# Perturbation Analysis of Neural Collapse

Bar-Ilan University

Tom Tirer*, Haoxiang Huang* & Jonathan Niles-Weed

New York University

## The Neural Collapse (NC) Phenomenon

▶ DNN-based classifiers (of $K$ classes) can be typically represented as

$$\psi_{\Theta}(x) = W h_{\theta}(x) + b$$

where $x \in \mathbb{R}^D$ is the sample, $h_{\theta}(\cdot): \mathbb{R}^D \to \mathbb{R}^d$ is the (deep) feature mapping, and $\{W \in \mathbb{R}^{K \times d}, b \in \mathbb{R}^K\}$ is the last layer classifier. Learnable params: $\Theta = \{W, b, \theta\}$.

▶ Common practice: Keep optimizing the network's parameters after the training error vanishes to further push the training loss toward zero.
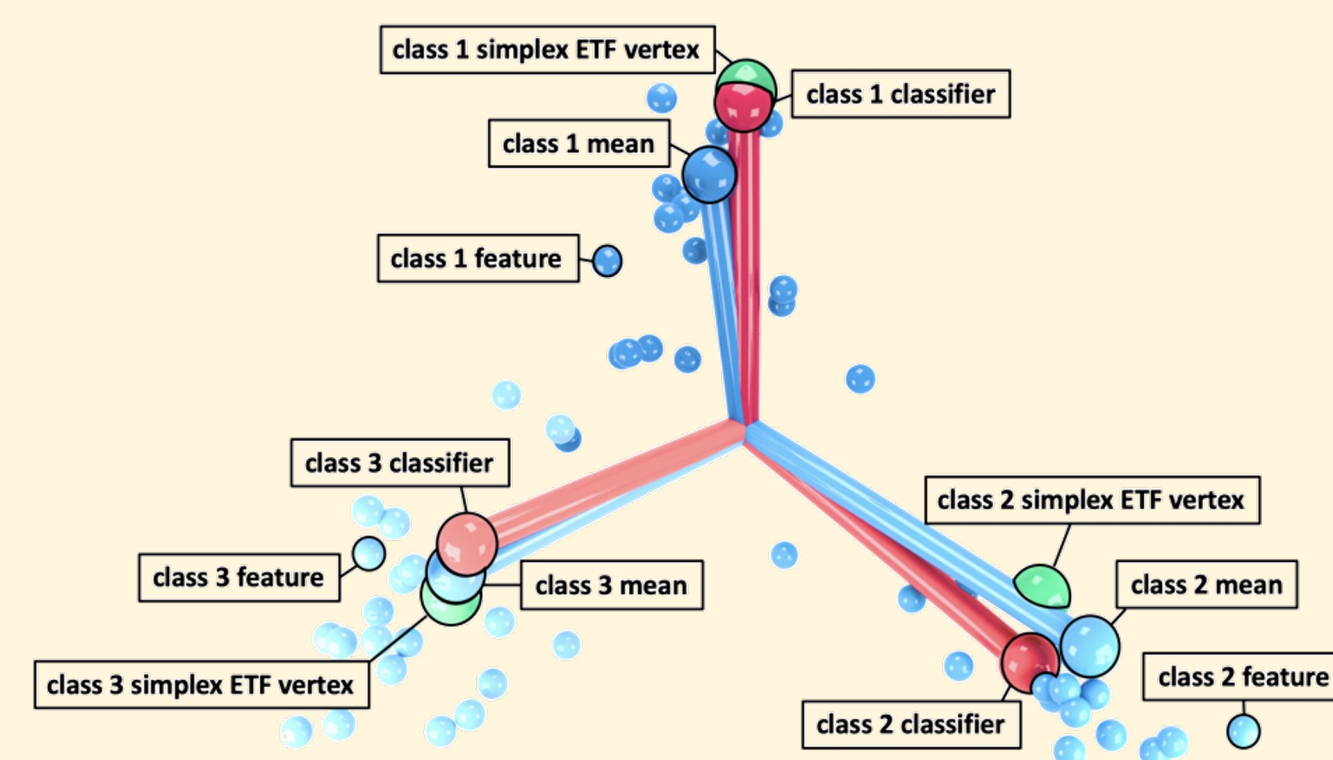
▶ The "Neural Collapse" (NC) phenomenon [Papyan et al. (2020)] has been empirically observed in this phase of training with CE loss (or MSE loss [Han et al. (2022)]): Let $H := [h_{\theta}(x_{1,1}), ..., h_{\theta}(x_{1,n}), ..., ..., h_{\theta}(x_{K,1}), ..., h_{\theta}(x_{K,n})] \in \mathbb{R}^{d \times Kn}$.

