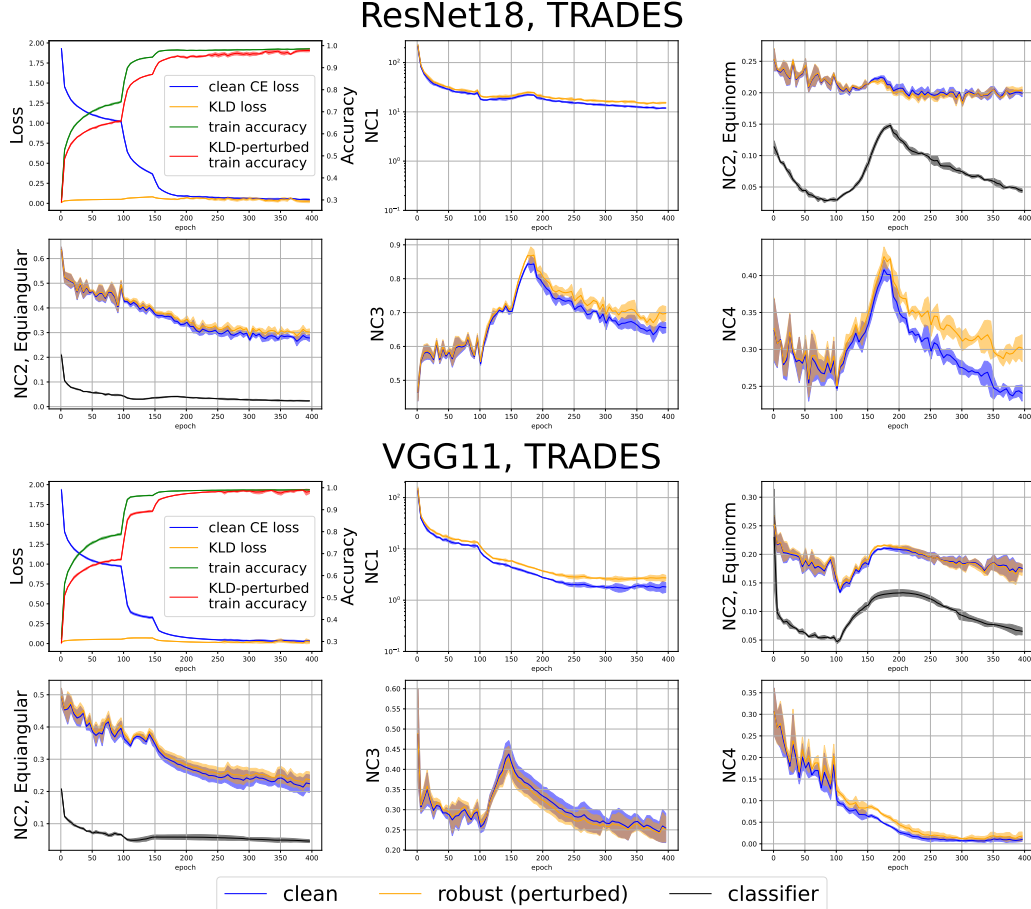


Figure 4: Accuracy, Loss and NC evolution with TRADES trained networks. *Upper*: ResNet18; *Lower*: VGG11. No simplices are formed with TRADES training. Setting: CIFAR-10,  $\ell_\infty$  adversary.

### 390 **F.3 CIFAR-100 $\ell_2$ AT and TRADES Results**

391 These results are shown in Figure 9. All observations are consistent with previous results.

### 392 **F.4 CIFAR-100 Simplex Similarity Results**

393 The Simplex Similarity and non-centered Angular Distance of the simplices formed by targeted  
 394 adversarial and clean examples with ST, and the simplices generated by clean and perturbed examples  
 395 with AT, are depicted in Figure 10. The result is the same as the one for CIFAR-10 in the main text,  
 396 Figure 3.

