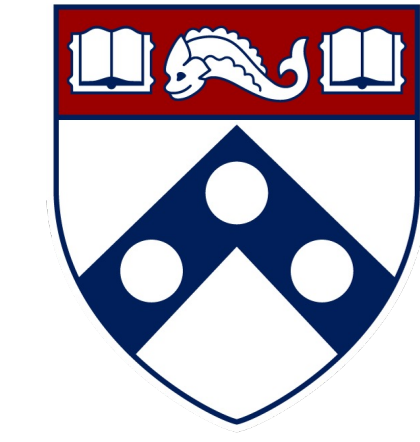


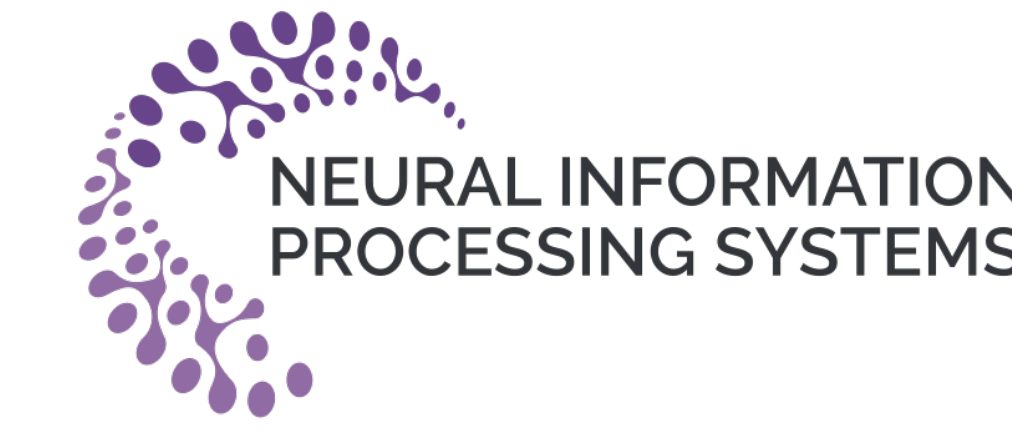
Neural Collapse under Gradient Flow on Shallow ReLU Networks for Orthogonally Separable Data



Hancheng Min
INS&SMS, SJTU

Zhihui Zhu
CSE, OSU

René Vidal
IDEAS, UPenn



INTRODUCTION

- **Neural Collapse (NC)** is a phenomenon where last-layer features and classifiers exhibit a highly structured, symmetric pattern
- Prior work on the **theory of NC** focuses on the **unconstrained feature model**: dynamics are simplified by treating all feature layers as one
- **We prove the emergence of NC for ReLU nets**:
 - With input data, only directional collapse is achieved, instead of singleton collapse
 - Gradient flow with small initialization provably converges to NC solution

