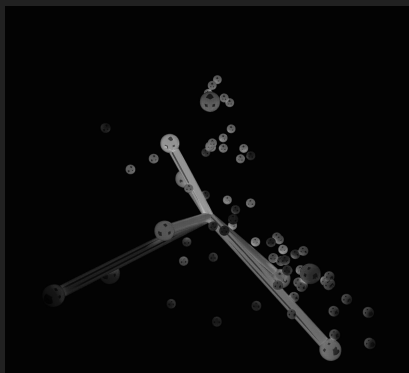
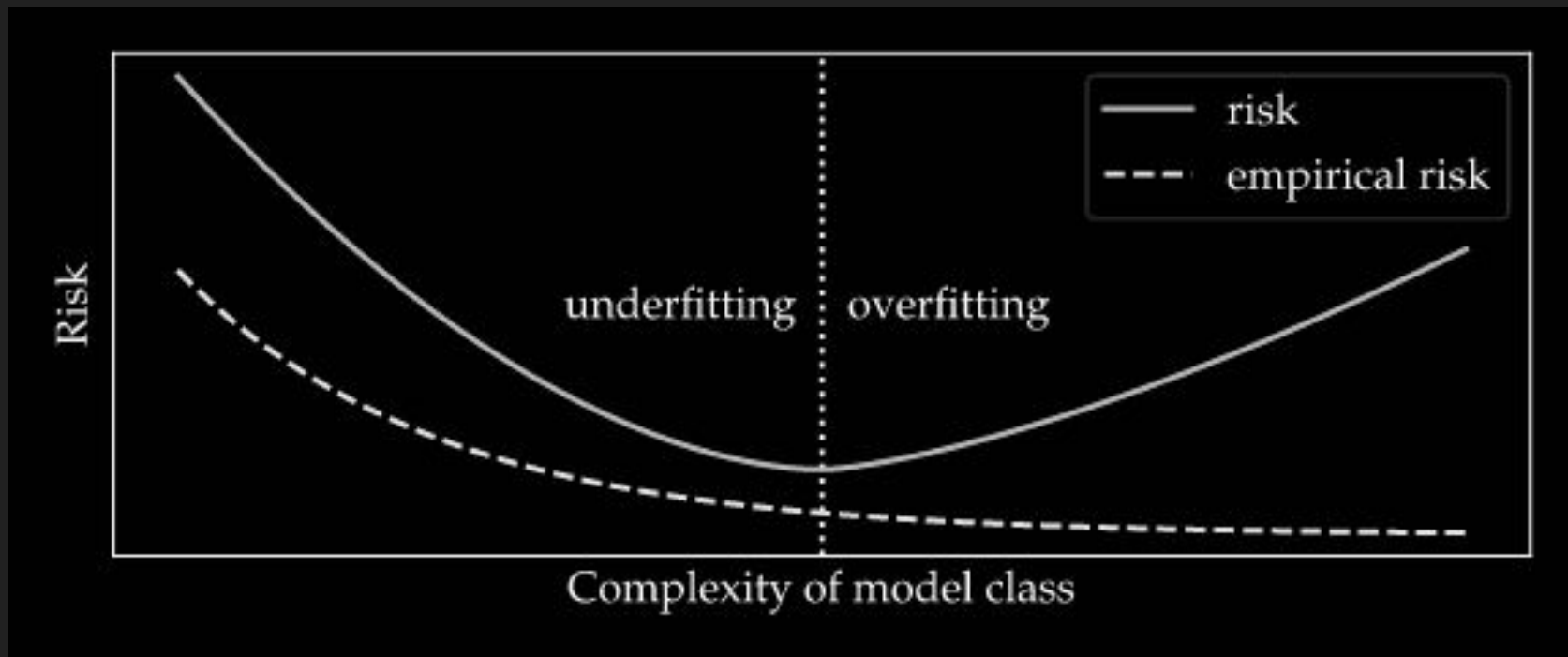


# Prevalence of neural collapse during the terminal phase of deep learning training

Vardan Papyan, X. Y. Han, and David L. Donoho



# Understanding Model Generalization



# Overparameterized Models

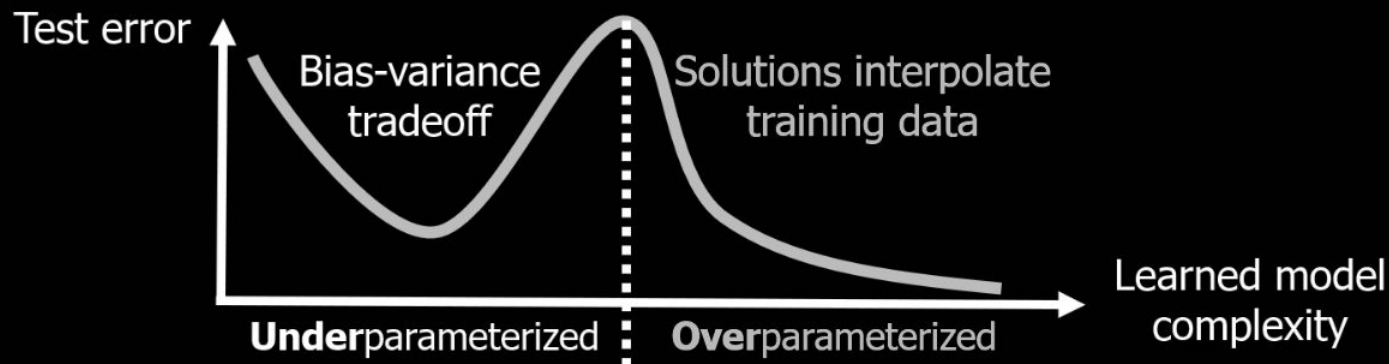
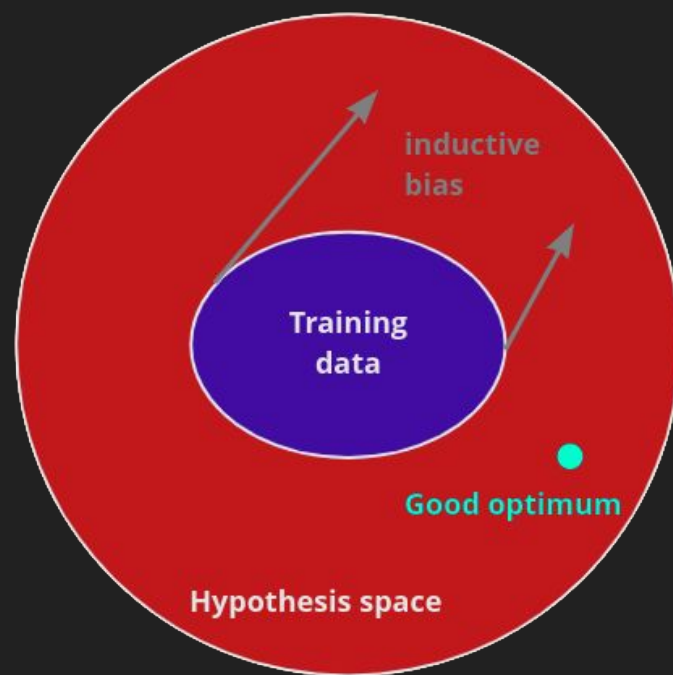
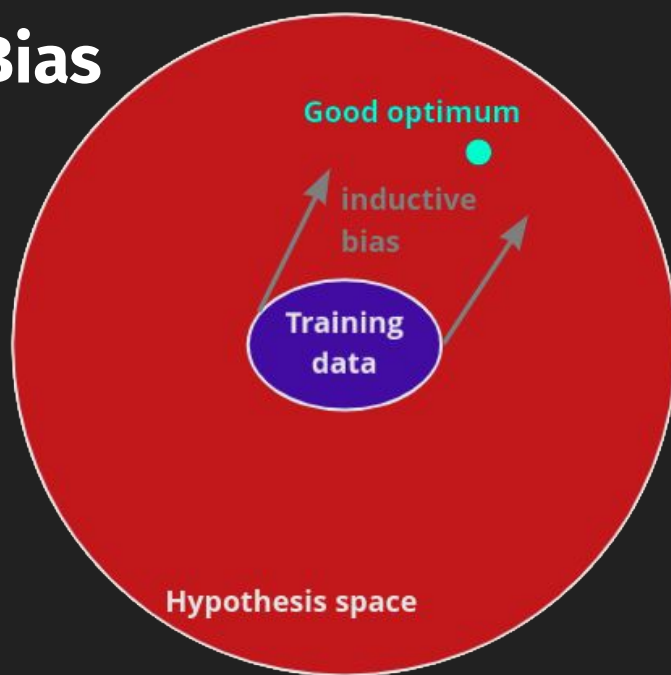


