

# Why Adversarial Interaction Creates Non-Homogeneous Patterns: A Pseudo-Reaction Diffusion Model for Turing Instability

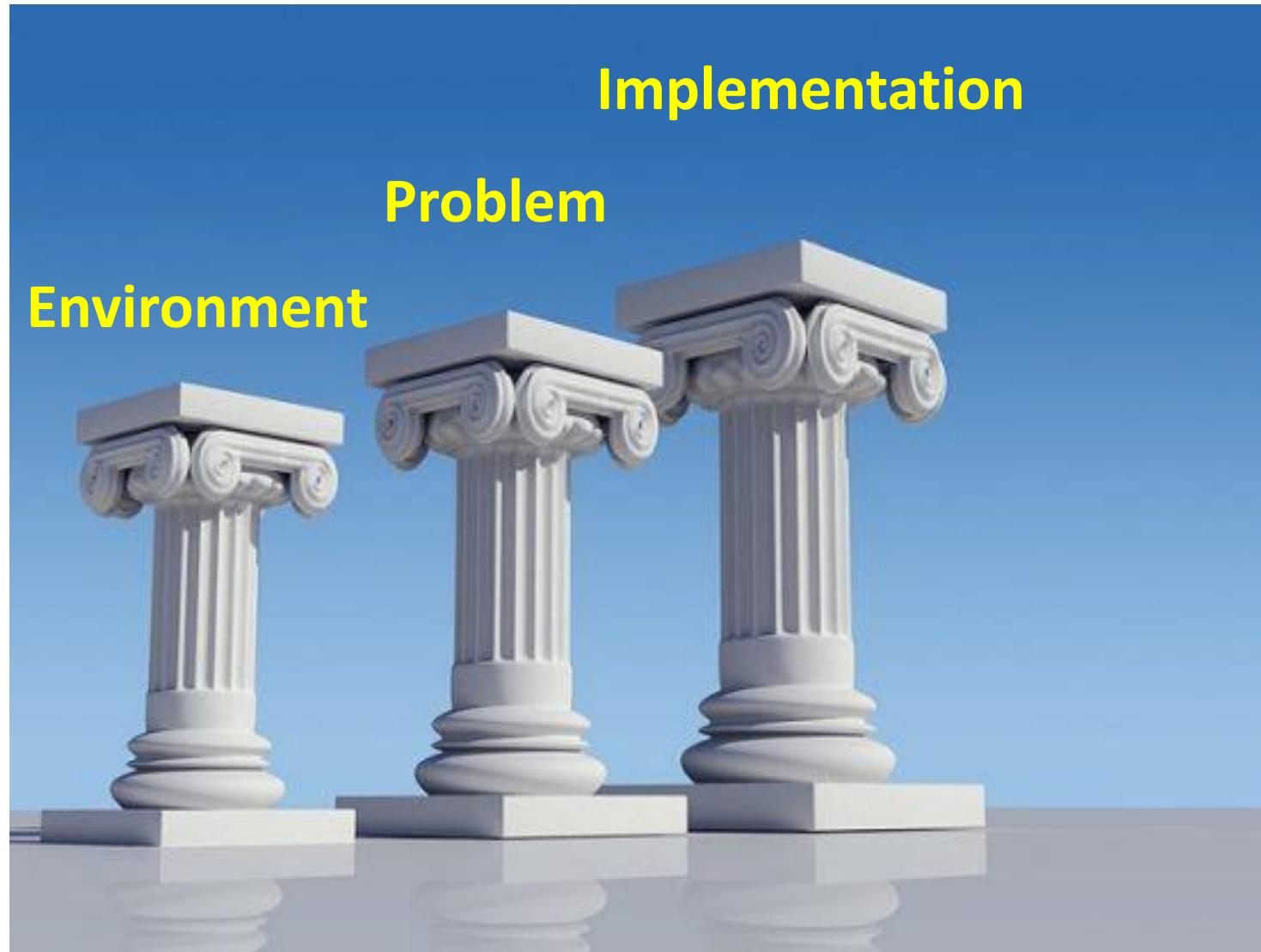
**Litu Rout**

Association for the Advancement of Artificial Intelligence (AAAI-21)

# Why Adversarial Interaction Creates Non-Homogeneous Patterns: A Pseudo-Reaction Diffusion Model for Turing Instability

- Adversarial Interaction
  - Generative Adversarial Networks (GANs)
  - Application of conditional GANs
- Non-Homogeneous Patterns
  - Homogeneous patterns
  - Supervised learning
- Reaction-Diffusion
  - Turing's RD model (1952)
  - Gray-Scott RD model (1984)
- Turing Instability
  - Reaction dynamics
  - Diffusion dynamics

# Three Pillars of Deep Learning



# Why Adversarial Interaction Creates Non-Homogeneous Patterns: A Pseudo-Reaction Diffusion Model for Turing Instability

- Adversarial Interaction
  - Generative Adversarial Networks (GANs)
  - Application of conditional GANs
- Non-Homogeneous Patterns
  - Homogeneous patterns
  - Supervised learning
- Reaction-Diffusion
  - Turing's RD model (1952)
  - Gray-Scott RD model (1984)
- Turing Instability
  - Reaction dynamics
  - Diffusion dynamics

# Outline

- Maximum Likelihood Estimation
- Reverse KL Divergence
- Introduction to GANs
- GANs in Remote Sensing
  - Problem Formulation

# Why Adversarial Interaction Creates Non-Homogeneous Patterns: A Pseudo-Reaction Diffusion Model for Turing Instability

- Adversarial Interaction
  - Generative Adversarial Networks (GANs)
  - Application of conditional GANs
- Non-Homogeneous Patterns
  - Homogeneous patterns
  - Supervised learning
- Reaction-Diffusion
  - Turing's RD model (1952)
  - Gray-Scott RD model (1984)
- Turing Instability
  - Reaction dynamics
  - Diffusion dynamics

# Understanding the Role of Adversarial Regularization in Supervised Learning

- Mitigating Vanishing Gradient
- Asymptotic Iteration Complexity
- Sub-optimality Gap
- Provable Convergence
- Generalization Error (Open Problems)
- Neural Topology Analysis (Open Problems)

# Problem Setup

- Supervised Learning

Input:  $x$

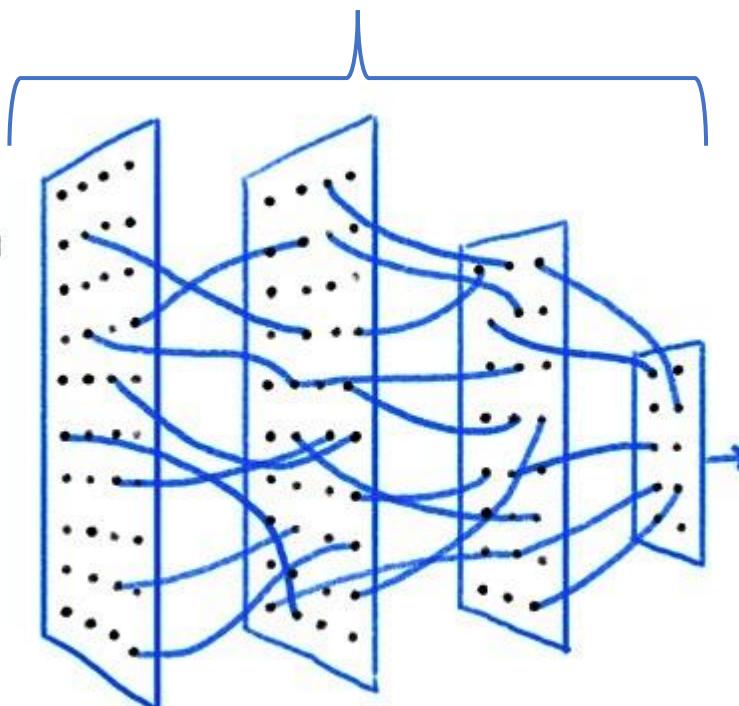
Cat



Dog



$$f(\theta; x)$$



Output:  $y$

$$\arg \min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{P}} [l(f(\theta; x); y)]$$