### 397 **G Layerwise Results**

398 While originally variability collapse and simplex formation were observed for the last layer repre-  
 399 sentations, follow-up studies extended the analysis to the intermediate layers of the neural network.  
 400 In particular, He and Su [2022] found that the amount of variability collapse measured at different  
 401 layers (at convergence) decreases smoothly as a function of the index of the layer. Further, Hui  
 402 et al. [2022] coined the term Cascading Neural Collapse to describe the phenomenon of cascading  
 403 variability collapse; starting from the end of the network, the collapse of one layer seemed to be  
 404 signaling the collapse of the previous layers (albeit to a lesser extent). Here, we replicate this study of  
 405 the intermediate layer computations, while also studying the representations of the perturbed points  
 406 (both in standard and adversarial training). In particular, we collect the input of either convolutional

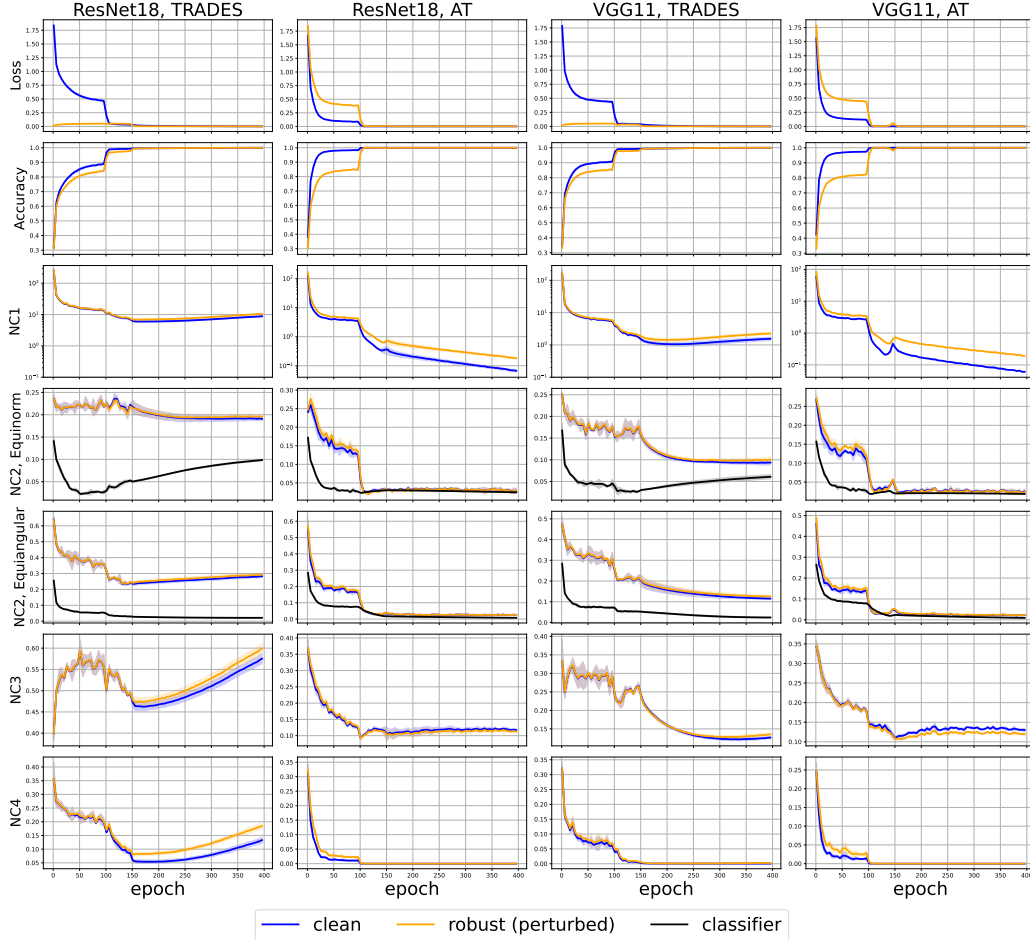


Figure 5: Accuracy, Loss and NC evolution with adversarially trained and TRADES trained networks. Setting: CIFAR-10,  $\ell_2$  adversary.

407 or linear layers of the network *at convergence*, order them by depth index, and compute the NC  
 408 quantities of Section B. The results are presented in Figure 11.