  ▸ (NC1): Decrease in within-class variability of features $h_{\theta}(x)$:
  $\|H - \overline{H} \otimes \mathbf{1}_n^\top\|_F$ decreases, where $\overline{H} := [\overline{h}_1, ..., \overline{h}_K] \in \mathbb{R}^{d \times K}$ are classes' mean features

  ▸ (NC2): Increase in the similarity of the mean features to a simplex ETF structure:
  $\left\|(\overline{H} - \overline{h}_G \mathbf{1}_K^\top)^\top (\overline{H} - \overline{h}_G \mathbf{1}_K^\top) - \rho\left(I_K - \frac{1}{K}\mathbf{1}_K\mathbf{1}_K^\top\right)\right\|_F$ decreases, for some $\rho > 0$

  ▸ (NC3): Increase in the alignment of the last weights $W^\top$ and the mean features $\overline{H}$:
  $\left\|W(\overline{H} - \overline{h}_G \mathbf{1}_K^\top) - \tilde{\rho}\left(I_K - \frac{1}{K}\mathbf{1}_K\mathbf{1}_K^\top\right)\right\|_F$ decreases, for some $\tilde{\rho} > 0$

**Empirical observations in practical settings:**

• "NC metrics" typically plateau above zero (even when reducing LR)

• The margin from exact NC depends on the dataset complexity, DNN architecture, hyperparameter tuning, etc.

• Interesting depthwise behavior: gradual reduction of within-class variability (NC1 metric)



## The Unconstrained Features Model (UFM)

▶ The typical way to optimize the DNN's parameters (empirical risk minimization):

$$\min_{\Theta} \frac{1}{Kn} \sum_{k=1}^{K} \sum_{i=1}^{n} \mathcal{L}(W h_{\theta}(x_{k,i}) + b, y_k) + \mathcal{R}(\Theta)$$

where $y_k \in \mathbb{R}^K$ is one-hot vector, $\mathcal{L}(\cdot, \cdot)$ is a loss function (e.g., CE or MSE) and $\mathcal{R}(\cdot)$ is a regularization term (e.g., squared $\ell_2$-norm)

▶ [Mixon et al. (2020)] suggested to explore NC via the Unconstrained Features Model (UFM) – the features $\{h_{k,i} := h_{\theta}(x_{k,i})\}$ are free optimization variables:

$$\min_{W, b, \{h_{k,i}\}} \frac{1}{Kn} \sum_{k=1}^{K} \sum_{i=1}^{n} \mathcal{L}(W h_{k,i} + b, y_k) + \mathcal{R}(W, b, \{h_{k,i}\})$$

▶ Most (if not all) of the existing theoretical works on NC consider UFM settings. The typical result: All the minimizers exhibit exact NC structures (zero NC metrics) with no effect of regularization hyperparameters on the structure

▶ UFMs limitations: cannot explain the aforementioned observations

## This Work Is About:

▶ Exploiting knowledge on gradient dynamics and minimizers of UFMs for studying practical (non-exact) NC behavior.