PROBLEM

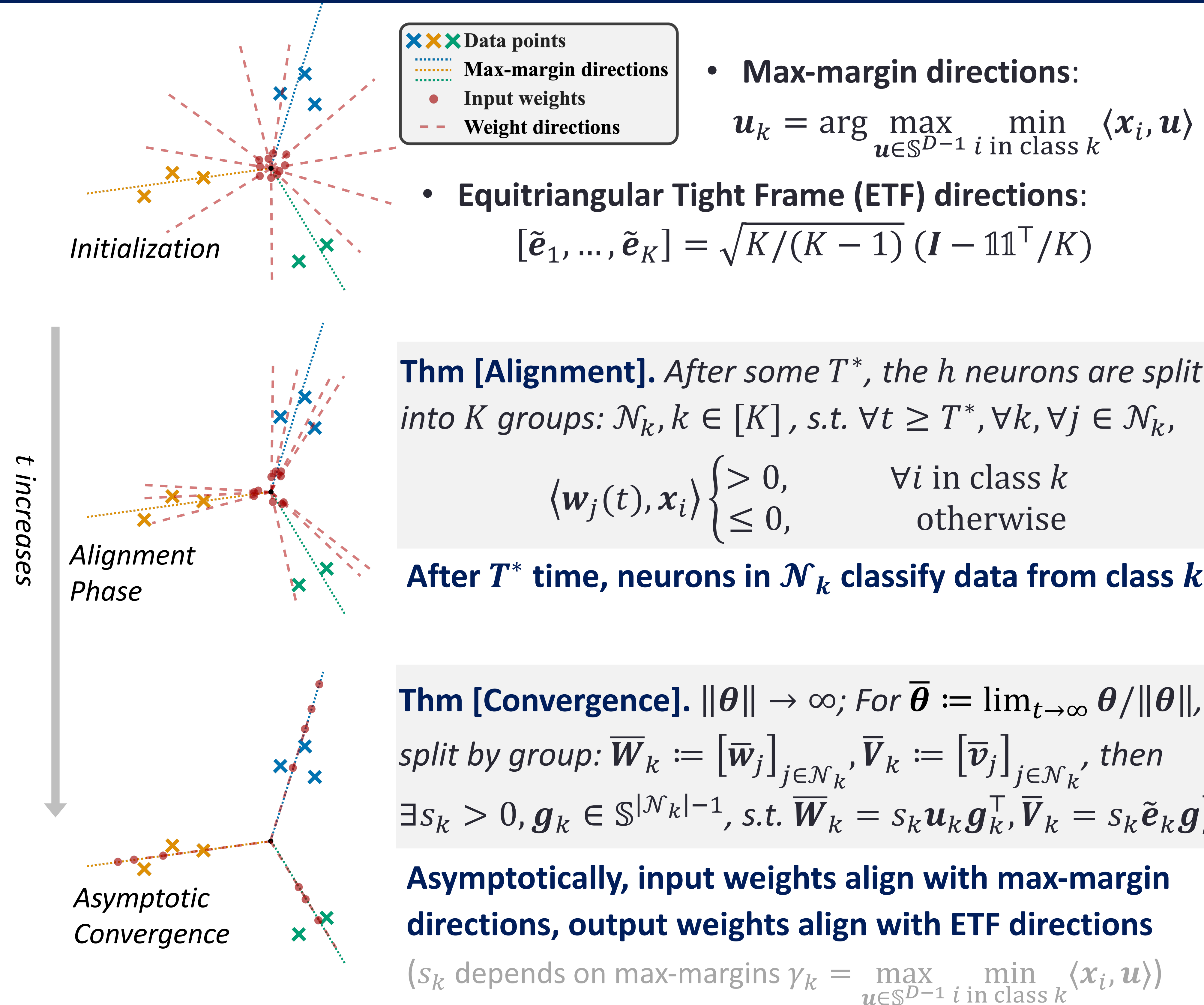
- **Orthogonally separable data from K classes**:

$$\langle x_i, x_j \rangle > 0 \quad \text{if } y_i = y_j \quad \text{where } x_i \in \mathbb{R}^D,$$

$$\langle x_i, x_j \rangle < 0 \quad \text{if } y_i \neq y_j \quad y_i \text{ 1-hot vector}$$
- **ReLU Network**: $f: \mathbb{R}^D \times \Theta \rightarrow \mathbb{R}^K, \theta = \{V, W\}$
 $f(x; \theta) = V\sigma(W^\top x) = \sum_{j=1}^h v_j \sigma(\langle w_j, x \rangle), \sigma: \text{ReLU}$
 - **Input weights**: w_j ; **Output weights**: v_j
 - **Neurons**: $(w_j, v_j), j = 1, \dots, h > K$
 - **Last-layer Feature**: $\phi_\theta(x) = \sigma(W^\top x)$
 - **Last-layer Classifier**: V
- **Cross-Entropy Loss**: $\mathcal{L} = \sum_{i=1}^n \ell_{\text{CE}}(y_i, f(x_i; \theta))$
- **Gradient flow (GF)** with small initialization:

$$\dot{\theta} = -\nabla_{\theta} \mathcal{L}, \quad \|\theta(0)\| \ll 1$$

CONVERGENCE OF GRADIENT FLOW



NEURAL COLLAPSE IN SHALLOW RELU NETWORKS

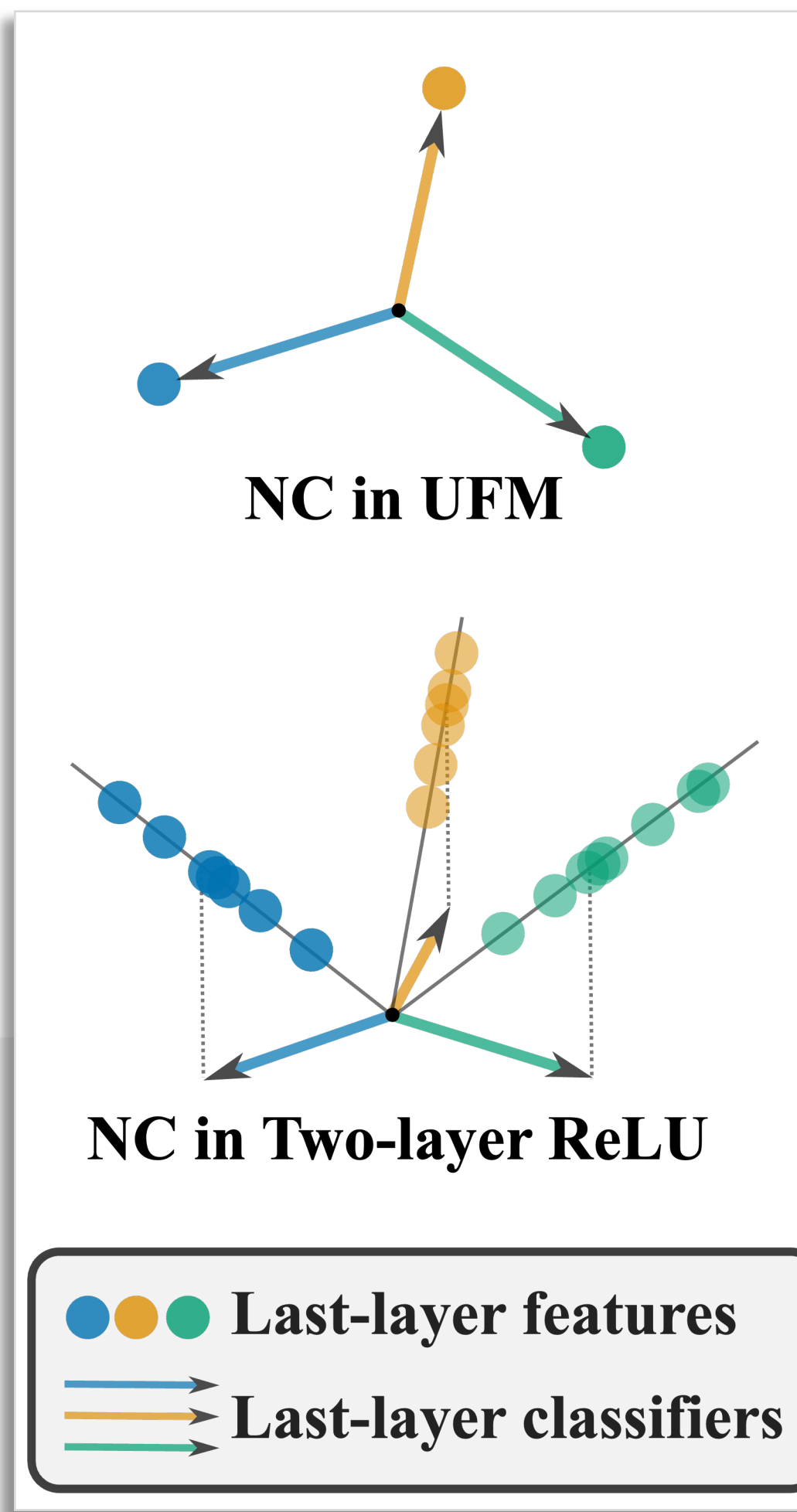
Recall the **NC** for the **unconstrained feature model**:

- *Intra-class variability collapse*: Last-layer features of same class collapse into a singleton
- *Maximal class separation*: Class means form an ETF
- *Self duality*: classifiers align with class means

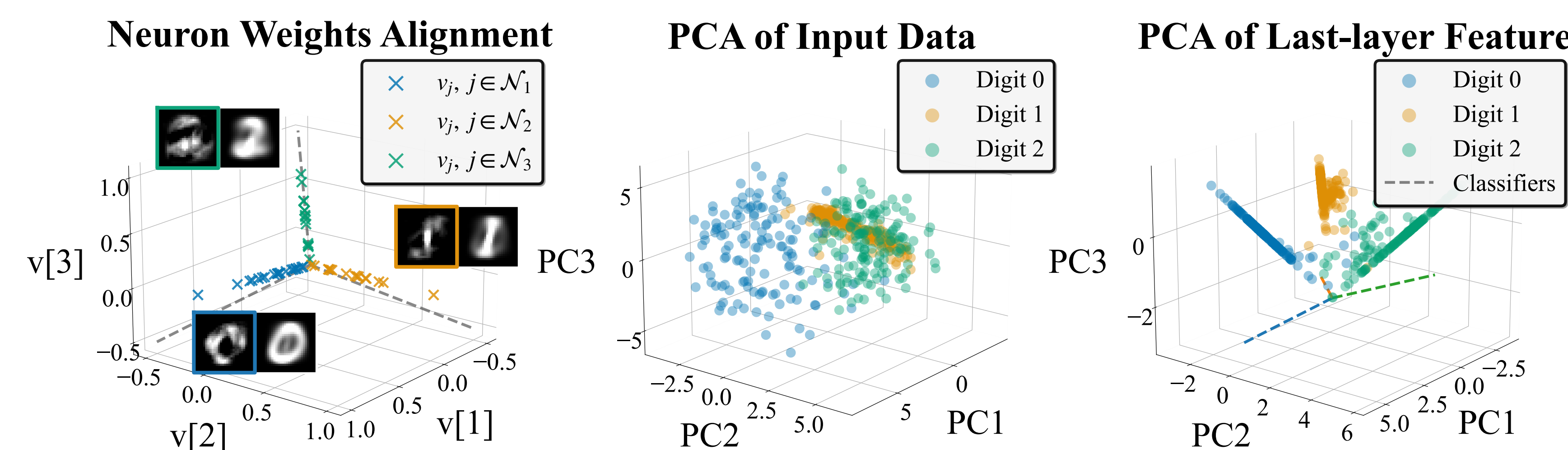
NC for shallow ReLU with orthogonally separable data:

Thm [NC]. For $\bar{\theta} := \lim_{t \rightarrow \infty} \theta / \|\theta\|$, $\exists \bar{\phi}_k \in \mathbb{S}^{h-1}$, s.t.

1. **(Intra-class directional collapse)**
 $\phi_{\bar{\theta}}(x_i) = \langle s_k u_k, x_k \rangle \cdot \bar{\phi}_k, \forall i \text{ in class } k, \forall k \in [K]$
Last-layer features collapse into 1-d subspaces
2. **(Orthogonal class means)**
 $\bar{\phi}_k \geq 0, \langle \bar{\phi}_k, \bar{\phi}_{k'} \rangle = 0, \forall k, k' \in [K], k \neq k'$
Class means form a non-negative orthogonal frame (when normalized)
3. **(Projected self-duality)** $\bar{V} = \sqrt{K/(K-1)} (I - \mathbb{1}\mathbb{1}^\top/K) [s_1 \bar{\phi}_1, \dots, s_K \bar{\phi}_K]^\top$
classifiers align with centered class means



EXPERIMENT ON MNIST DIGITS



Training a shallow ReLU net for classifying digits 0, 1, 2

- Output weights align with ETF directions (determines neurons' group)
- Average input weight per group aligns with average digit image
- PCA of raw data X has appr. error of 61%
- PCA of collapsed features $\phi_\theta(X)$ has appr. error of 0.2%



Full Paper