409 Both for ST and AT models, we reproduce the power law behavior observed in [He and Su, 2022]  
 410 for clean data; the feature variability collapses progressively, and, interestingly, undergoes a slower  
 411 decrease in the case of adversarial training. The adversarial data representations for ST models,  
 412 however, while failing to collapse at the final layer (as already established in Figure 2), exhibit the  
 413 same amount of “clustering” as those of the original data for the earlier layers. This hints that from the  
 414 viewpoint of the earlier layers, clean and adversarial data are indistinguishable. And, this, is indeed  
 415 the case! Looking at the first and third column of Figure 12, we observe that the simple classifier  
 416 formed by the centers of the early layers is quite robust ( $\sim 40\%$ ) to these adversarial examples (both  
 417 train and test). Curiously, this robustness is higher than the one of the simple classifiers defined by  
 418 layers of an adversarially trained model (although the two numbers are not directly comparable). This  
 419 is, undeniably, a peculiar phenomenon of standardly trained models that is worth more exploration;  
 420 could it be that the lesser variability exhibited in the earlier layers is actually beneficial for robustness  
 421 or is it just the stability of the feature space that makes prediction more robust?

422 In Figure 13 and Figure 14, we perform the same computations on CIFAR-100. We arrive at the same  
 423 conclusions.

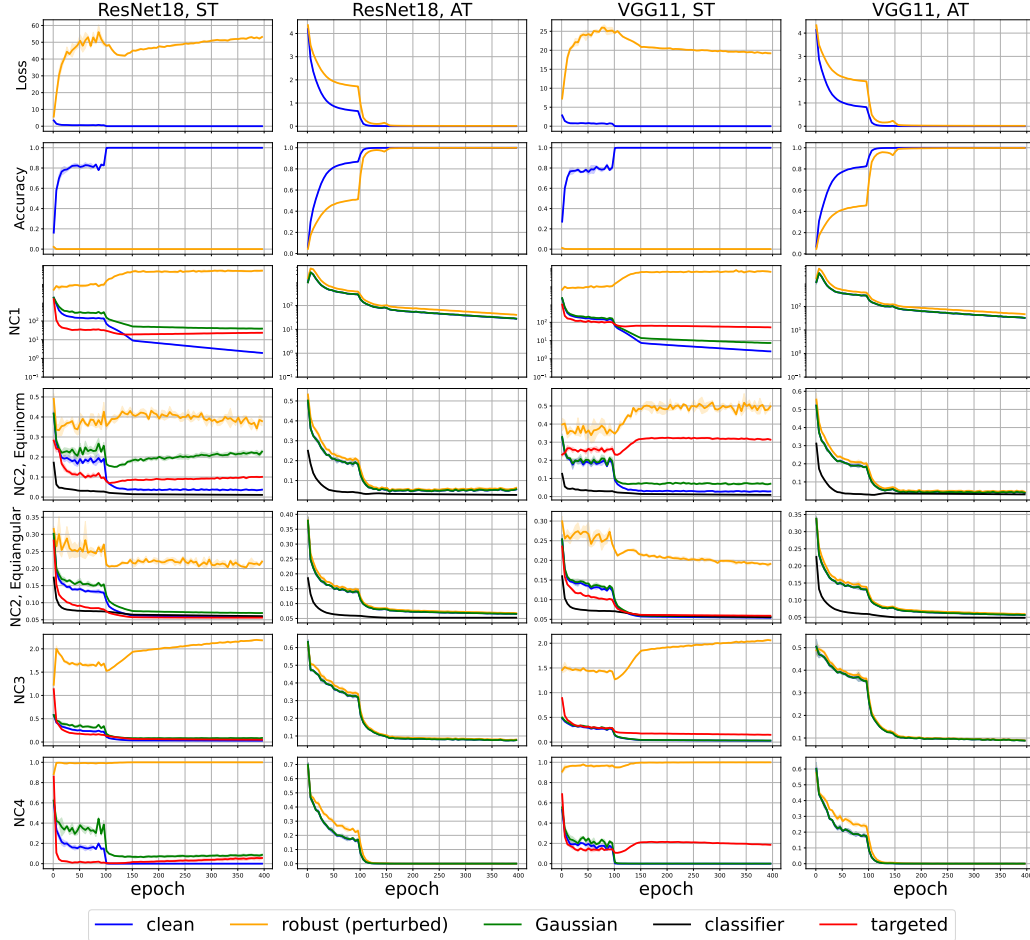


Figure 6: Accuracy, Loss and NC evolution with standardly trained networks. Setting: CIFAR-100,  $\ell_\infty$  adversary.