## Existing and New UFM Results

### Theorem 3.1 in [Tirer & Bruna, 2022] (characterization of minimizers)

Let $d \geq K$, $c := \sqrt{\lambda_H \lambda_W}$ and $\rho := \max\{(1-c)\sqrt{\lambda_W/\lambda_H}, 0\}$. Any global minimizer $(W^*, H^*)$ of

$$\min_{W \in \mathbb{R}^{K \times d}, H \in \mathbb{R}^{d \times Kn}} \mathcal{L}(W, H) := \frac{1}{2Kn}\|WH - Y\|_F^2 + \frac{\lambda_W}{2K}\|W\|_F^2 + \frac{\lambda_H}{2Kn}\|H\|_F^2$$

obeys that $H^* = \overline{H} \otimes \mathbf{1}_n^\top$ for some $\overline{H} := [h_1^*, ..., h_K^*] \in \mathbb{R}^{d \times K}$, $W^{*\top} = \sqrt{\lambda_H/\lambda_W}\overline{H}$, and

$$\overline{H}^\top \overline{H} = \rho I_K \implies (\overline{H} - h_G^* \mathbf{1}_K^\top)^\top (\overline{H} - h_G^* \mathbf{1}_K^\top) = \rho(I_K - \frac{1}{K}\mathbf{1}_K\mathbf{1}_K^\top)$$

▶ **New & useful NC1 metric**: $\boxed{\widetilde{NC}_1(H) := \text{trace}(\Sigma_W(H))/\text{trace}(\Sigma_B(H))}$

$\Sigma_W(H)$ and $\Sigma_B(H)$ are the within- and between-class covariance matrices

▶ More amenable for theoretical analysis than $NC_1(H) := \frac{1}{K}\text{trace}(\Sigma_W(H)\Sigma_B^\dagger(H))$

▶ For fixed $H$, the minimizer w.r.t. $W$: $W^*(H) = YH^\top(HH^\top + n\lambda_W I_d)^{-1}$

▶ [Han et al. (2022)] empirically showed that $\|WH - Y\|_F^2 - \|W^*(H)H - Y\|_F^2$ is small during MSE minimization of practical DNN classifiers

### Theorem (NC1 metric decreases along the gradient flow)

Assume that $\lambda_W > 0$, $\lambda_H \geq 0$, and that $H_0$ is non-collapsed (i.e., $\Sigma_W(H_0) \neq 0$). Then, along the gradient flow: $\frac{dH_t}{dt} = -Kn\nabla\mathcal{L}(W^*(H_t), H_t)$

  ● $\widetilde{NC}_1(H_t)$ strictly decreases along the flow until it reaches zero.
  ● $t \mapsto e^{2\lambda_H t}\text{trace}(\Sigma_W(H_t))$ decreases along the flow. In particular, when $\lambda_H > 0$, $\text{trace}(\Sigma_W(H_t))$ decays exponentially.
  ● $t \mapsto e^{2\lambda_H t}\text{trace}(\Sigma_B(H_t))$ strictly increases along the flow.

▶ We got with minimal assumptions : separation between the behavior of $\Sigma_W$ and $\Sigma_B$ along the flow, $\widetilde{NC}_1 \to 0$ exponentially if $\lambda_H > 0$,

## Analysis of the Near-Collapse Regime

### Theorem (Perturbation analysis around collapse for $\beta \gg 1$)

Let $d > K$, $\lambda_H \lambda_W < 1$, and $H_0 = H^*$ where $(W^*, H^*)$ is a minimizer of $\mathcal{L}$ (i.e., collapsed). Set $\delta H_0$, and let $(\tilde{W}^*, \tilde{H}^*)$ be the minimizer of $f(\cdot, \cdot; \tilde{H}_0 = H_0 + \delta H_0)$. Define $\delta H := \tilde{H}^* - H^*$.
For $\beta \gg \max\{1, \lambda_H\}$ we have (with approximation error of $O(\beta^{-2}, \|\delta H_0\|^2)$)

$$\begin{bmatrix} \text{vec}(\delta H^{(1)}) \\ \vdots \\ \text{vec}(\delta H^{(K)}) \end{bmatrix} \approx \begin{bmatrix} F_{1,1} & \dots & F_{1,K} \\ \vdots & \ddots & \vdots \\ F_{K,1} & \dots & F_{K,K} \end{bmatrix} \begin{bmatrix} \text{vec}(\delta H_0^{(1)}) \\ \vdots \\ \text{vec}(\delta H_0^{(K)}) \end{bmatrix}$$