Figure 1: Double descent of test errors (i.e., generalization errors) with respect to the complexity of the learned model. TOPML studies often consider settings in which the learned model complexity is expressed as the number of (independently tunable) parameters in the model. In this qualitative demonstration, the global minimum of the test error is achieved by maximal overparameterization.

# Inductive Bias



In a low data setting, right inductive bias may help to find good optimum, but in a rich data setting, it may lead to constraints that harm generalization. [\[Image Source\]](#)



(a) Texture image

81.4%	<b>Indian elephant</b>
10.3%	indri
8.2%	black swan



(b) Content image

71.1%	<b>tabby cat</b>
17.3%	grey fox
3.3%	Siamese cat



(c) Texture-shape cue conflict

63.9%	<b>Indian elephant</b>
26.4%	indri
9.6%	black swan

# In Brief: What is Neural Collapse?

## NC1: Collapse of variability:

For data samples belonging to the same class, their final hidden layer features **concentrate around their class mean**. Thus, the variability of intra class features during training is lost as they collapse to a point.

## NC2: Preference towards a simplex ETF:

The class means of the penultimate layer features tend to form a **simplex equiangular tight frame** (simplex ETF). A simplex ETF is a symmetric structure whose **vertices lie on a hyper-sphere, are linearly separable and are placed at the maximum possible distance** from each other.

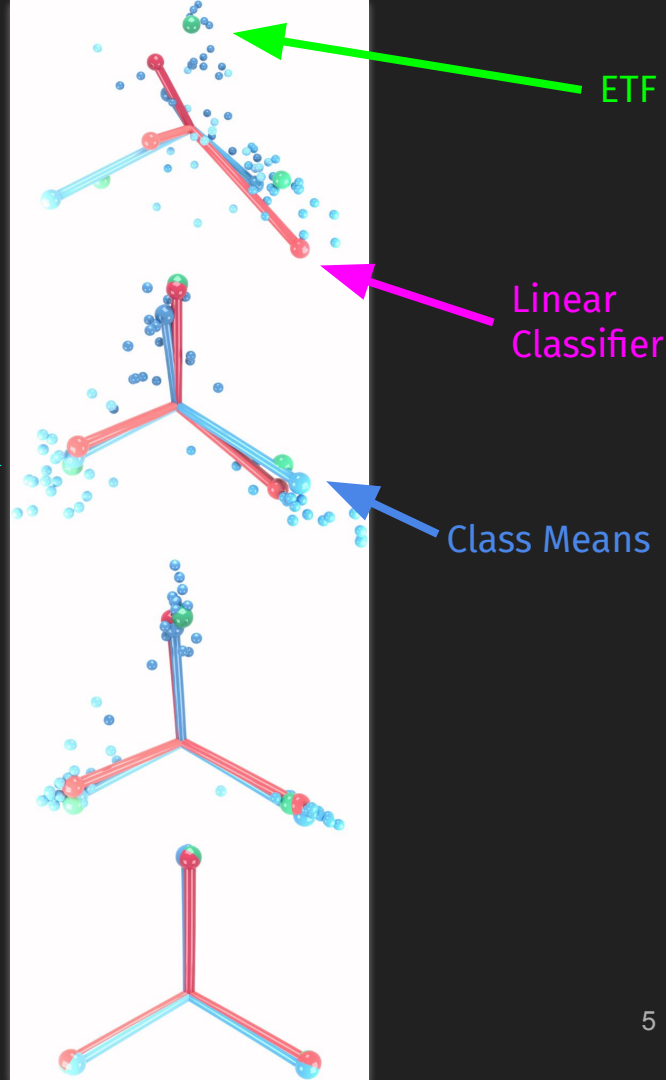
## NC3: Self-dual alignment:

The **columns of the last layer linear classifier matrix** also form a simplex ETF in their dual vector space and converge to the simplex ETF (up to rescaling) of the penultimate layer features.

## NC4: Choose the nearest class mean:

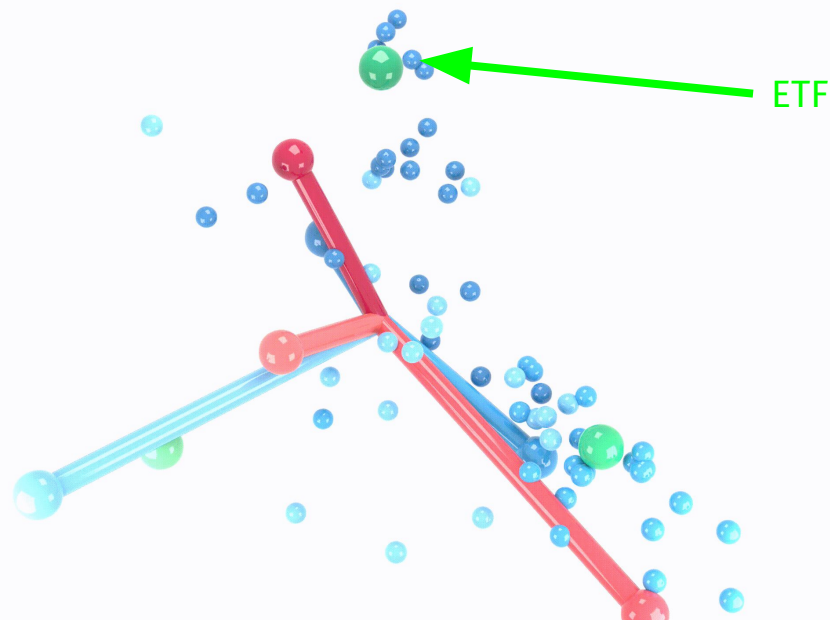
When a test point is to be classified, the last-layer classifier essentially acts as a **nearest (train)-class mean decision rule w.r.t penultimate layer features**.