# Problem Setup

- Adversarial Learning

Input:  $x$

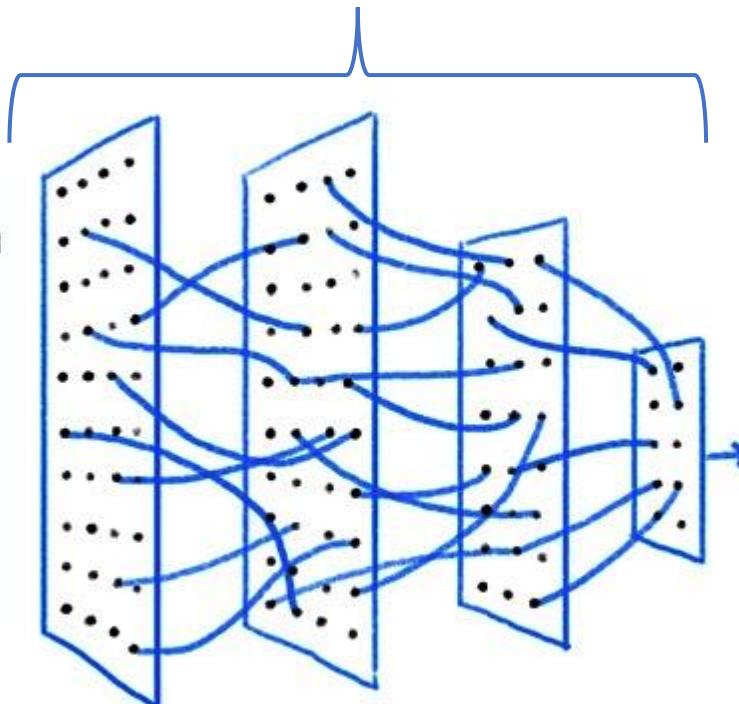
Cat



Dog



$$f(\theta; x)$$



Output:  $y$

# Problem Setup

- Adversarial Learning

Input:  $x$

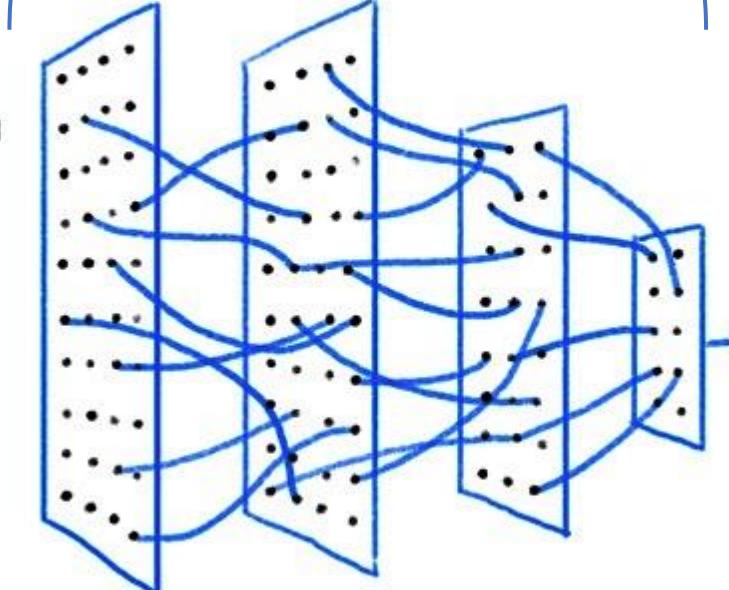
Cat



Dog



$$f(\theta; x)$$



Output:  $y$

Detective  
(Discriminator:  $g(\cdot)$ )

$$\arg \min_{\theta} -\mathbb{E}_{x \sim \mathcal{P}_X} [g(\psi; f(\theta; x))]$$

Detective  
(Discriminator:  $g(\cdot)$ )

# Problem Setup

- Adversarial Regularization

Input:  $x$

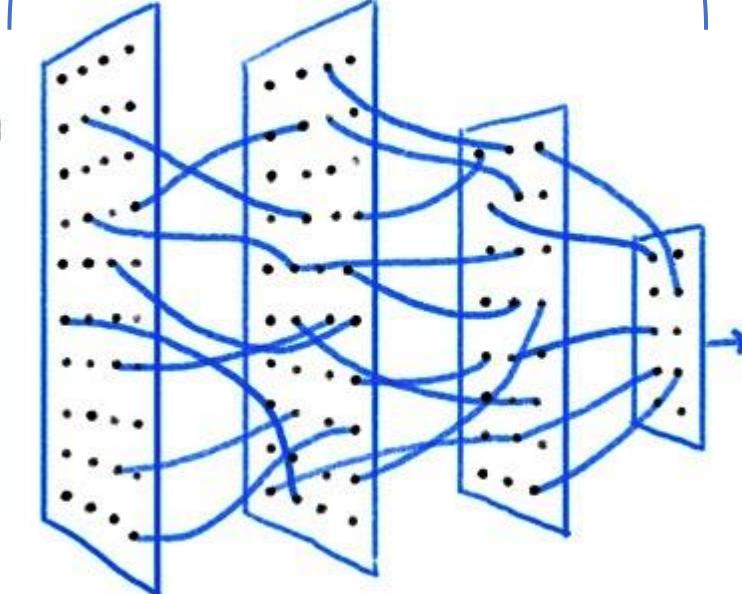
Cat



Dog



$$f(\theta; x)$$



Output:  $y$

$$\arg \min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{P}} [l(f(\theta; x); y) - g(\psi; f(\theta; x))]$$

# Problem Setup

- Supervised Learning

$$\arg \min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{P}} [l(f(\theta; x); y)]$$

# Problem Setup

- Supervised Learning

$$\arg \min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{P}} [l(f(\theta; x); y)]$$

- Adversarial Learning

$$\arg \min_{\theta} -\mathbb{E}_{x \sim \mathcal{P}_X} [g(\psi; f(\theta; x))]$$

# Problem Setup

- Supervised Learning

$$\arg \min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{P}} [l(f(\theta; x); y)]$$

- Adversarial Learning

$$\arg \min_{\theta} -\mathbb{E}_{x \sim \mathcal{P}_X} [g(\psi; f(\theta; x))]$$

- Adversarial Regularization

$$\arg \min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{P}} [l(f(\theta; x); y) - g(\psi; f(\theta; x))]$$

# Problem Setup

- Supervised Learning

$$\arg \min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{P}} [l(f(\theta; x); y)]$$

- Adversarial Learning

$$\arg \min_{\theta} -\mathbb{E}_{x \sim \mathcal{P}_X} [g(\psi; f(\theta; x))]$$

- Adversarial Regularization

$$\arg \min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{P}} [l(f(\theta; x); y) - g(\psi; f(\theta; x))]$$

# Understanding the Role of Adversarial Regularization in Supervised Learning

- Mitigating Vanishing Gradient
- Asymptotic Iteration Complexity
- Sub-optimality Gap
- Provable Convergence
- Generalization Error (Open Problems)
- Neural Topology Analysis (Open Problems)

# Mitigating Vanishing Gradient

## Supervised Objective

$$\theta(t + 1) = \theta(t) - \eta \frac{\partial L(\theta)}{\partial \theta(t)}$$

**Assumption 1.** *The function  $f(\theta; x)$  is  $L$ -Lipschitz in  $\theta$ .*

**Assumption 2.** *The loss function  $l(p; y)$ , where  $p = f(\theta; x)$ , is  $\beta$ -smooth in  $p$ .*

**Lemma 1.** *Let Assumption 1 and Assumption 2 hold. If  $\|\theta - \theta^*\| \leq \epsilon$ , then  $\|\nabla_\theta \mathbb{E}_{(x,y) \sim \mathcal{P}} [l(f(\theta; x); y)]\| \leq L^2 \beta \epsilon$ .*

# Mitigating Vanishing Gradient

## Augmented Objective

$$\theta(t + 1) = \theta(t) - \eta \frac{\partial L(\theta)}{\partial \theta(t)}$$

### Theorem 1 (Informal).

*Under valid assumptions, the adversarial regularization mitigates vanishing gradient in the near optimal region.*

**Theorem 1.** Let us suppose **Assumption 1** and **Assumption 2** hold. If  $\|\theta - \theta^*\| \leq \epsilon$  and  $\|g - g^*\| \leq \delta$ , then  $\|\nabla_\theta \mathbb{E}_{(x,y) \sim \mathcal{P}} [l(f(\theta; x); y) - g(\psi; f(\theta; x))] \| \leq (L^2 \beta \epsilon + L \delta)$ .