The $dn \times dn$ blocks have closed-form expressions made of $W^*, H^*, \lambda_W, \lambda_H, \beta$
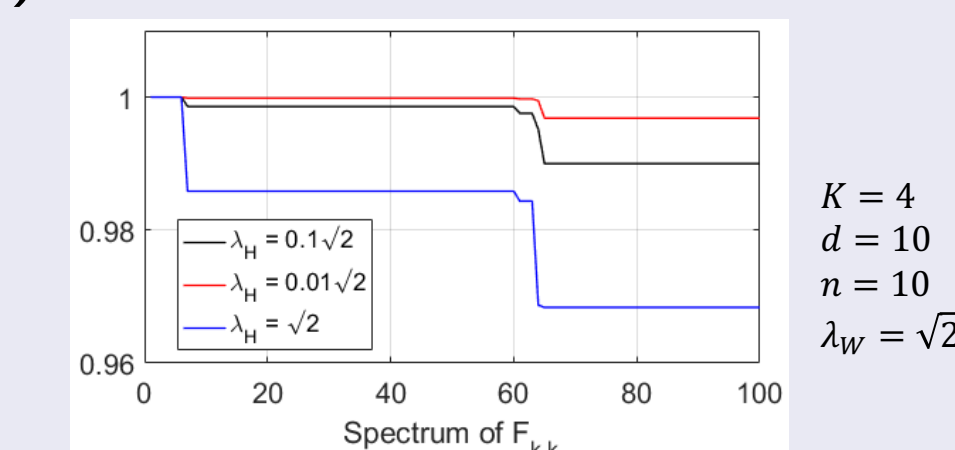
### Theorem (Spectral analysis of inter/intra class blocks)

Consider the setting of the previous theorem and let $k, \tilde{k} \in [K]$ with $k \neq \tilde{k}$. We have that $F_{k,k}$ is full rank, $F_{k,\tilde{k}}$ is rank-1, $\sigma_{max}(F_{k,k}) = 1$ and

$$\sigma_{min}(F_{k,k}) = 1 - \beta^{-1}\sqrt{\lambda_H/\lambda_W}$$

$$\sigma_{max}(F_{k,\tilde{k}}) = 2\beta^{-1}\lambda_H(1 - \sqrt{\lambda_H\lambda_W})$$

(*Actually, we compute the entire spectra)



## New Model: Constraining the UFM

$$\min_{W,H} f(W, H; H_0) := \frac{1}{2Kn}\|WH - Y\|_F^2 + \frac{\lambda_W}{2K}\|W\|_F^2 + \frac{\lambda_H}{2Kn}\|H\|_F^2 + \frac{\beta}{2Kn}\|H - H_0\|_F^2$$

▶ The $\beta \gg 1$ case: can be interpreted as simple architecture between $H_0$ and $H$ that significantly constrains $H$ (e.g., $H_0$ are features one layer before $H$)

▶ Practical DL motivation for $H \approx H_0$: some ResNets, neural ODE, and DEQ

### Corollary (Transferring orthogonal collapse properties from $H_0$)

Let $d \geq K$, $\lambda_H \lambda_W < 1$, and let $(W^*, H^*)$ be a minimizer of $\mathcal{L}(W, H)$. Then, the minimizer of $f(W, H; H_0 = H^*)$ is unique and it is given by $(W^*, H^*)$.

▶ Since we know a lot on $(W^*, H^*)$ minimizer of UFM — we can explore the near-collapse regime via perturbation analysis