## 424 H Small Epsilon Results

425 Here, we illustrate how AT indeed progressively induces more robust NC metrics and sim-  
 426 plex ETFs with respect to the perturbation radius  $\epsilon$ . Figure 15 shows the NC metrics over  
 427  $8/255$ -perturbed data. Conversely, using an ST model, the NC metrics when evaluating on  
 428  $(2/255, 4/255, 8/255)$ -perturbed data also increases monotonically with adversarial strength. This  
 429 is illustrated in Figure 16.<sup>5</sup>

<sup>5</sup>For small radius AT and small radius adversarial attack for ST, we scale the PGD step size  $\alpha$  linearly with  $\epsilon$  to ensure PGD to work properly.



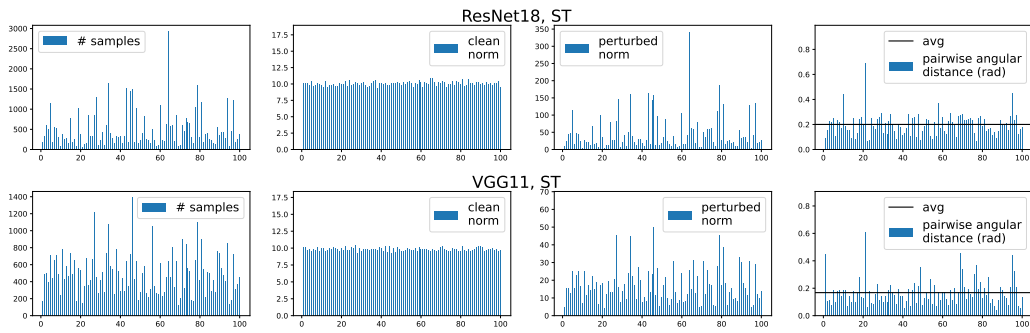


Figure 7: Illustration of untargeted adversarial attacks on standardly trained, converged, models that correspond to one random seed. (CIFAR-100,  $\ell_\infty$ ). *Left*: Number of examples with a certain predicted label. *Inner Left*: The norms of clean class-means. *Inner Right*: The norms of predicted class-means with perturbed data. *Right*: Angular distance between clean and predicted class-mean with perturbed data. *Upper*: ResNet18; *Lower*: VGG11. For 100 classes, the between-class angular distance is  $\arccos(-\frac{1}{99}) = 1.58$  rad = 90.58 degrees, while 0.2 rad is only 11.4 degrees.