# Mitigating Vanishing Gradient

## Experimental Result

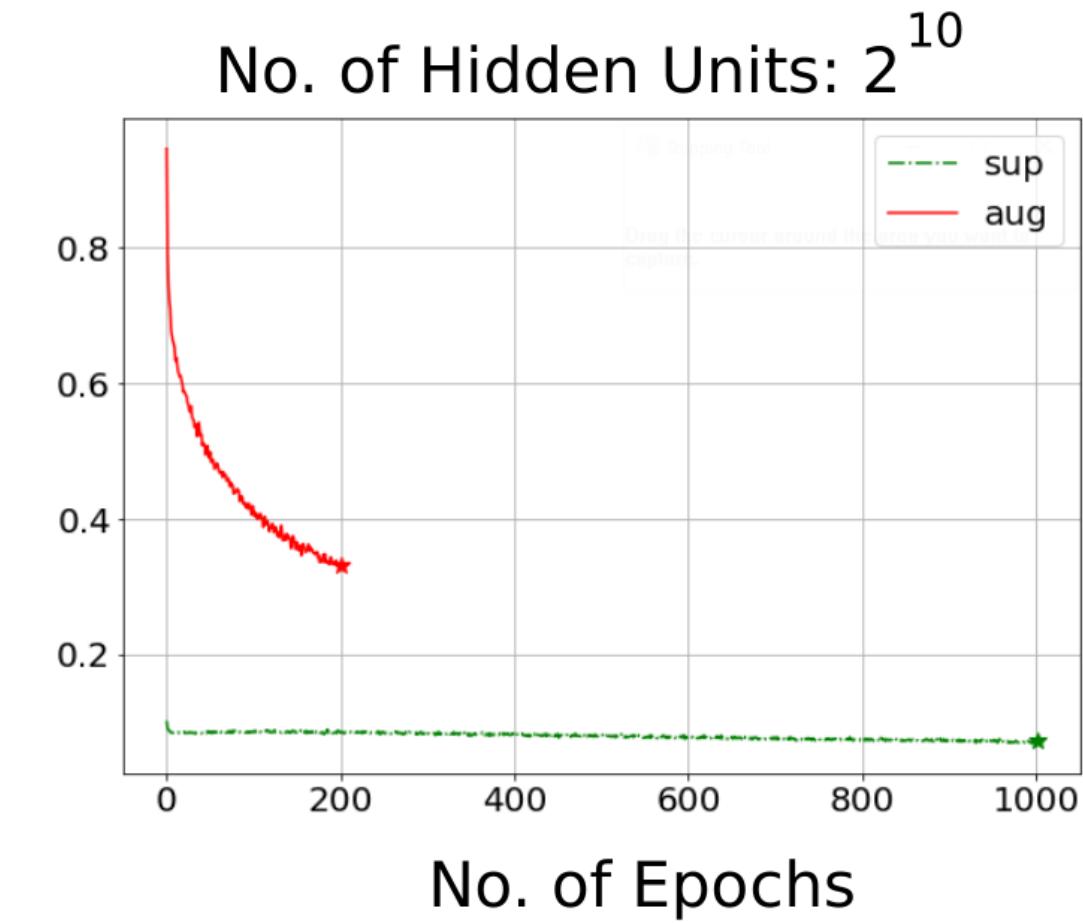
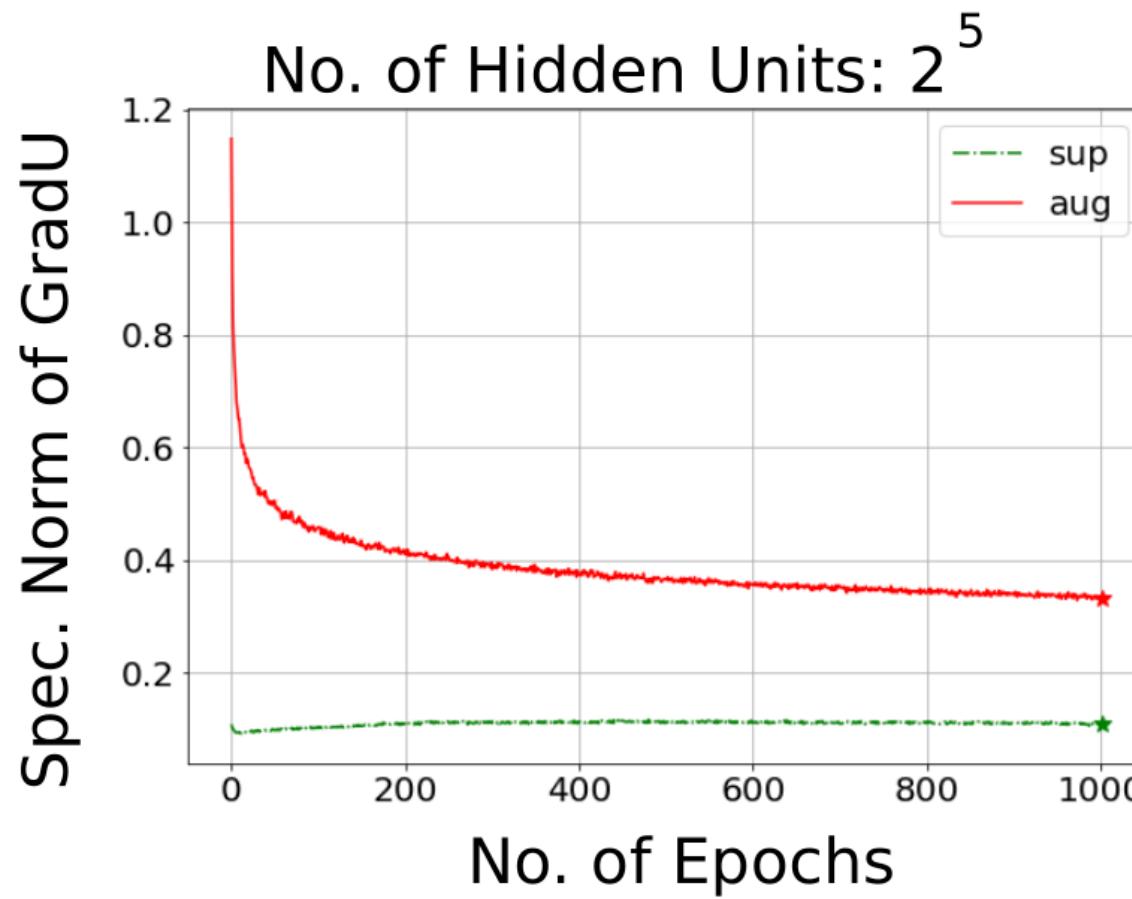
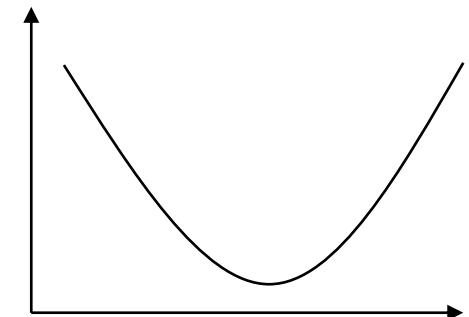
- Gradient norm in hidden layer
- Gradient norm in top layer
- Application of adversarial regularization

# Mitigating Vanishing Gradient

## Experimental Result

- Gradient norm in hidden layer
- Gradient norm in top layer
- Application of adversarial regularization

# Mitigating Vanishing Gradient

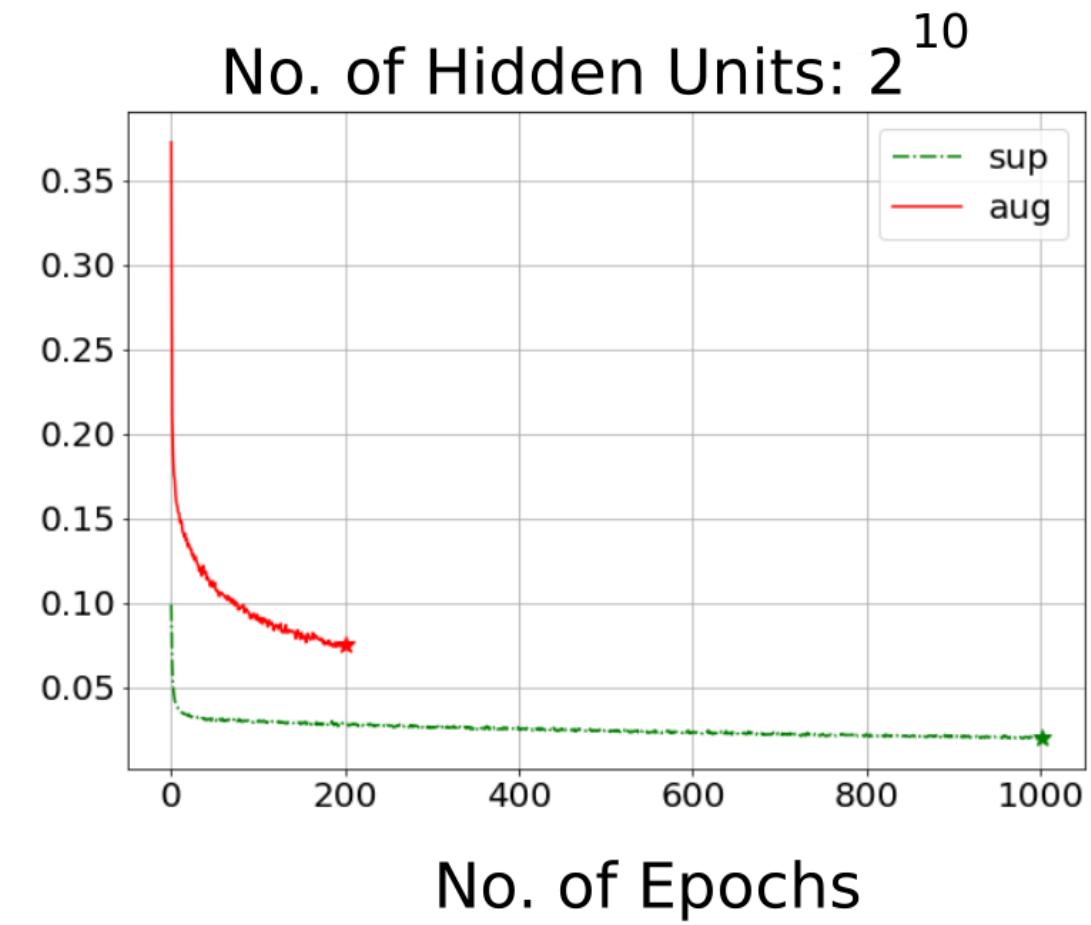
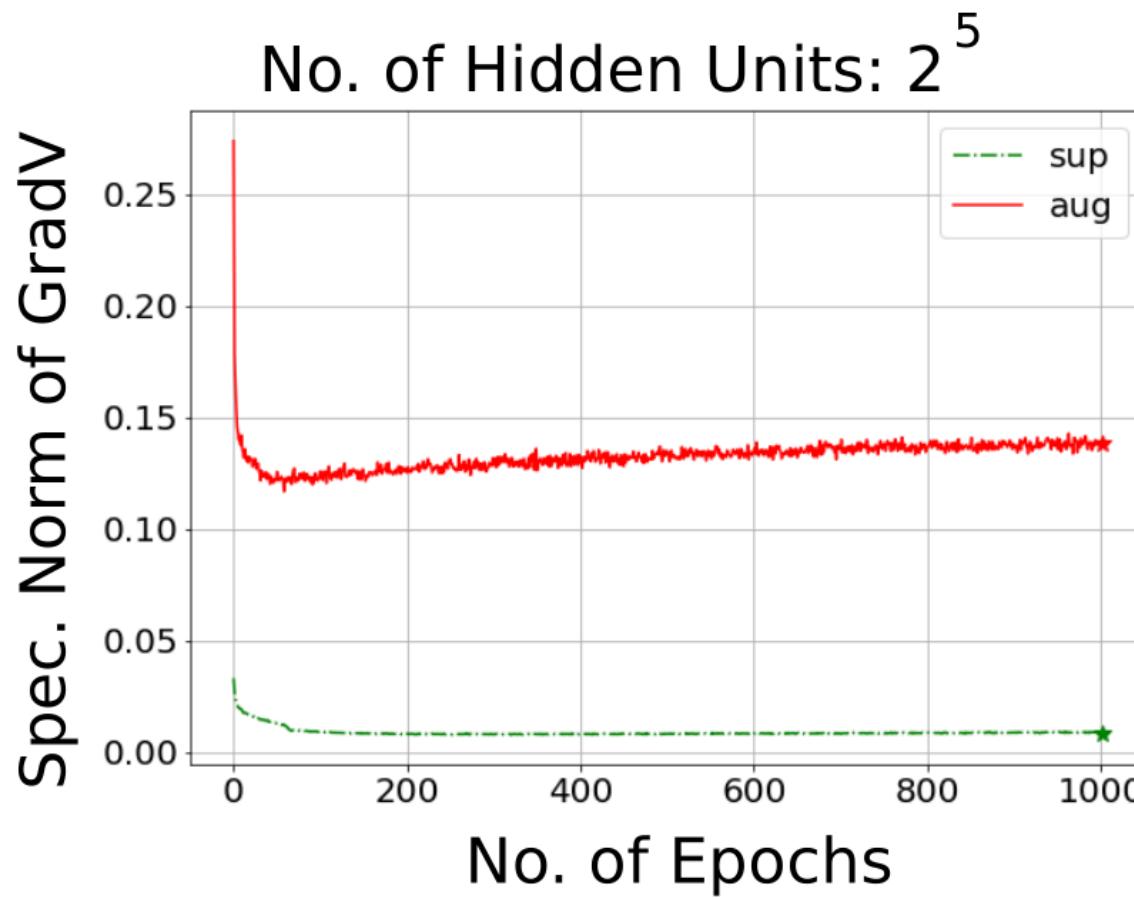
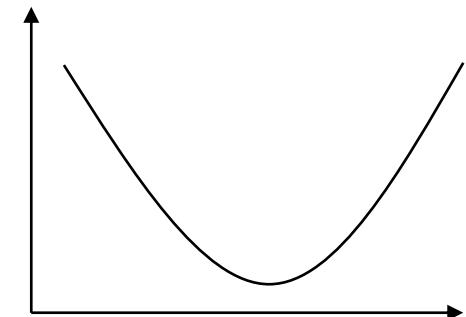