Activations



Class 0

Class 1



Final  
Layer

**Visualization of NC.** The figure depicts, in three dimensions, NC as training proceeds from top to bottom. Green spheres represent the vertices of the standard simplex, see Definition (Simplex ETF), red ball and sticks represent linear classifiers, blue ball and sticks represent class means, and small blue spheres represent last-layer features. For all objects, we distinguish different classes via the shade of the color. As training proceeds, last-layer features collapse onto their class means (NC1), class means converge to the vertices of the simplex ETF (NC2), and the linear classifiers approach their corresponding class means (NC3). An animation can be found at <https://purl.stanford.edu/br193mh4244>.

## Definition: Simplex ETF

$$\mathbf{M}^\star = \sqrt{\frac{C}{C-1}} \left( \mathbf{I} - \frac{1}{C} \mathbf{1}_C \mathbf{1}_C^\top \right),$$

where  $\mathbf{I} \in \mathbb{R}^{C \times C}$  is the identity matrix, and  $\mathbf{1}_C \in \mathbb{R}^C$  is the ones vector. In this paper, we allow other poses, as well as rescaling, so the general simplex ETF consists of the points specified by the columns of  $\mathbf{M} = \alpha \mathbf{U} \mathbf{M}^\star \in \mathbb{R}^{p \times C}$ , where  $\alpha \in \mathbb{R}_+$  is a scale factor, and  $\mathbf{U} \in \mathbb{R}^{p \times C}$  ( $p \geq C$ ) is a partial orthogonal matrix ( $\mathbf{U}^\top \mathbf{U} = \mathbf{I}$ ).

# Notation

Weights/Bias

$$\mathbf{W} \in \mathbb{R}^{C \times p} \quad \mathbf{b} \in \mathbb{R}^C$$

Activations from Encoder

$$\mathbf{h} = \mathbf{h}_{\theta}(\mathbf{x}).$$

Linear Classifier

$$\mathbf{W}\mathbf{h}(\mathbf{x}) + \mathbf{b}.$$

Objective

$$\min_{\theta, \mathbf{W}, \mathbf{b}} \sum_{c=1}^C \sum_{i=1}^N \mathcal{L}(\mathbf{W}\mathbf{h}_{\theta}(\mathbf{x}_{i,c}) + \mathbf{b}, \mathbf{y}_c).$$



# Notation

Train Global Mean

$$\boldsymbol{\mu}_G \triangleq \text{Ave}_{i,c} \{ \mathbf{h}_{i,c} \}$$

objective

$$\min_{\boldsymbol{\theta}, \mathbf{W}, \mathbf{b}} \sum_{c=1}^C \sum_{i=1}^N \mathcal{L}(\mathbf{W} \mathbf{h}_{\boldsymbol{\theta}}(\mathbf{x}_{i,c}) + \mathbf{b}, \mathbf{y}_c).$$

Train class means

$$\boldsymbol{\mu}_c \triangleq \text{Ave}_i \{ \mathbf{h}_{i,c} \}, \quad c = 1, \dots, C,$$

Train total covariance

$$\boldsymbol{\Sigma}_T \triangleq \text{Ave}_{i,c} \left\{ (\mathbf{h}_{i,c} - \boldsymbol{\mu}_G)(\mathbf{h}_{i,c} - \boldsymbol{\mu}_G)^\top \right\},$$

Inter-class COV

$$\boldsymbol{\Sigma}_B \triangleq \text{Ave}_c \left\{ (\boldsymbol{\mu}_c - \boldsymbol{\mu}_G)(\boldsymbol{\mu}_c - \boldsymbol{\mu}_G)^\top \right\},$$

Intra-class COV

$$\boldsymbol{\Sigma}_W \triangleq \text{Ave}_{i,c} \left\{ (\mathbf{h}_{i,c} - \boldsymbol{\mu}_c)(\mathbf{h}_{i,c} - \boldsymbol{\mu}_c)^\top \right\}.$$

# Formalizations

(NC1) Variability collapse:  $\Sigma_W \rightarrow \mathbf{0}$

(Recall what \_\_\_\_ is)

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

(NC2) Convergence to Simplex ETF:

$$\begin{aligned} & | \|\boldsymbol{\mu}_c - \boldsymbol{\mu}_G\|_2 - \|\boldsymbol{\mu}_{c'} - \boldsymbol{\mu}_G\|_2 | \rightarrow 0 \quad \forall c, c' \\ & \langle \tilde{\boldsymbol{\mu}}_c, \tilde{\boldsymbol{\mu}}_{c'} \rangle \rightarrow \frac{C}{C-1} \delta_{c,c'} - \frac{1}{C-1} \quad \forall c, c'. \end{aligned}$$

$$\tilde{\boldsymbol{\mu}}_c = (\boldsymbol{\mu}_c - \boldsymbol{\mu}_G) / \|\boldsymbol{\mu}_c - \boldsymbol{\mu}_G\|_2$$

$$\mathbf{M}^* = \sqrt{\frac{C}{C-1}} \left( \mathbf{I} - \frac{1}{C} \mathbb{H} \mathbb{H}^\top \right),$$

# Formalizations

## (NC3) Convergence to Self-Duality:

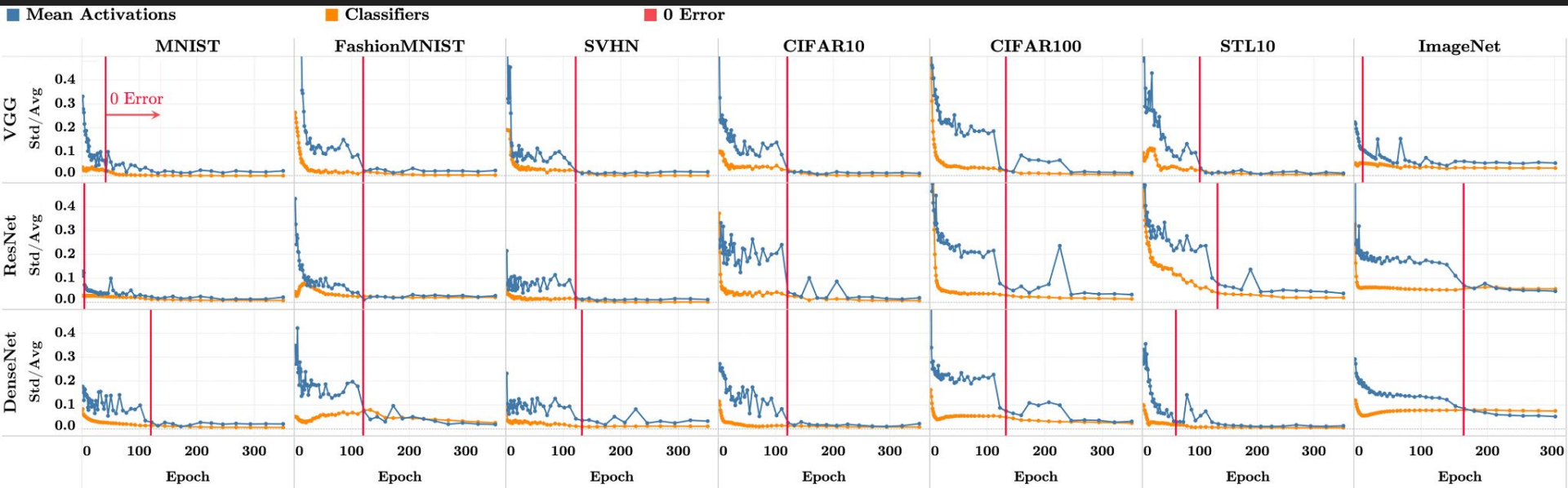
$$\left\| \frac{\mathbf{W}^\top}{\|\mathbf{W}\|_F} - \frac{\dot{\mathbf{M}}}{\|\dot{\mathbf{M}}\|_F} \right\|_F \rightarrow 0.$$

$$\dot{\mathbf{M}} = [\boldsymbol{\mu}_c - \boldsymbol{\mu}_G, c = 1, \dots, C] \in \mathbb{R}^{p \times C}$$