▶ First order optimality condition:
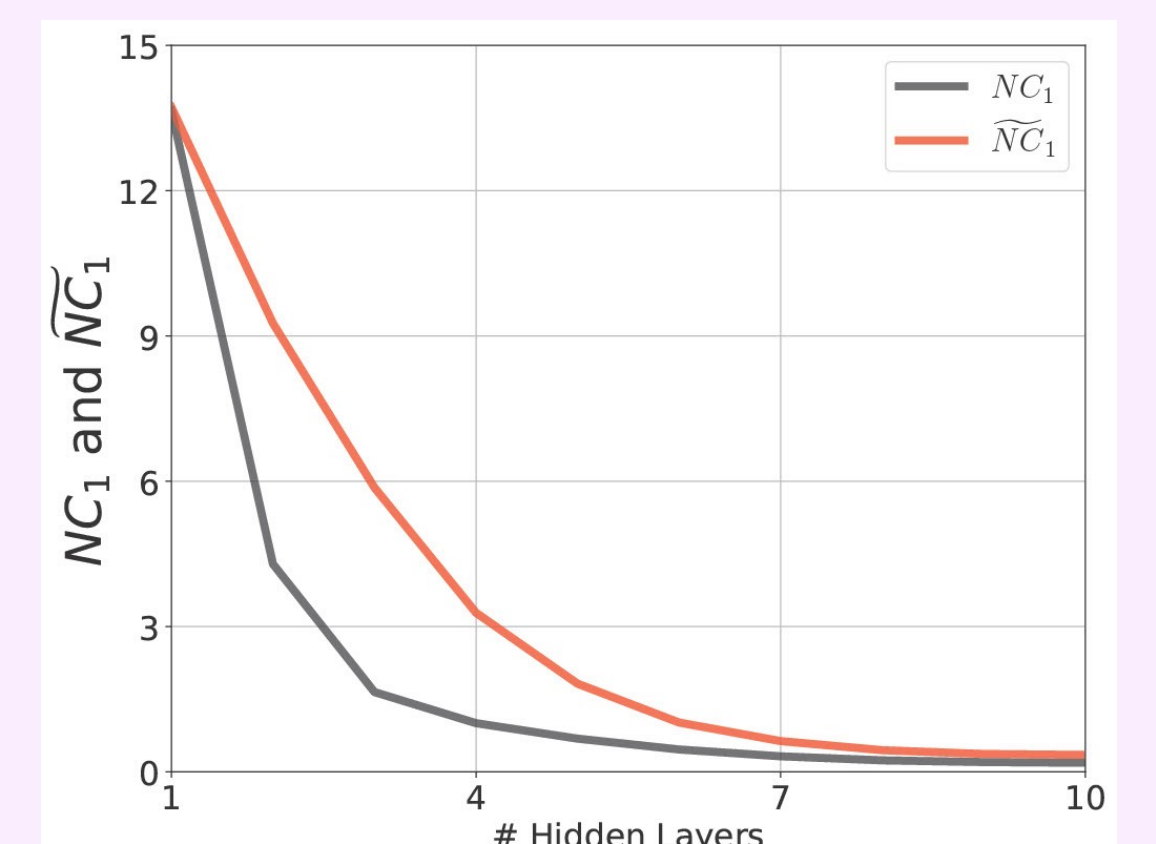$$\frac{H_{1/\beta} - H_0}{1/\beta} = -Kn\nabla\mathcal{L}(W^*(H_{1/\beta}), H_{1/\beta})$$
where $H_{1/\beta} = \min_H f(W^*(H), H; H_0) = \min_H \mathcal{L}(W^*(H), H) + \frac{\beta}{2Kn}\|H - H_0\|_F^2$

### Corollary (Depthwise decrease in NC1 – via gradient flow theory)

Assume that $H_0$ is non-collapsed (i.e., $\Sigma_W(H_0) \neq 0$). For $\beta > C = C(H_0)$, the minimizer of $f$, $H_{1/\beta}$, obeys $\widetilde{NC}_1(H_{1/\beta}) < \widetilde{NC}_1(H_0)$.

▶ Numerical results:

Training an MLP on CIFAR-10 in layer-wise fashion (akin to updating $H_0$ in our model with the previous $H_{1/\beta}$)



▶ Insights gained from the model:

• Increasing $\lambda_H$: increasing the intra-class (diagonal) blocks attenuation
• Increasing $\lambda_W$: increasing the inter-class "interference" blocks attenuation
• Main insight: the intra-class blocks (the effect of perturbation in a certain class in $H_0$ on the features of the same class in $H$) are the dominant. So $\lambda_H$ plays the major role.
• NC1 metric is less affected by the perturbations than other NC metrics (e.g., NC2)

▶ Numerical results: (*more results in the paper, including an "interference" study)

Training ResNet18 on CIFAR-10 with various weight decay (WD) settings – Modifying WD of feature mapping: more deviation from the baseline than modifying WD of last layer