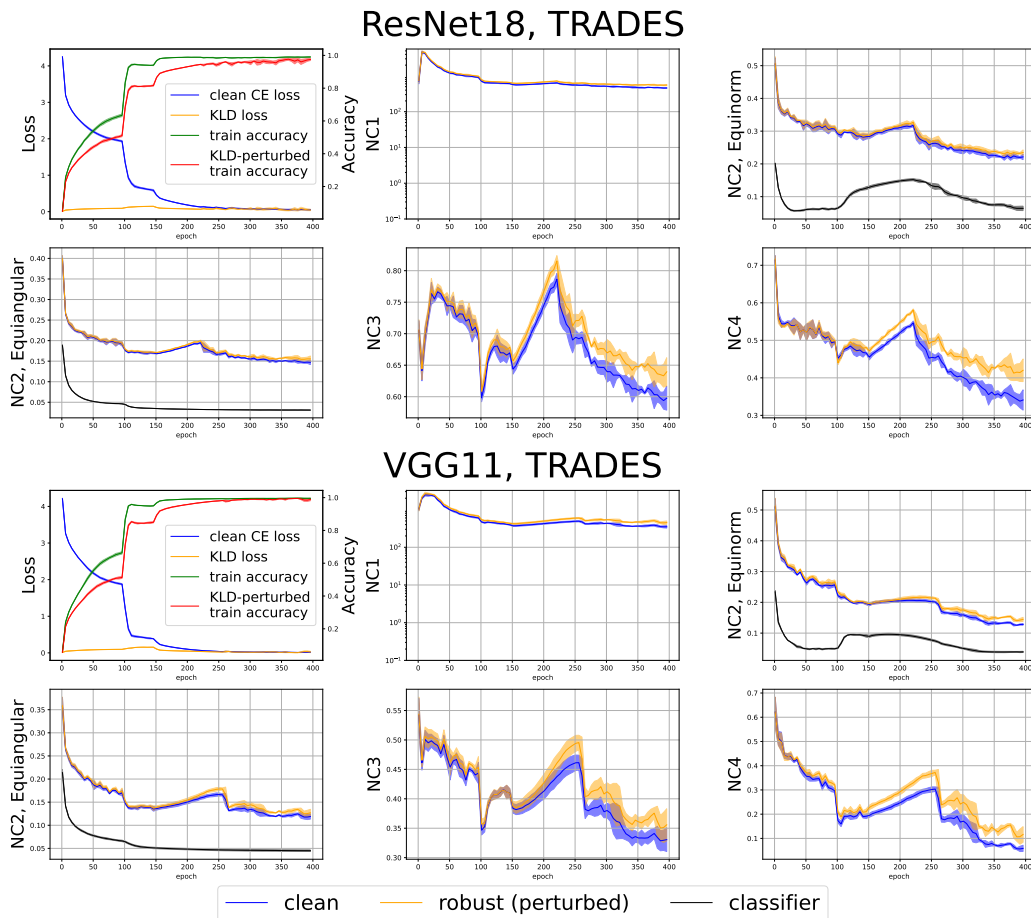


Figure 8: Accuracy, Loss and NC evolution with TRADES trained networks. *Upper*: ResNet18; *Lower*: VGG11. Results indicate AT boosts Neural Collapse so that it also happens on adversarially-perturbed data. Setting: CIFAR-100,  $\ell_\infty$  adversary.