## (NC4): Simplification to Nearest Class Decision:

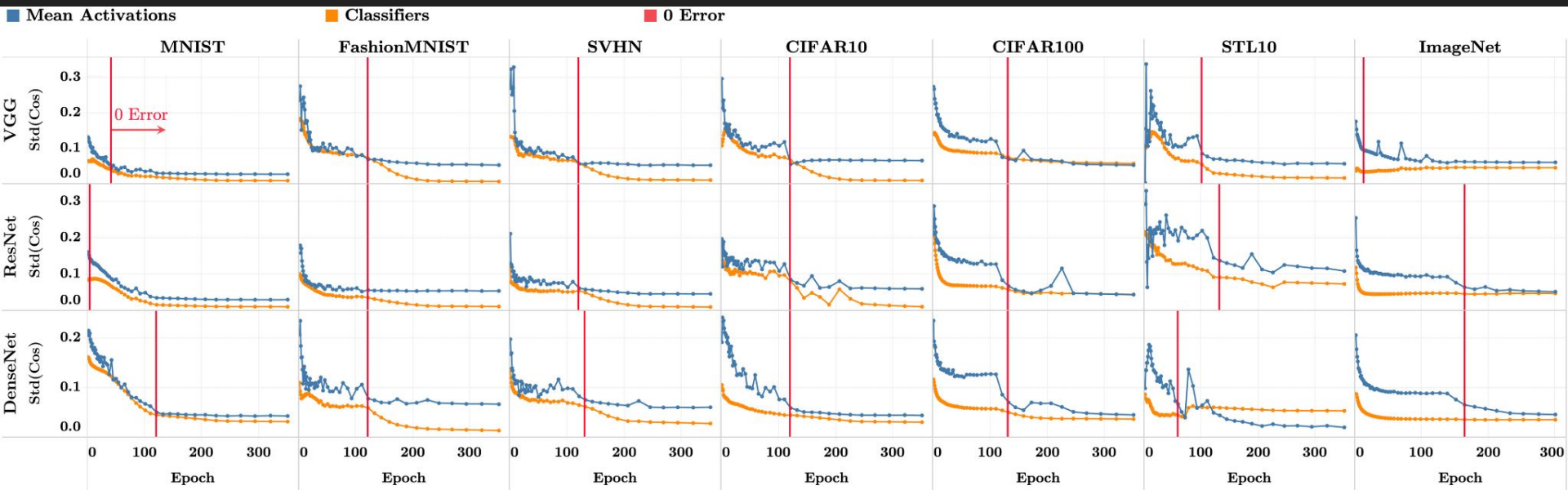
$$\arg \max_{c'} \langle \mathbf{w}_{c'}, \mathbf{h} \rangle + b_{c'} \rightarrow \arg \min_{c'} \|\mathbf{h} - \boldsymbol{\mu}_{c'}\|_2$$

# Figure 2: Means and classifiers become equinorm



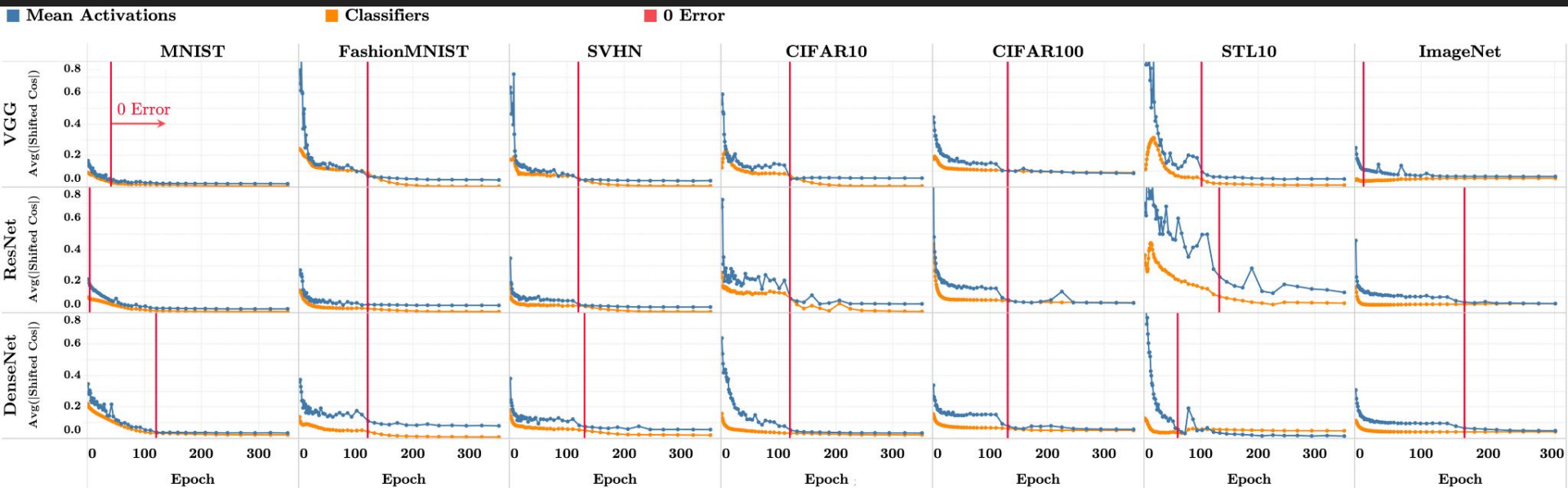
**Train class means become equinorm.** The formatting and technical details are as described in Section 3. In each array cell, the vertical axis shows the coefficient of variation of the centered class-mean norms as well as the network classifiers norms. In particular, the blue lines show  $\text{Std}_c(\|\mu_c - \mu_G\|_2) / \text{Avg}_c(\|\mu_c - \mu_G\|_2)$  where  $\{\mu_c\}$  are the class means of the last-layer activations of the training data and  $\mu_G$  is the corresponding train global mean; the orange lines show  $\text{Std}_c(\|w_c\|_2) / \text{Avg}_c(\|w_c\|_2)$  where  $w_c$  is the last-layer classifier of the  $c$ th class. As training progresses, the coefficients of variation of both class means and classifiers decrease.

# Figures 3,4: Means and classifiers become maximally equiangular



**Classifiers and train class means approach equiangularity.** The formatting and technical details are as described in Section 3. In each array cell, the vertical axis shows the SD of the cosines between pairs of centered class means and classifiers across all distinct pairs of classes  $c$  and  $c'$ . Mathematically, denote  $\cos_{\mu}(c, c') = \langle \mu_c - \mu_G, \mu_{c'} - \mu_G \rangle / (\|\mu_c - \mu_G\|_2 \|\mu_{c'} - \mu_G\|_2)$  and  $\cos_w(c, c') = \langle w_c, w_{c'} \rangle / (\|w_c\|_2 \|w_{c'}\|_2)$  where  $\{w_c\}_{c=1}^C$ ,  $\{\mu_c\}_{c=1}^C$ , and  $\mu_G$  are as in Fig. 2. We measure  $\text{Std}_{c, c' \neq c}(\cos_{\mu}(c, c'))$  (blue) and  $\text{Std}_{c, c' \neq c}(\cos_w(c, c'))$  (orange). As training progresses, the SDs of the cosines approach zero, indicating equiangularity.