# Mitigating Vanishing Gradient

## Experimental Result

- Gradient norm in hidden layer
- Gradient norm in top layer
- Application of adversarial regularization

# Mitigating Vanishing Gradient



# Mitigating Vanishing Gradient

## Experimental Result

- Gradient norm in hidden layer
- Gradient norm in top layer
- Application of adversarial regularization

# Mitigating Vanishing Gradient

bicubic  
(21.59dB/0.6423)



SRResNet  
(23.53dB/0.7832)



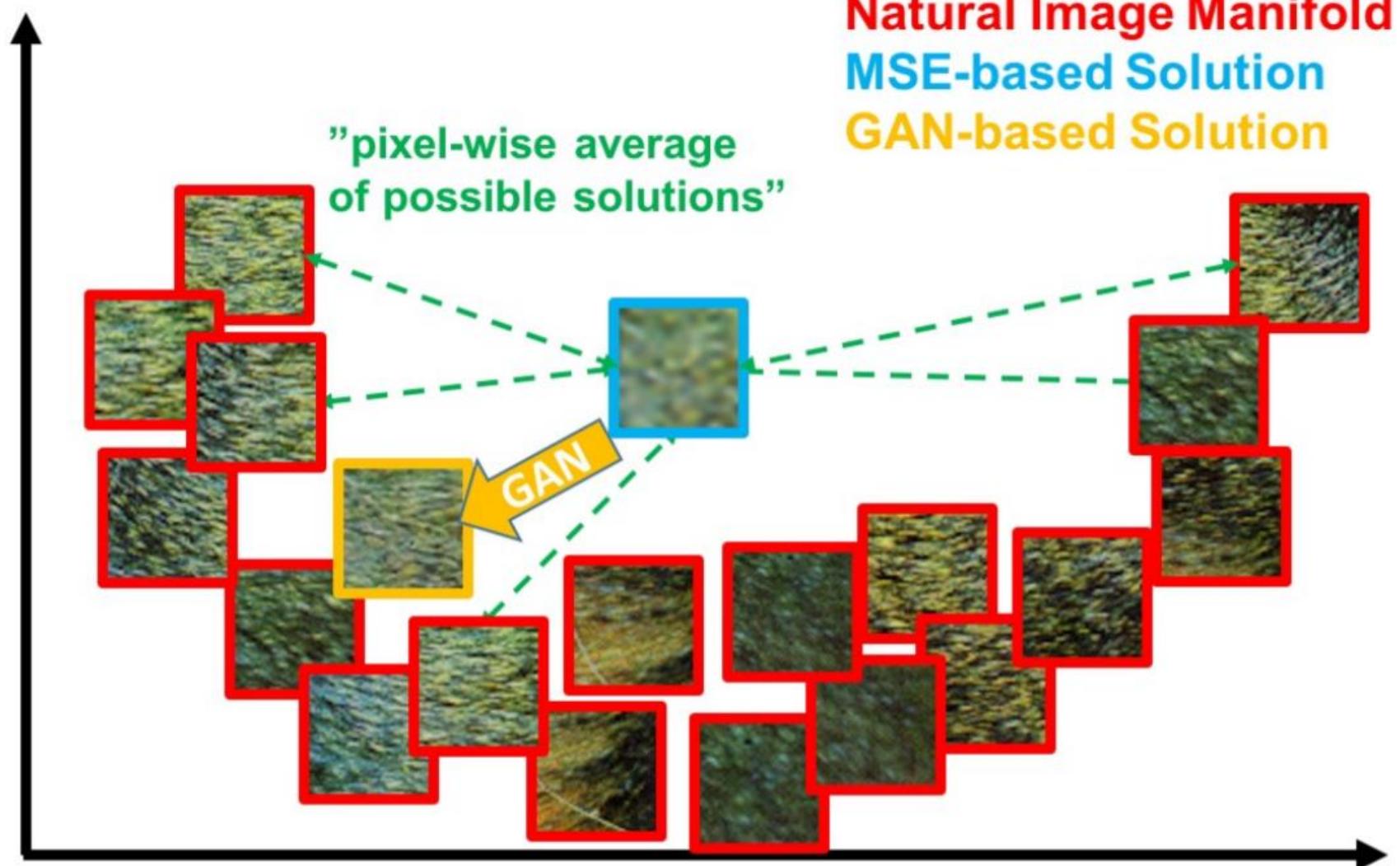
SRGAN  
(21.15dB/0.6868)