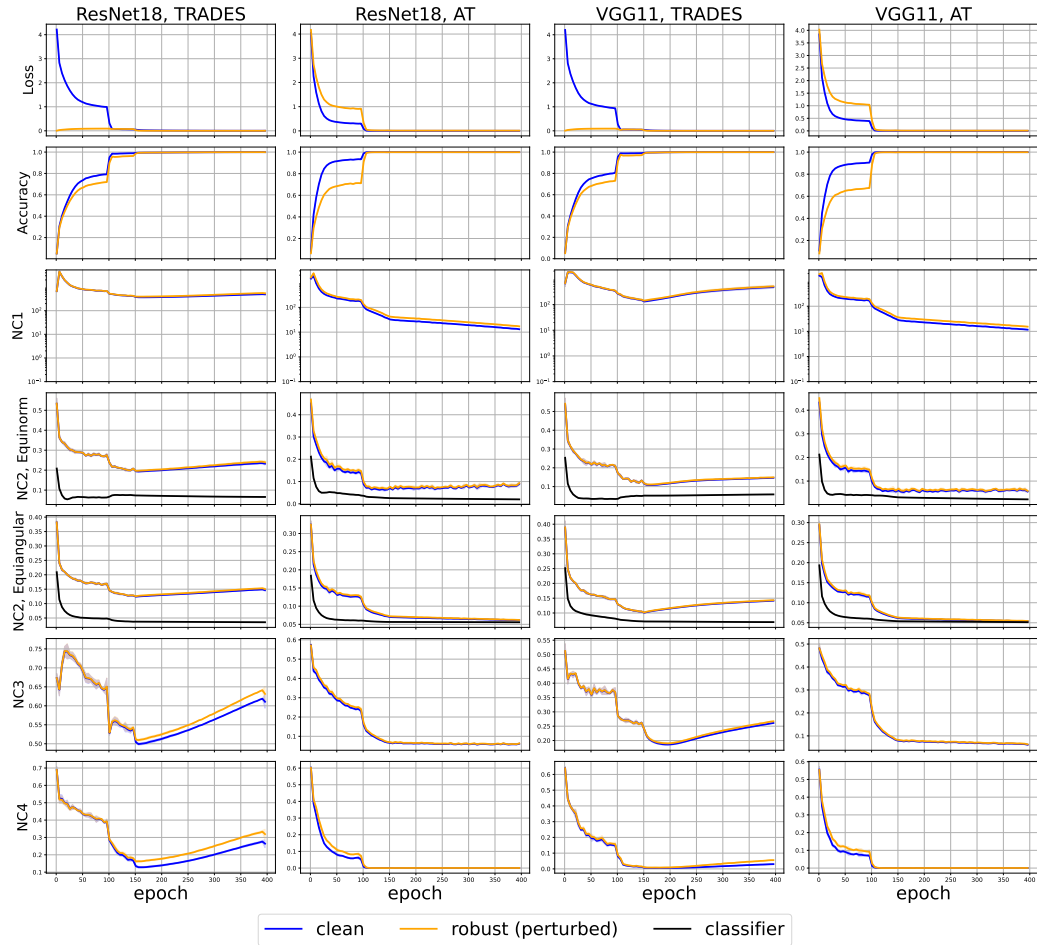


Figure 9: Accuracy, Loss and NC evolution with  $\ell_2$  robust models on CIFAR-100. Setting: CIFAR-100,  $\ell_2$  adversary.

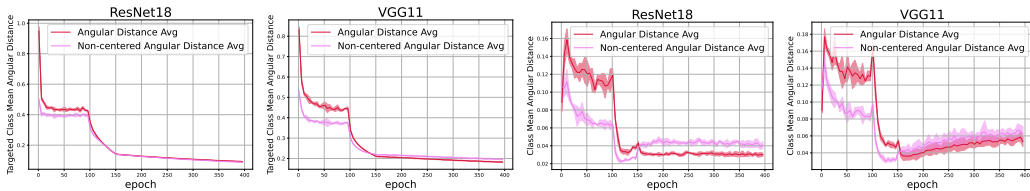


Figure 10: Angular distance. *Left and Inner Left*: Average between targeted attack class-means and clean class-means on **ST** network. *Inner Right and Right*: Average between perturbed class-means and clean class-means on **AT** network. Setting: CIFAR-100,  $\ell_\infty$  adversary.

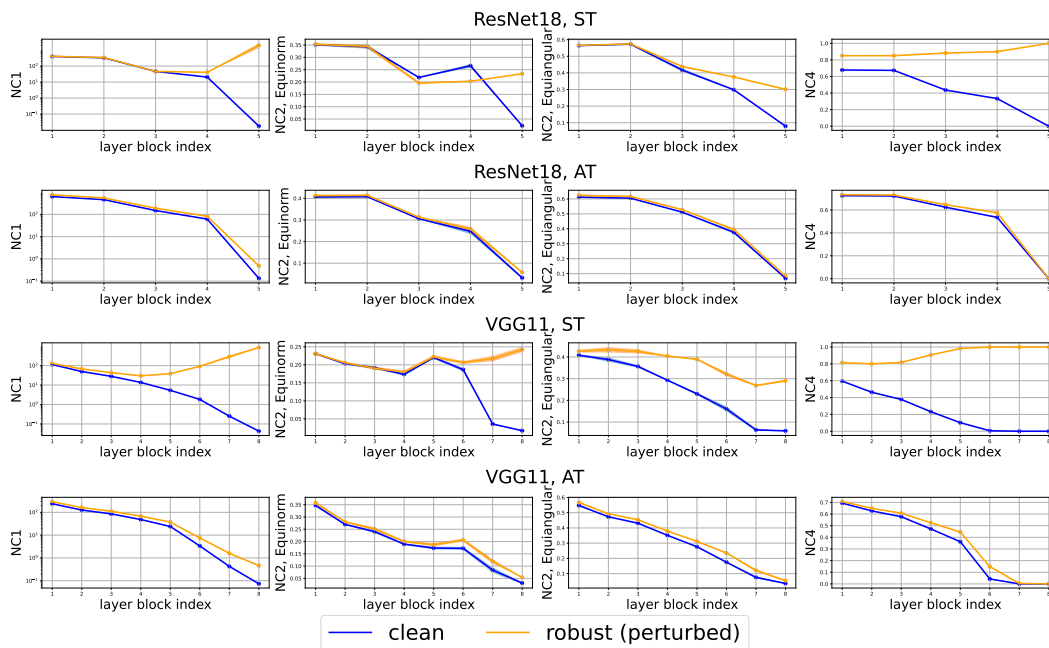


Figure 11: Layerwise evolution of NC1, NC2 and NC4 for ST and AT networks. NC metrics for perturbed data tend to undergo some amount of clustering in the earlier layers. For AT, collapse undergoes a slower decrease through layers than for ST. Setting: CIFAR-10,  $\ell_\infty$  adversary.