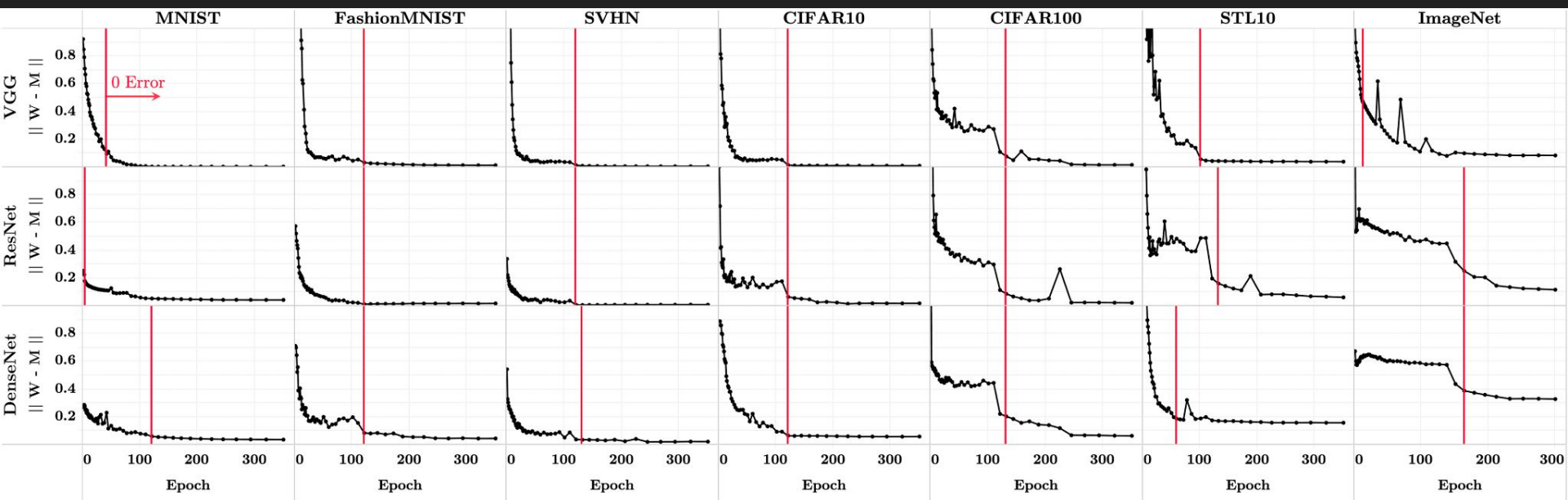
# Figures 3, 4: Means and classifiers become maximally equiangular



**Classifiers and train class means approach maximal-angle equiangularity.** The formatting and technical details are as described in Section 3. We plot in the vertical axis of each cell the quantities  $\text{Avg}_{c,c'} |\cos_{\mu}(c, c') + 1/(C-1)|$  (blue) and  $\text{Avg}_{c,c'} |\cos_w(c, c') + 1/(C-1)|$  (orange), where  $\cos_{\mu}(c, c')$  and  $\cos_w(c, c')$  are as in Fig. 3. As training progresses, the convergence of these values to zero implies that all cosines converge to  $-1/(C-1)$ . This corresponds to the maximum separation possible for globally centered, equiangular vectors.

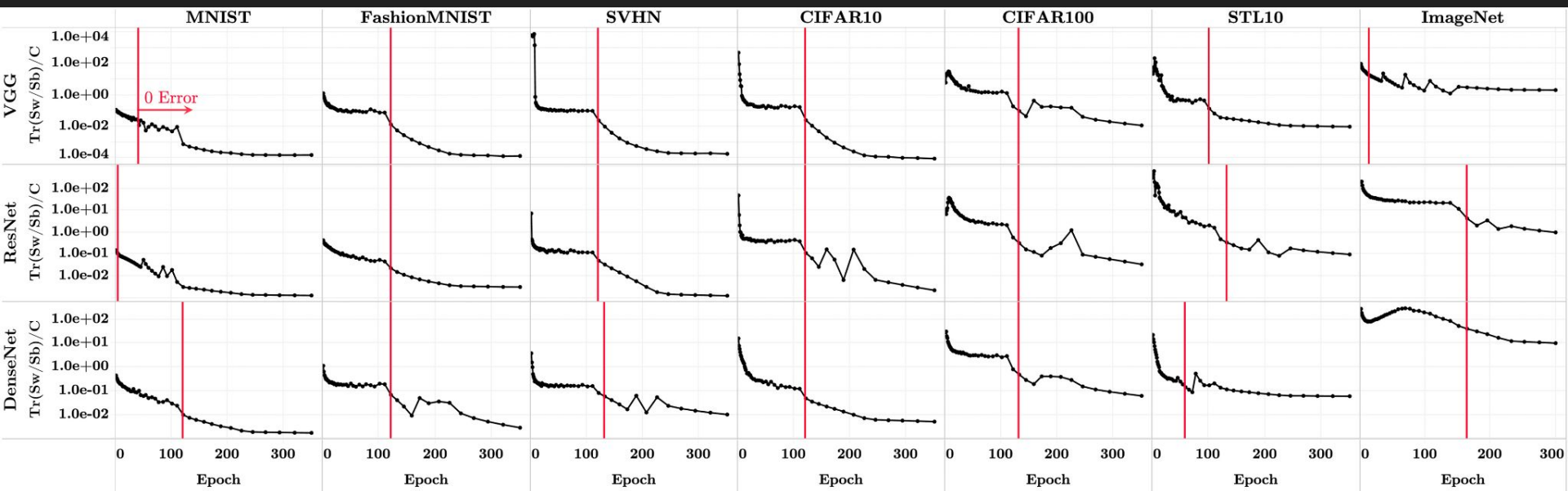


# Figure 5: Means and classifiers become self-dual



**Classifier converges to train class means.** The formatting and technical details are as described in Section 3. In the vertical axis of each cell, we measure the distance between the classifiers and the centered class means, both rescaled to unit norm. Mathematically, denote  $\tilde{\mathbf{M}} = \dot{\mathbf{M}} / \|\dot{\mathbf{M}}\|_F$  where  $\dot{\mathbf{M}} = [\mu_c - \mu_G : c = 1, \dots, C] \in \mathbb{R}^{p \times C}$  is the matrix whose columns consist of the centered train class means; denote  $\tilde{\mathbf{W}} = \mathbf{W} / \|\mathbf{W}\|_F$  where  $\mathbf{W} \in \mathbb{R}^{C \times p}$  is the last-layer classifier of the network. We plot the quantity  $\|\tilde{\mathbf{W}}^T - \tilde{\mathbf{M}}\|_F^2$  on the vertical axis. This value decreases as a function of training, indicating that the network classifier and the centered-means matrices become proportional to each other (self-duality).