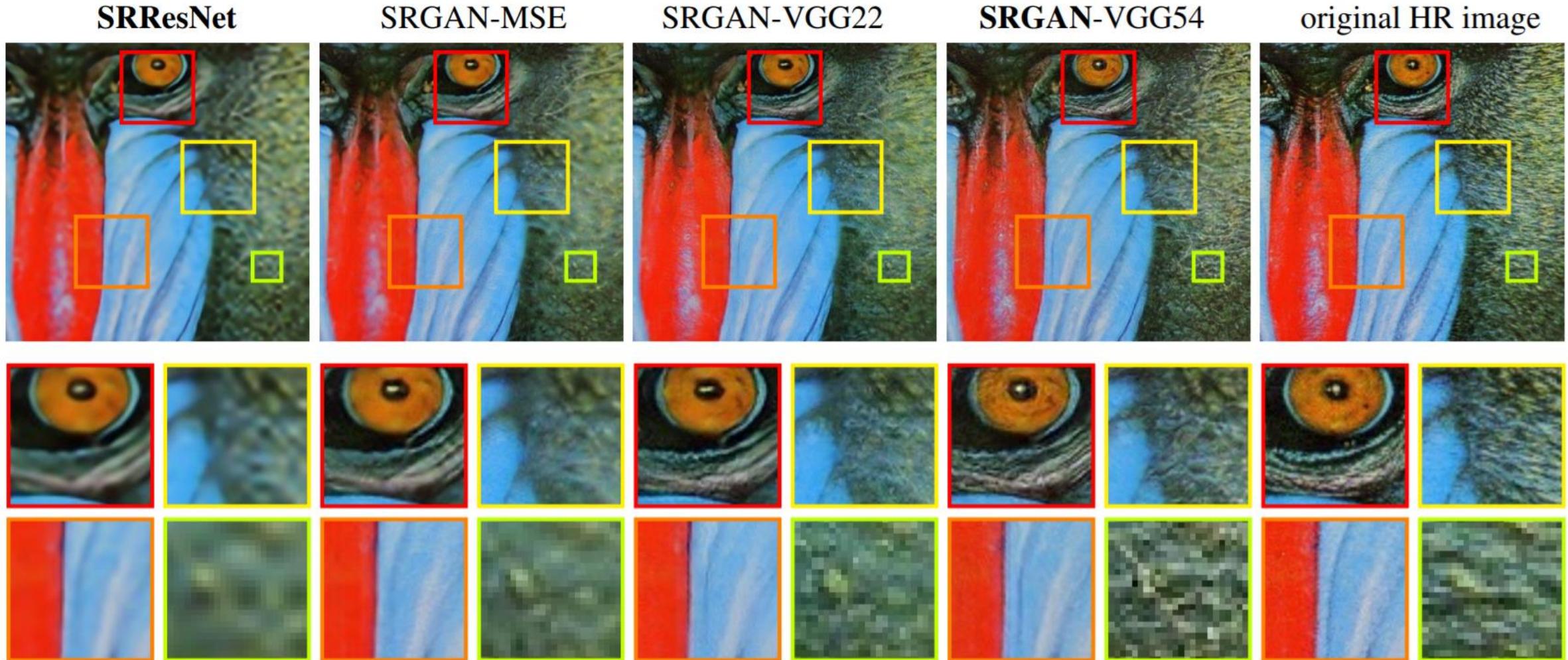
original



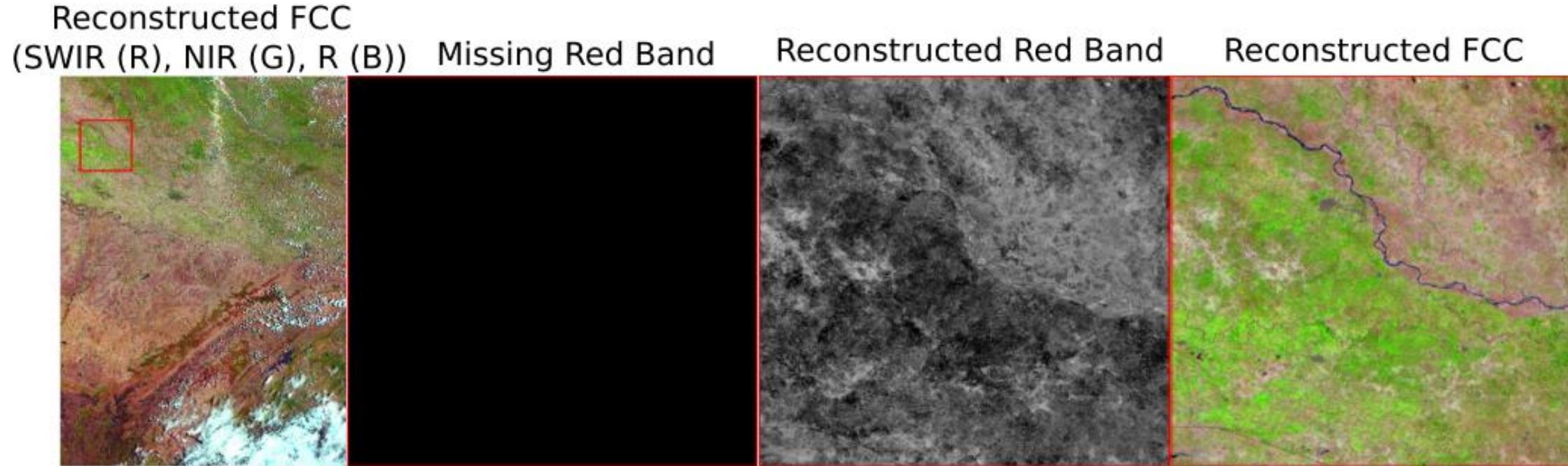
# Mitigating Vanishing Gradient



# Mitigating Vanishing Gradient



# Mitigating Vanishing Gradient



Metrics	Parameters	Training epochs	Training RMSE	Inference Time
AeroGAN	2 M	100	41.57	64 s
DeepAWiFS	4.7 M	1000	18.2	495 s
DSen2	4.7 M	1000	21.2	456 s
ALERT	2.8 M	50	10.3	70 s
% gain	40.42	95.00	51.41	84.64

# Understanding the Role of Adversarial Regularization in Supervised Learning

- Mitigating Vanishing Gradient
- Asymptotic Iteration Complexity
- Sub-optimality Gap
- Provable Convergence
- Generalization Error (Open Problems)
- Neural Topology Analysis (Open Problems)

# Asymptotic Iteration Complexity

## Supervised Objective

**Assumption 5**  $((L_0, L_1)\text{-Smoothness})$ . *The function is  $(L_0, L_1)$ -smooth, i.e., there exist positive constants  $L_0$  and  $L_1$  such that  $\|\nabla^2 l(\theta)\| \leq L_0 + L_1 \|\nabla l(\theta)\|$ .*

**Theorem 2.** *Suppose the functions in  $\mathcal{L}$  satisfy Assumption 3, 4 and 5. Given  $\epsilon > 0$ , the iteration complexity in sole supervision is upper bounded by*

$$\mathcal{O}\left(\frac{(l(\theta_0) - l^*)(L_0 + L_1 L^2 \beta \epsilon)}{\epsilon^2}\right).$$

# Asymptotic Iteration Complexity

## Augmented Objective

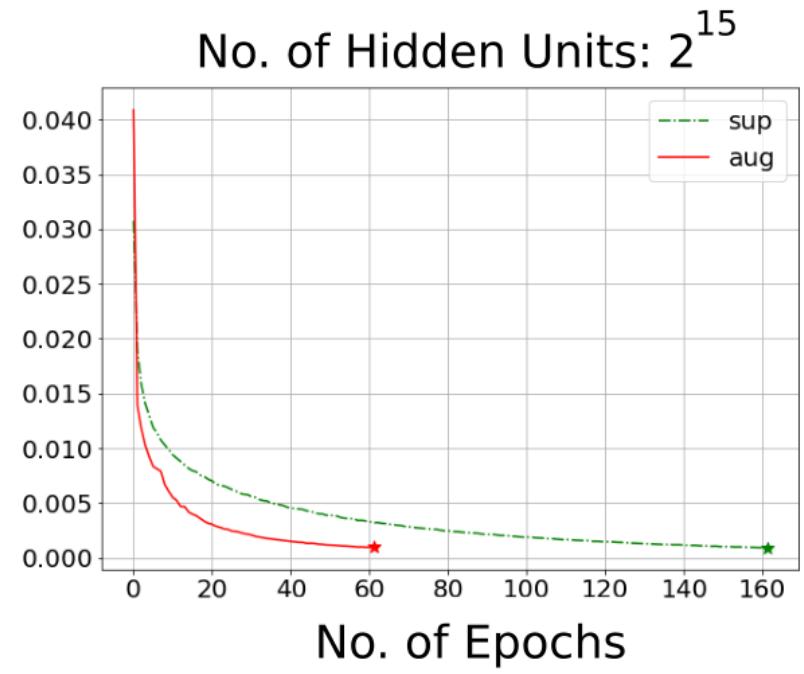
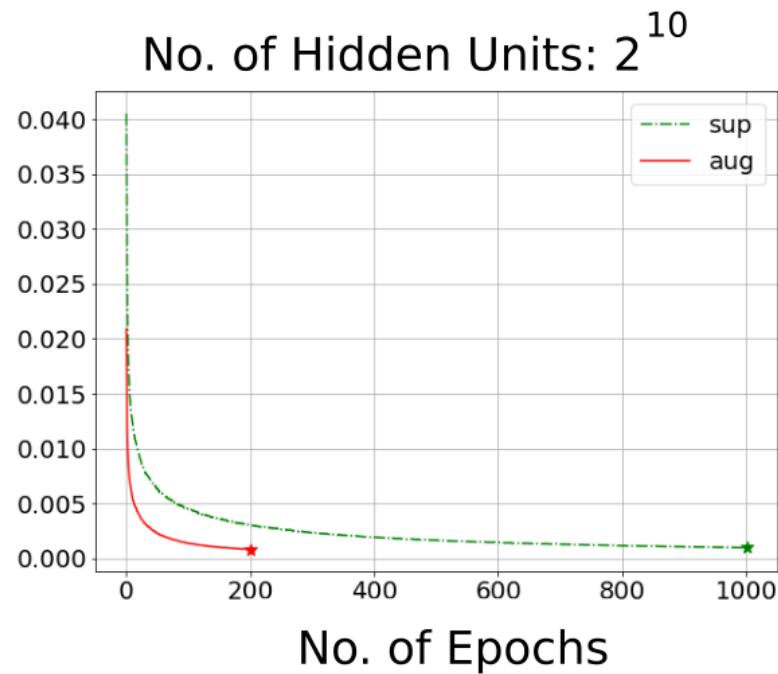
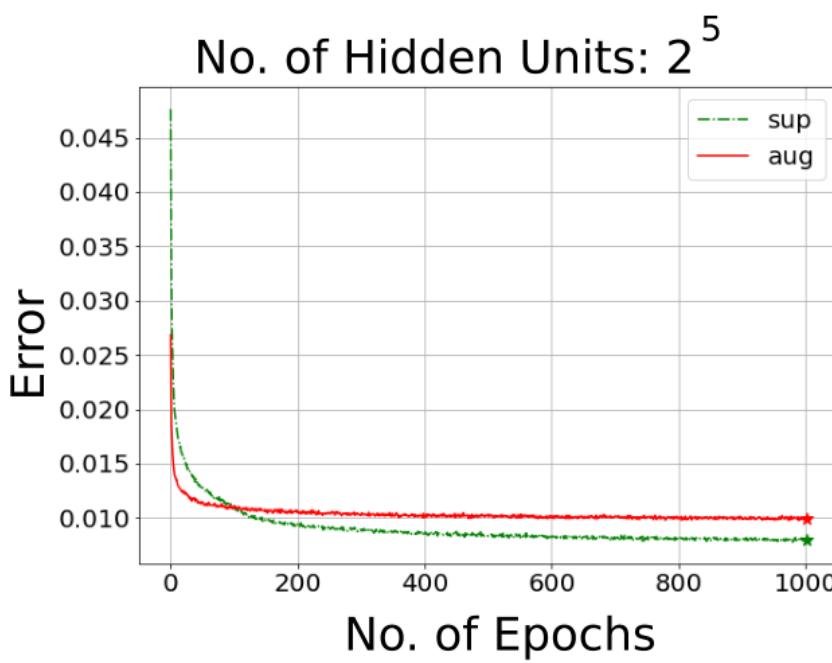
### Theorem 2 and 3 (Informal).