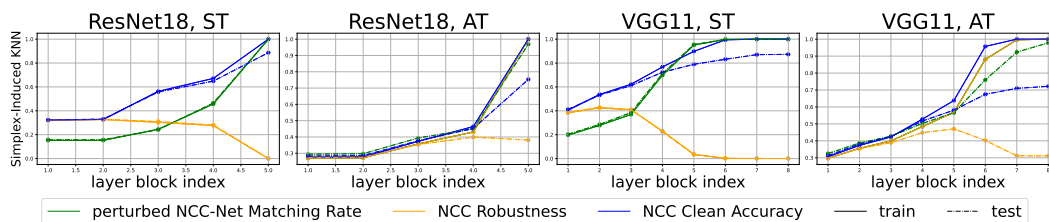


Figure 12: Layerwise NCC classifier. We measure the performance of the NCC classifier obtained from (training) class means on both train and test data. NCC Robustness refers to NCC Accuracy on perturbed data. Note that, on training data, the NCC Robustness and the perturbed NCC-Net Matching Rate curves overlap. Early layers give a surprisingly robust NCC classifier (NCC Robustness) for both train and test data. Setting: CIFAR-10,  $\ell_\infty$  adversary.

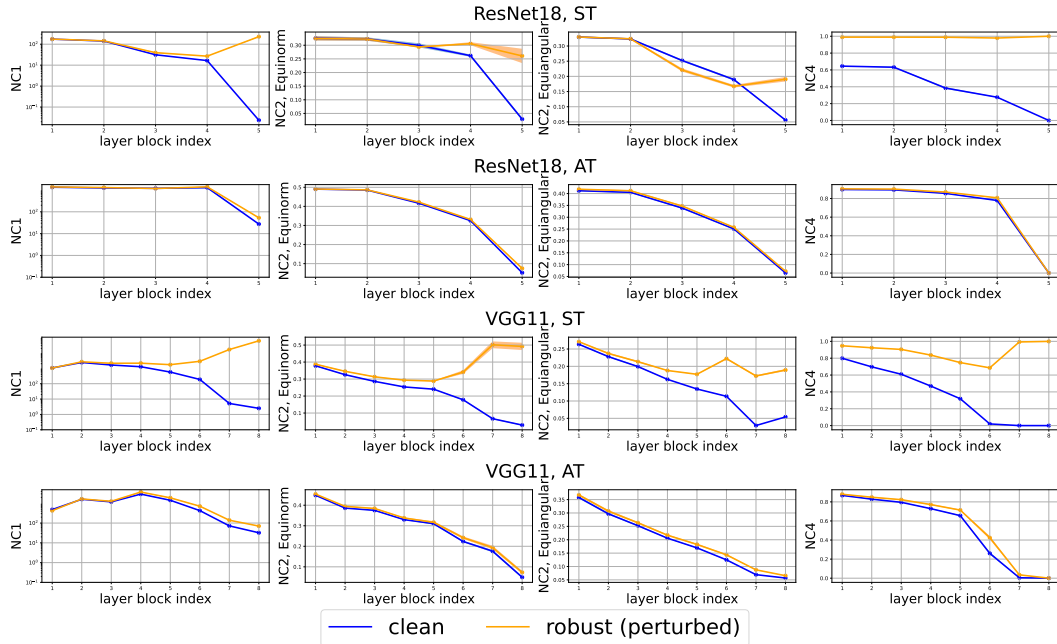


Figure 13: Layerwise evolution of NC1, NC2 and NC4 for ST and AT networks. NC metrics for perturbed data tend to undergo some amount of clustering in the earlier layers. For AT, collapse undergoes a slower decrease through layers than for ST. Setting: CIFAR-100,  $\ell_\infty$  adversary.