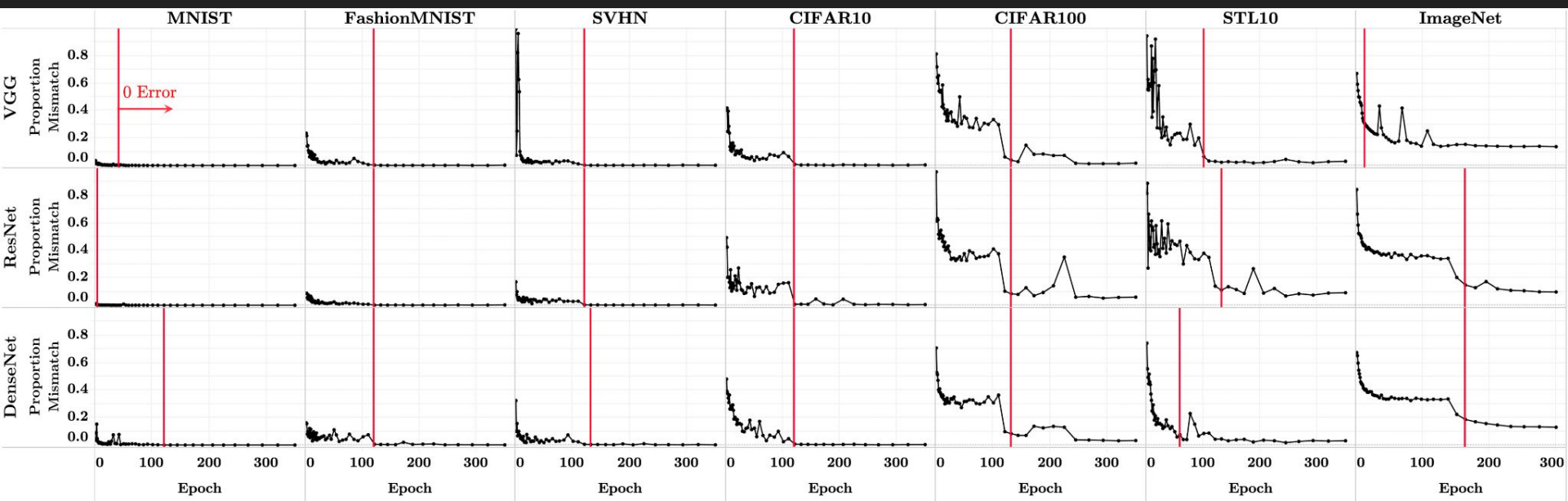
# Figure 6: Train within-class covariance collapses



**Training within-class variation collapses.** The formatting and technical details are as described in Section 3. In each array cell, the vertical axis (log scaled) shows the magnitude of the between-class covariance compared with the within-class covariance of the train activations. Mathematically, this is represented by  $\text{Tr}(\Sigma_W \Sigma_B^\dagger / C)$  where  $\text{Tr}(\cdot)$  is the trace operator,  $\Sigma_W$  is the within-class covariance of the last-layer activations of the training data,  $\Sigma_B$  is the corresponding between-class covariance,  $C$  is the total number of classes, and  $[\cdot]^\dagger$  is Moore–Penrose pseudoinverse. This value decreases as a function of training—indicating collapse of within-class variation.

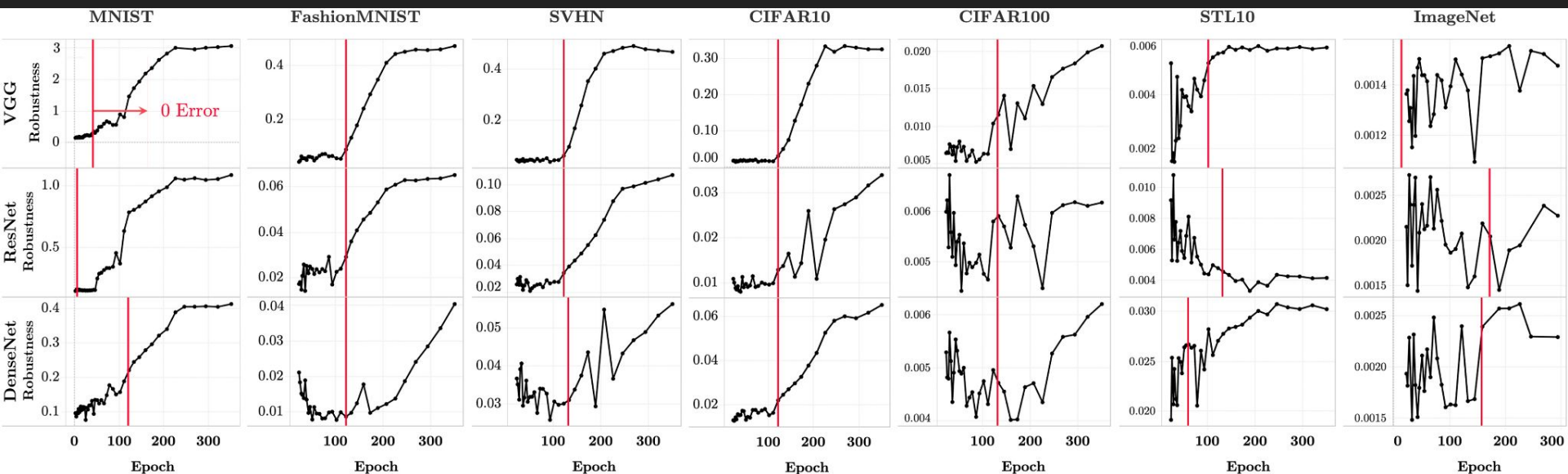


# Figure 7: Classifier approaches NCC



**Classifier behavior approaches that of NCC.** The formatting and technical details are as described in Section 3. In each array cell, we plot the proportion of examples (vertical axis) in the testing set where network classifier disagrees with the result that would have been obtained by choosing  $\arg \min_c \|\mathbf{h} - \mu_c\|_2$  where  $\mathbf{h}$  is a last-layer test activation, and  $\{\mu_c\}_{c=1}^C$  are the class means of the last-layer train activations. As training progresses, the disagreement tends to zero, showing the classifier's behavioral simplification to the nearest train class-mean decision rule.

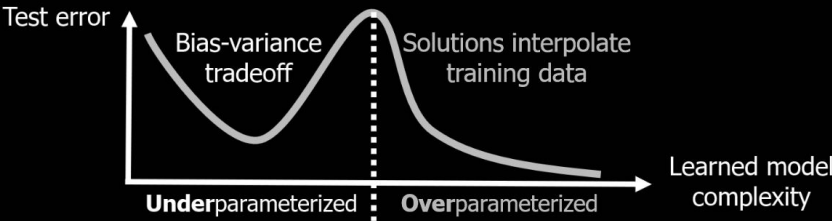
# Fig 8: TPT improves robustness



**Training beyond zero error improves adversarial robustness.** The formatting and technical details are as described in Section 3. For each dataset and network, we sample without replacement 100 test images—constructing for each an adversarial example using the DeepFool method proposed in ref. 23. In each array cell, we plot on the vertical axis the robustness measure,  $\text{Ave}_i \|\mathbf{r}(\mathbf{x}_i)\|_2 / \|\mathbf{x}_i\|_2$ , from the same paper—where  $\mathbf{r}(\mathbf{x}_i)$  is the minimal perturbation required to change the class predicted by the classifier, for a given input image  $\mathbf{x}_i$ . As training progresses, larger  $\mathbf{x}_i$ . As training progresses, larger perturbations are required to fool the deepnet. Across all array cells, the median improvement in the robustness measure in the last epoch over the first epoch achieving zero training error is 0.0252; the mean improvement is 0.2452.