*Under valid assumptions, the number of iterations required to find a good solution in the augmented objective is equal to or less than supervised objective.*

**Theorem 3.** Suppose the functions in  $\mathcal{L}$  satisfy **Assumption 3, 4 and 5**. Given **Assumption 6** holds,  $\epsilon > 0$  and  $\delta \leq \frac{\sqrt{2\epsilon\zeta}}{L}$ , the iteration complexity in adversarial regularization is upper bounded by  $\mathcal{O} \left( \frac{(l(\theta_0) - l^*) (L_0 + L_1 L^2 \beta \epsilon)}{\epsilon^2 + 2\epsilon\zeta - L^2 \delta^2} \right)$ .

# Asymptotic Iteration Complexity

## Experimental Result



# Asymptotic Iteration Complexity

## Experimental Result: Generator Network Configurations

Architecture	No. Layer	Activation	No. ResBlock	No. DenseBlock	No. Epoch Sup	No. Epoch Aug	Hypothesis
MLP-Deep	6	ELU	2	0	391	55	✓
CNN-ResNet	6	ReLU	2	0	215	41	✓
CNN-DenseNet	6	ReLU	2	1	163	39	✓
CNN-DenseNet-L1	6	ReLU	2	1	1000	39	✓
CNN-DenseNet-L2	6	ReLU	2	1	155	39	✓
CNN-ResNet-AvgPool	6	ReLU	2	0	109	29	✓

Adversarial regularization is at least as good as sole supervision.

# Understanding the Role of Adversarial Regularization in Supervised Learning

- Mitigating Vanishing Gradient
- Asymptotic Iteration Complexity
- Sub-optimality Gap
- Provable Convergence
- Generalization Error (Open Problems)
- Neural Topology Analysis (Open Problems)

# Sub-optimality Gap

## Supervised Objective

**Theorem 4.** *In purely supervised learning, the sub-optimality gap at the average over all iterates in a trajectory of  $T$  time steps is upper bounded by  $\mathcal{O}\left(\frac{\|\theta(0)-\theta^*\|^2}{2T}\right)$ .*

Generator:  $\kappa(t) = \kappa(\theta(t)) := l(\theta(t)) - l(\theta^*)$

# Sub-optimality Gap

## Augmented Objective

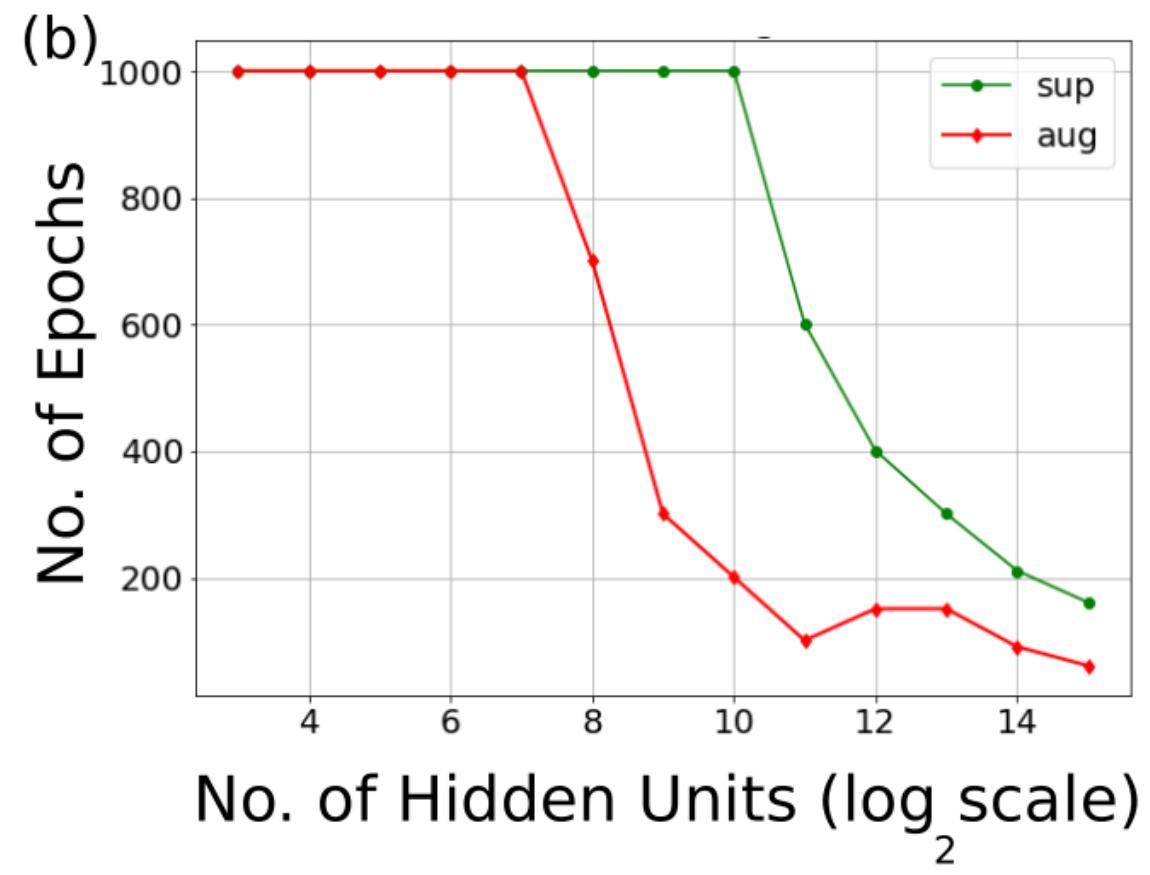
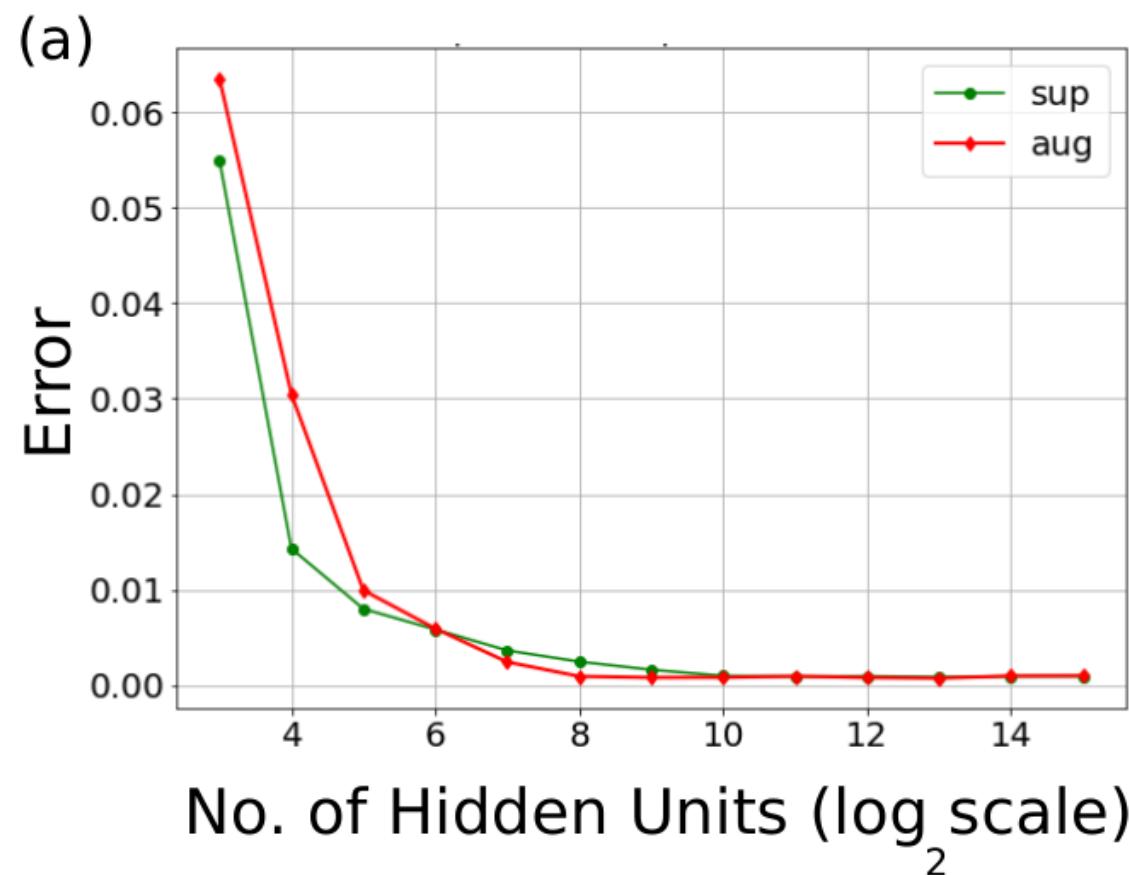
**Theorem 5.** *In supervised learning with adversarial regularization, the sub-optimality gap at the average over all iterates in a trajectory of  $T$  time steps is upper bounded by*

$$\mathcal{O} \left( \frac{\|\theta(0) - \theta^*\|^2}{2T} - \pi \left( \frac{1}{T} \int_0^T \theta(t) dt \right) \right).$$