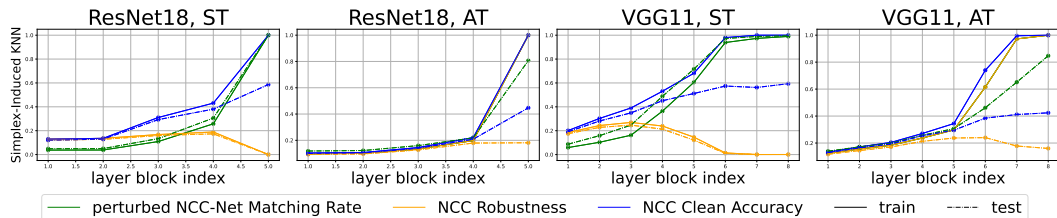


Figure 14: Layerwise NCC classifier. We measure the performance of the NCC classifier obtained from (training) class means on both train and test data. NCC Robustness refers to NCC Accuracy on perturbed data. Note that, on training data, the NCC Robustness and the perturbed NCC-Net Matching Rate curves overlap. Early layers give a surprisingly robust NCC classifier (NCC Robustness) for both train and test data. Setting: CIFAR-100,  $\ell_\infty$  adversary.

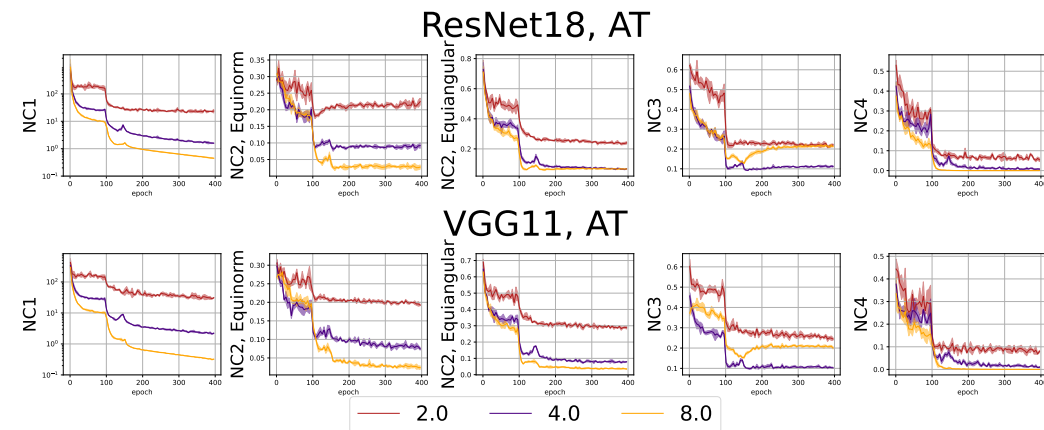


Figure 15: Progressive Loss and NC evolution, AT with varying strength. The color indicates the epsilon used for **training**. Setting: CIFAR-10,  $\ell_\infty$  adversary.

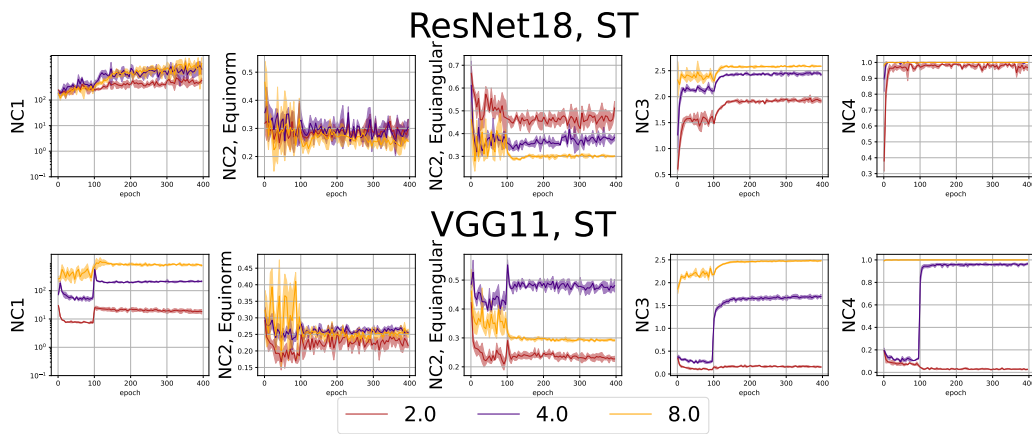


Figure 16: Progressive Loss and NC evolution, ST with varying attacking strength. The color indicates the epsilon used for **evaluation**. Setting: CIFAR-10,  $\ell_\infty$  adversary.