Table 1: TPT improves test error

Dataset and network	Test accuracy at zero error	Test accuracy at last epoch
MNIST		
VGG	99.40	99.56
ResNet	99.32	99.71
DenseNet	99.65	99.70
FashionMNIST		
VGG	92.92	93.31
ResNet	93.29	93.64
DenseNet	94.18	94.35
SVHN		
VGG	93.82	94.53
ResNet	94.64	95.70
DenseNet	95.87	95.93
CIFAR10		
VGG	87.85	88.65
ResNet	88.72	89.44
DenseNet	91.14	91.19
CIFAR100		
VGG	63.03	63.85
ResNet	66.19	66.21
DenseNet	77.19	76.56
STL10		
VGG	65.15	68.00
ResNet	69.99	70.24
DenseNet	67.79	70.81
ImageNet		
VGG	47.26	50.12
ResNet	65.41	64.45
DenseNet	65.04	62.38



# Discussion

Fig2: Training Progress -> Variation

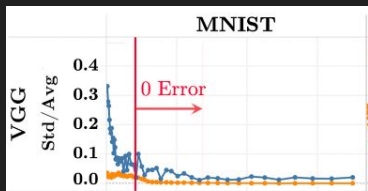


Fig3: All pairs of class means form equal-sized angles

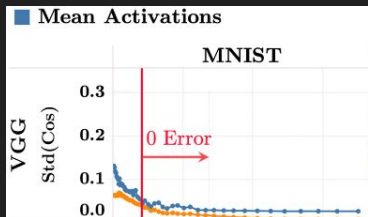


Fig4: Angles Converge to max angle

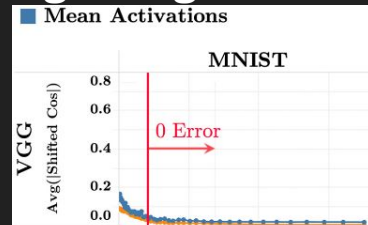


Fig5: Converge up to same ETF (Dual)

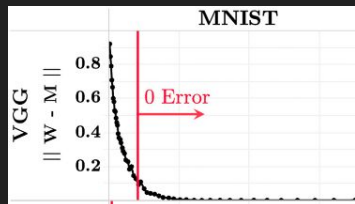


Fig6: Normalized Variation Converges

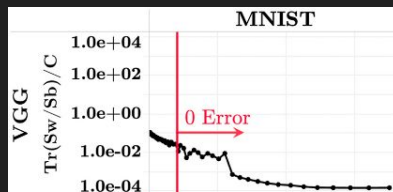
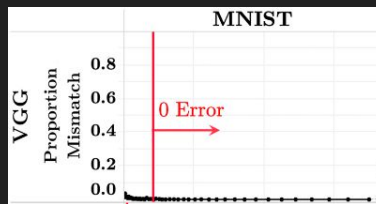
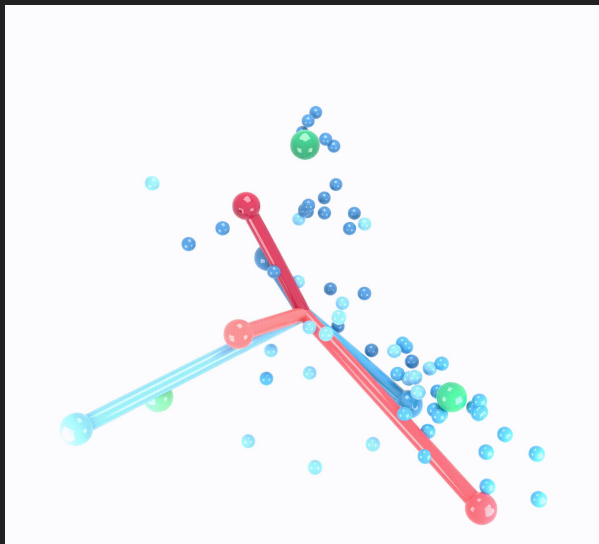


Fig7: Approach to NCC



# Discussion

- Recurring themes:
  - NC continues after zero error
  - This explains TPT - induces significant changes in underlying network structure



# Relation to Previous Work: [Web and Lowe](#)

- “The optimised internal representation of multilayer classifier networks performs nonlinear discriminant analysis”

**Proposition 1** (Section 3 in ref. 28). *Fix the deepnet architecture and the underlying tuning parameters  $\theta$ , so that the activations  $\mathbf{h}_\theta(\mathbf{x})$  involve no training and so that only the classifier weights  $\mathbf{W}$  and biases  $\mathbf{b}$  need to be trained. Maintain the same definitions— $\Sigma_T$ ,  $\mu_c$ ,  $\mu_G$ , etc.—as in Section 2. Adopting the mean squared error loss in place of the cross-entropy loss, the optimal classifier weights and biases are given by*

$$\begin{aligned}\mathbf{W} &= \frac{1}{C} \dot{\mathbf{M}}^\top \Sigma_T^\dagger \\ \mathbf{b} &= \frac{1}{C} \mathbf{1}_C - \frac{1}{C} \dot{\mathbf{M}}^\top \Sigma_T^\dagger \mu_G,\end{aligned}\tag{6}$$

where  $^\dagger$  denotes the Moore–Penrose pseudoinverse,

$$\dot{\mathbf{M}} = [\mu_c - \mu_G, c = 1, \dots, C] \in \mathbb{R}^{p \times C}$$

# Theorem 1

**Theorem 1** [**Proposition 1** +  $\overrightarrow{(\text{NC1} - 2)}$  **ImPLY**  $\overrightarrow{(\text{NC3} - 4)}$ ]. Adopt the framework and assumptions of Proposition 1, as well as the end state implied by  $\overrightarrow{(\text{NC1})}$  and  $\overrightarrow{(\text{NC2})}$  [i.e.,  $\overrightarrow{(\text{NC1})}$  and  $\overrightarrow{(\text{NC2})}$ ]. The Webb-Lowe classifier [6], in this setting, has the additional properties  $\overrightarrow{(\text{NC3})}$  and  $\overrightarrow{(\text{NC4})}$ .

Sketch of Proof:

By prop 1

$$\begin{aligned}\mathbf{W} &= \frac{1}{C} \dot{\mathbf{M}}^\top \boldsymbol{\Sigma}_B^\dagger \\ \mathbf{b} &= \frac{1}{C} \mathbf{1}_C - \frac{1}{C} \dot{\mathbf{M}}^\top \boldsymbol{\Sigma}_B^\dagger \boldsymbol{\mu}_G.\end{aligned}$$

Inter-class variance rewrite as  
Thus

$$\boldsymbol{\Sigma}_B = \frac{1}{C} \dot{\mathbf{M}} \dot{\mathbf{M}}^\top$$

$$\begin{aligned}\mathbf{W} &= \dot{\mathbf{M}}^\top \left( \dot{\mathbf{M}} \dot{\mathbf{M}}^\top \right)^\dagger = \dot{\mathbf{M}}^\dagger \\ \mathbf{b} &= \frac{1}{C} \mathbf{1}_C - \dot{\mathbf{M}}^\top \left( \dot{\mathbf{M}} \dot{\mathbf{M}}^\top \right)^\dagger \boldsymbol{\mu}_G = \frac{1}{C} \mathbf{1}_C - \dot{\mathbf{M}}^\dagger \boldsymbol{\mu}_G.\end{aligned}$$

M has C-1 nonzero SVs so  $\dot{\mathbf{M}}^\dagger = \alpha \dot{\mathbf{M}}^\top$

Therefore, we have

$$\begin{aligned}\mathbf{W} &= \alpha \dot{\mathbf{M}}^\top \\ \mathbf{b} &= \frac{1}{C} \mathbf{1}_C - \alpha \dot{\mathbf{M}}^\top \boldsymbol{\mu}_G;\end{aligned}$$

This is NC3 up to rescaling.