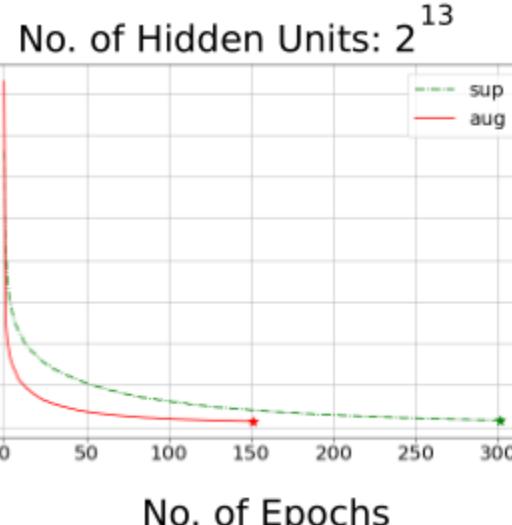
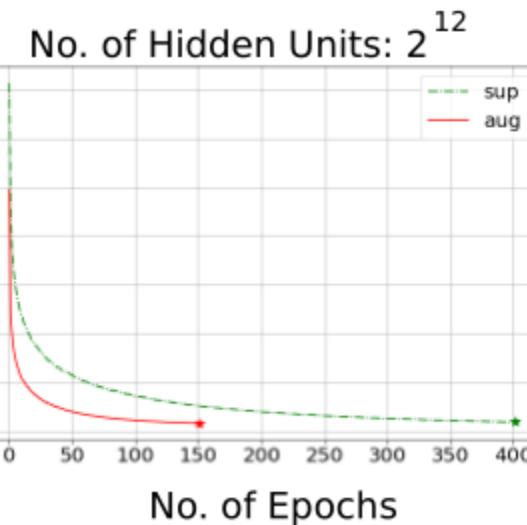
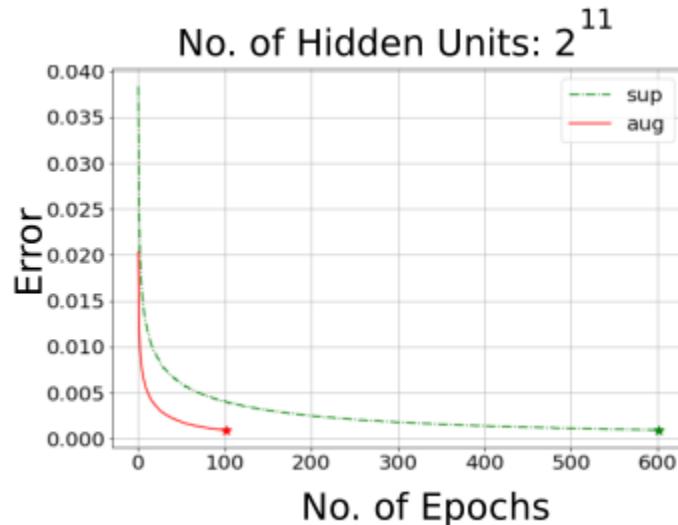
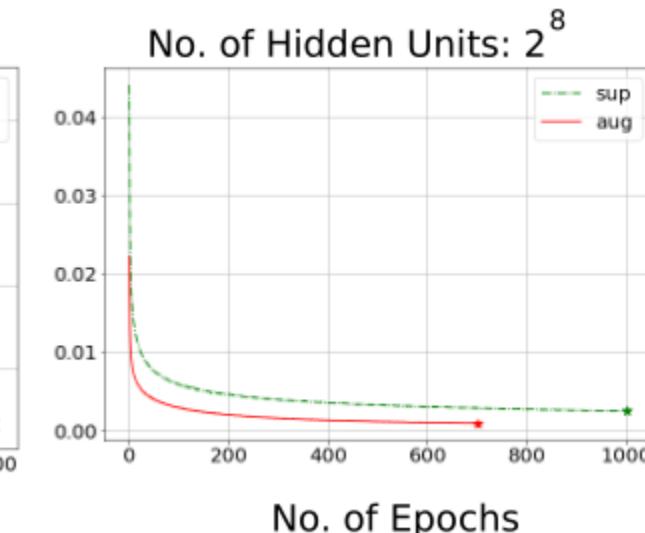
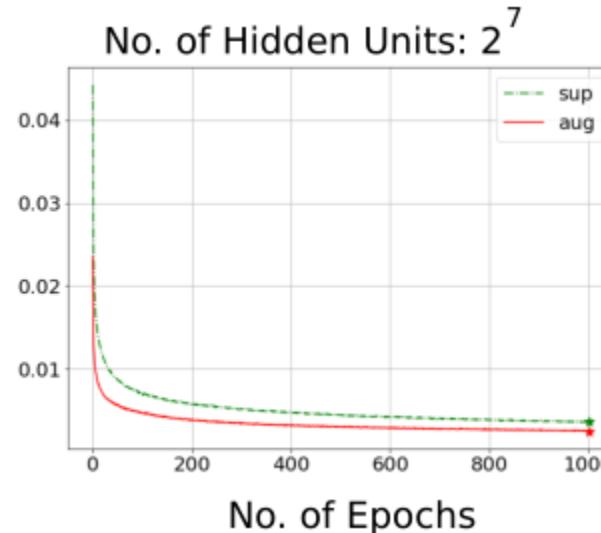
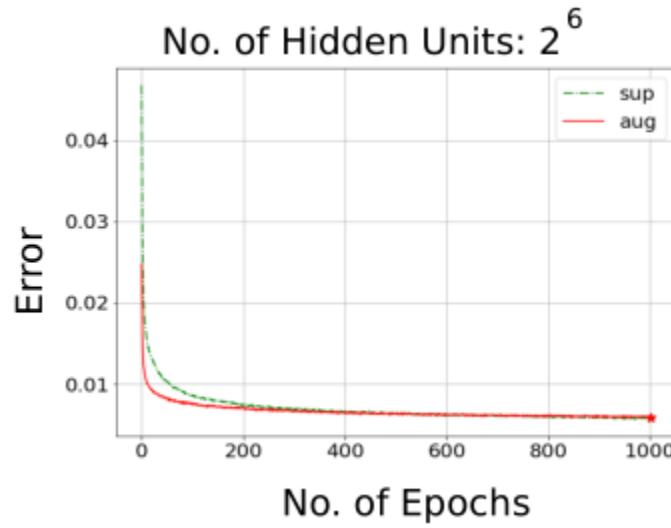
Discriminator:  $\pi(t) = \pi(\theta(t)) := g(\theta^*) - g(\theta(t))$

# Sub-optimality Gap

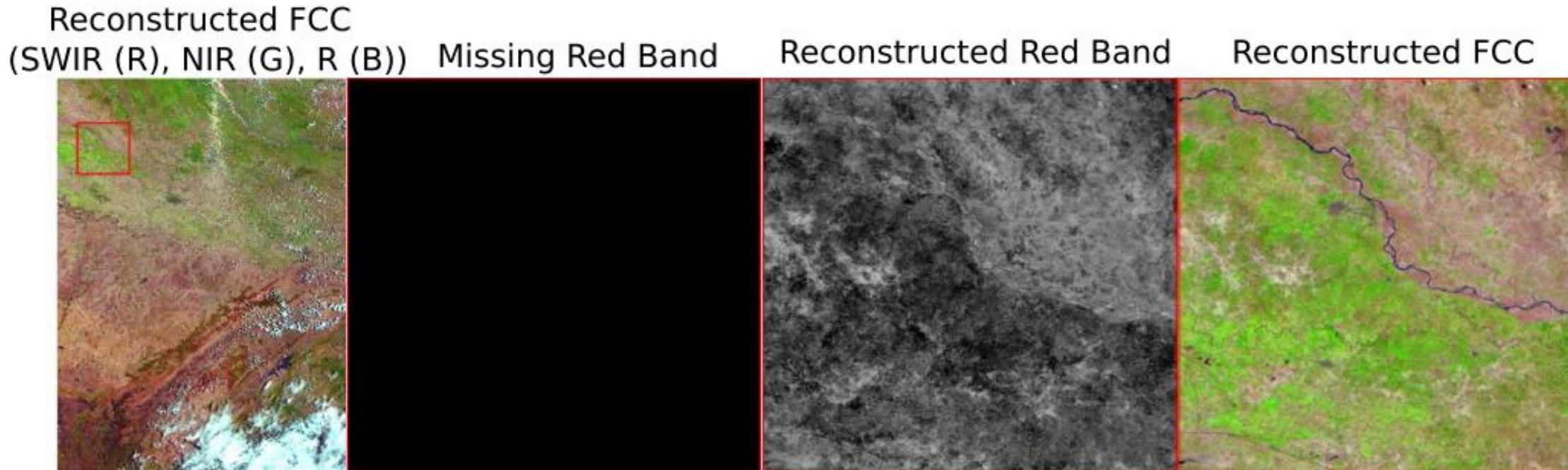
## Experimental Result



# Sub-optimality Gap



# Sub-optimality Gap



Metrics	Parameters	Training epochs	Training RMSE	Inference Time
AeroGAN	2 M	100	41.57	64 s
DeepAWiFS	4.7 M	1000	18.2	495 s
DSen2	4.7 M	1000	21.2	456 s
ALERT	2.8 M	50	10.3	70 s
% gain	40.42	95.00	51.41	84.64

# Understanding the Role of Adversarial Regularization in Supervised Learning

- Mitigating Vanishing Gradient
- Asymptotic Iteration Complexity
- Sub-optimality Gap
- Provable Convergence
- Generalization Error (Open Problems)
- Neural Topology Analysis (Open Problems)