To get NC4, the class predicted by this is

$$= \arg \min_{c'} \|\mathbf{h} - \boldsymbol{\mu}_{c'}\|_2.$$



# Relation to Previous Work: [Soudry et al.](#)

Constrained form of trained classifier

With fixed activations

Inductive bias is reason for success of deep learning

According to this result inductive bias is even more constraining

Modern deepnets produce linear classifiers constrained to simplex ETFs

**Proposition 2 (Theorem 7 in ref. 29).** Let  $\mathbb{R}^{NC \times p}$  denote the vector space spanning all last-layer activation datasets,  $\mathbf{H} = (\mathbf{h}_{i,c} : 1 \leq i \leq N, 1 \leq c \leq C)$ , and let  $\mathcal{H}$  denote the measurable subset of  $\mathbb{R}^{NC \times p}$  consisting of linearly separable datasets: that is, consisting of datasets  $\mathbf{H}$  where, for some linear classifiers  $\{\mathbf{w}_c\}_{c=1}^C$  (possibly depending on dataset  $\mathbf{H}$ ), separability holds:

$$\langle \mathbf{w}_c - \mathbf{w}'_c, \mathbf{h}_{i,c} \rangle \geq 1, \quad \mathbf{h}_{i,c} \in \mathbf{H}.$$

For (Lebesgue) almost every dataset  $\mathbf{H} \in \mathcal{H}$ , gradient descent minimizing the cross-entropy loss, as a function of the classifier weights, tends to a limit. This limit is identical to the solution of the max-margin classifier problem:

$$\min_{\{\mathbf{w}_c\}_{c=1}^C} \sum_{c=1}^C \|\mathbf{w}_c\|_2^2 \text{ s.t. } \forall i, c, c' \neq c : \langle \mathbf{w}_c - \mathbf{w}_{c'}, \mathbf{h}_{i,c} \rangle \geq 1. \quad [9]$$



# Theorem 2

**Theorem 2** (**Proposition 2** +  $\overrightarrow{(\text{NC1} - 2)}$  **ImPLY**  $\overrightarrow{(\text{NC3} - 4)}$ ). *Adopt the framework and assumptions of Proposition 2, as well as the end state implied by  $(\text{NC1})$  and  $(\text{NC2})$  [i.e.,  $\overrightarrow{(\text{NC1})}$  and  $\overrightarrow{(\text{NC2})}$ ]. The Soudry et al. (29) classifier [9], in this setting, has the additional properties  $\overrightarrow{(\text{NC3})}$  and  $\overrightarrow{(\text{NC4})}$ .*

I am not going through the proof since it's long. But the sketch is that we do some linear optimization and algebra to show that the max margin classifier also has NC3 and NC4 properties.

This is similar to result of theorem 1, but now for Cross Entropy Loss

# Deriving ETF Emergence

Abstract Feature Engineer chooses activations to minimize classification error

Assume we are given an observation  $\mathbf{h} = \boldsymbol{\mu}_\gamma + \mathbf{z} \in \mathbb{R}^C$ , where  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$  and  $\gamma \sim \text{Unif}\{1, \dots, C\}$  is an unknown class index, distributed independently from  $\mathbf{z}$ . Our goal is to recover  $\gamma$  from  $\mathbf{h}$ , with as small an error rate as possible. We constrain ourselves to use a linear classifier,  $\mathbf{W}\mathbf{h} + \mathbf{b}$ , with weights  $\mathbf{W} = [\mathbf{w}_c : c = 1, \dots, C] \in \mathbb{R}^{C \times C}$  and biases  $\mathbf{b} = (b_c) \in \mathbb{R}^C$ ; our decision rule is

$$\hat{\gamma}(\mathbf{h}) = \hat{\gamma}(\mathbf{h}; \mathbf{W}, \mathbf{b}) = \arg \max_c \langle \mathbf{w}_c, \mathbf{h} \rangle + b_c.$$

Our task is to design the classifier  $\mathbf{W}$  and bias  $\mathbf{b}$ , as well as a matrix

$\mathbf{M} = [\boldsymbol{\mu}_c : c = 1, \dots, C] \in \mathbb{R}^{C \times C}$ , subject to the norm constraints  $\|\boldsymbol{\mu}_c\|_2 \leq 1$  for all  $c$ .

We can rewrite this as an optimal coding problem. We design a codebook and decoder that allows optimal retrieval of class identity from noisy information.

And we are interested in the large deviations error exponent

$$\beta(\mathbf{M}, \mathbf{W}, \mathbf{b}) = -\lim_{\sigma \rightarrow 0} \sigma^2 \log P_\sigma \{ \hat{\gamma}(\mathbf{h}) \neq \gamma \}$$

# Theorem 3 - Simplex ETF Emergence

We engineer a collection of codebooks to optimize the vanishingly small theoretical misclassification, we obtain the standard simplex ETF or a rotation of it

**Theorem 3.** Under the model assumptions just given in subsections A, B, and C, the optimal error exponent is

$$\begin{aligned}\beta^* &= \max_{\mathbf{M}, \mathbf{W}, \mathbf{b}} \beta(\mathbf{M}, \mathbf{W}, \mathbf{b}) \quad \text{s.t.} \quad \|\boldsymbol{\mu}_c\|_2 \leq 1 \quad \forall c \\ &= \frac{C}{C-1} \cdot \frac{1}{4},\end{aligned}$$

where the maximum is over  $C \times C$  matrices  $\mathbf{M}$  with at most unit-norm columns, and over  $C \times C$  matrices  $\mathbf{W}$  and  $C \times 1$  vectors  $\mathbf{b}$ .

Moreover, denote  $\mathbf{M}^* = \sqrt{\frac{C}{C-1}} \left( \mathbf{I} - \frac{1}{C} \mathbf{1}\mathbf{1}^\top \right)$  (i.e.,  $\mathbf{M}^*$  is the standard simplex ETF). The optimal error exponent is precisely achieved by  $\mathbf{M}^*$ :

$$\beta(\mathbf{M}^*, \mathbf{M}^*, \mathbf{0}) = \beta^*.$$

# Conclusion

- This is a study in TPT, in which a process called Neural Collapse emerges
- The four interconnected phenomena of Neural Collapse are distinct but related.
- Last layer classifier of a trained deepnet exhibits clear mathematical structure: **INDUCTIVE BIAS**
- Last-layer features collapse to an ETF
- Last-layer classifier is equivalent to NCC decision rule
- Standard work flow for empirical deep learning:
  - Series of arbitrary steps that helped win prediction challenge contest
  - Careful analysis was never the point!
  - TPT benefits today standard deep learning training by leading to Neural Collapse
- Doors open to new formal insights