# Provable Convergence

## Augmented Objective

**Theorem 6** (Strongly-convex-strongly-concave convergence). Suppose **Assumption 7** holds. Let  $\ell(\theta_k) \triangleq l(\theta_k) - g(\theta_k)$  is a  $\mu$ -strongly convex function. Let  $\{\theta_k\}$  be the sequence of iterates obtained using global clipping on SGD with zero momentum. Define the output to be  $k$ -weighted combination of iterates:  $\bar{\theta} = \frac{\sum_{k=1}^T k \theta_{k-1}}{\sum_{k=1}^T k}$ . If adaptive clipping  $\tau_k = G k^{\frac{1}{\alpha}} \mu^{\frac{1}{\alpha}}$  and step size  $\eta_k = \frac{5}{2\mu(k+1)}$ , then the output iterate  $\bar{\theta}$  satisfies

$$\mathbb{E} [l(\bar{\theta})] - l(\theta^*) \leq \mathcal{O} \left( G^2 (\mu(T+1))^{\frac{2-2\alpha}{\alpha}} - (g(\theta^*) - \mathbb{E}[g(\bar{\theta})]) \right).$$

# Provable Convergence

## Augmented Objective

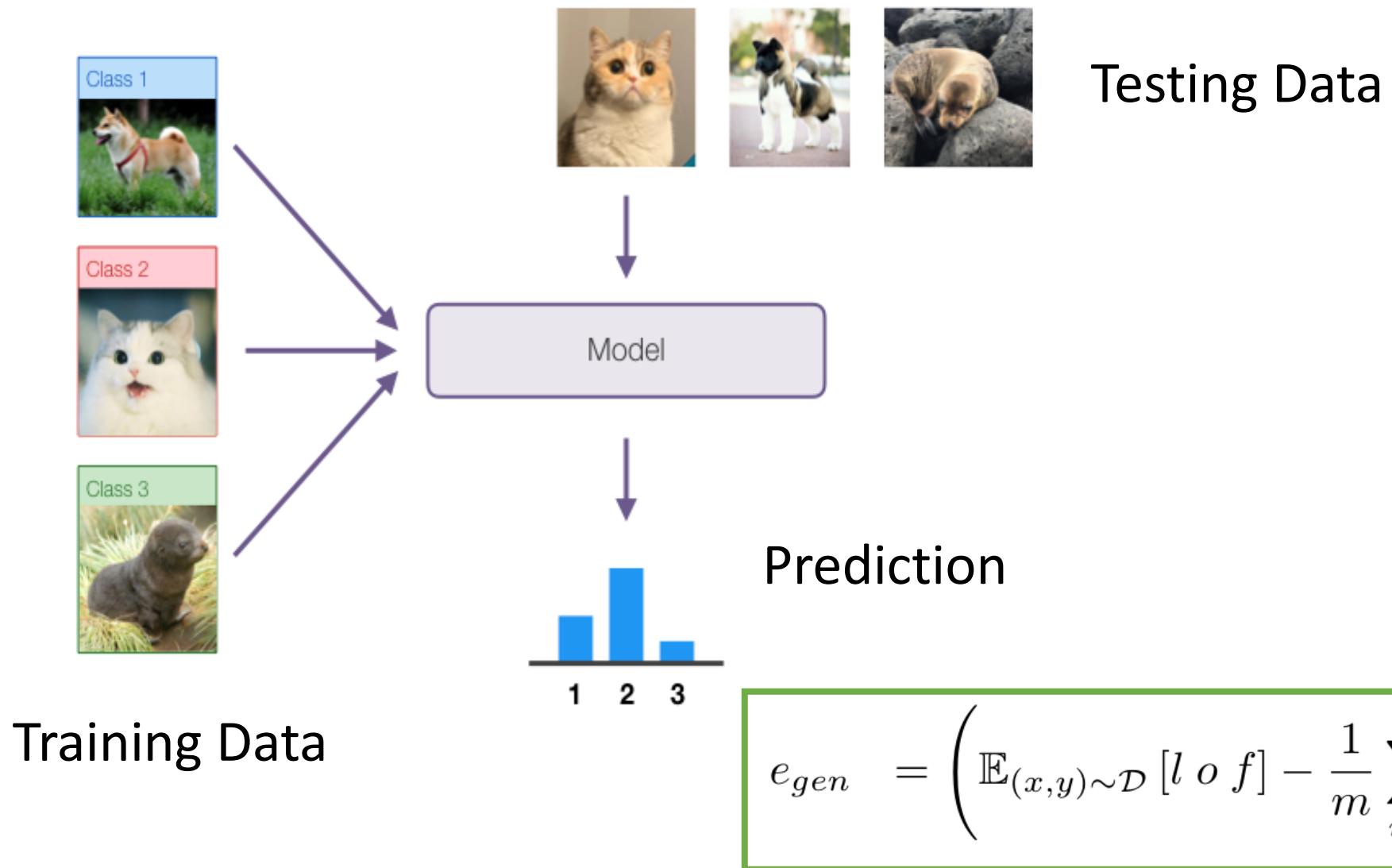
**Theorem 7** (Nonconvex-nonconcave convergence). Suppose **Assumption 7** holds. Let  $\mathfrak{l}(\theta_k) \triangleq l(\theta_k) - g(\theta_k)$  is a possible  $L$ -smooth function and  $\{\theta_k\}$  be the sequence of iterates obtained using global clipping on SGD with zero momentum. Given constant clipping  $\tau_k = G(\eta_k L)^{\frac{-1}{\alpha}}$  and constant step size  $\underline{\eta_k} = \left( \frac{R_0^\alpha L^{2-2\alpha}}{G^2 T^\alpha} \right)^{\frac{1}{3\alpha-2}}$ , where  $R_0 = l(\theta_0) - l(\theta^*)$ , the sequence  $\{\theta_k\}$  satisfies

$$\frac{1}{T} \sum_{k=1}^T \mathbb{E} \left[ \|\nabla l(\theta_{k-1})\|^2 \right] \leq \mathcal{O} \left( G^{\frac{2\alpha}{3\alpha-2}} \left( \frac{R_0 L}{T} \right)^{\frac{2\alpha-2}{3\alpha-2}} - \frac{1}{T} \sum_{k=1}^T \mathbb{E} \left[ \|\nabla g(\theta_{k-1})\|^2 \right] \right).$$

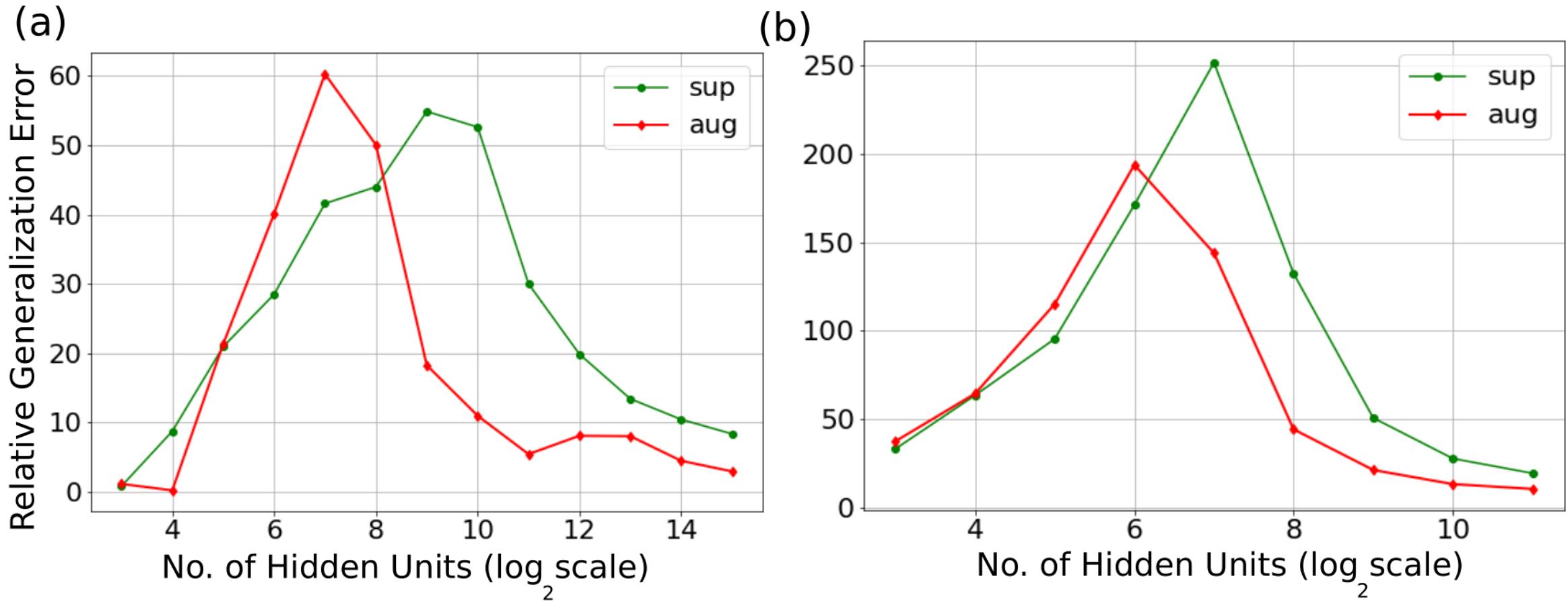
# Understanding the Role of Adversarial Regularization in Supervised Learning

- Mitigating Vanishing Gradient
- Asymptotic Iteration Complexity
- Sub-optimality Gap
- Provable Convergence
- Generalization Error (Open Problems)
- Neural Topology Analysis (Open Problems)

# Generalization Error

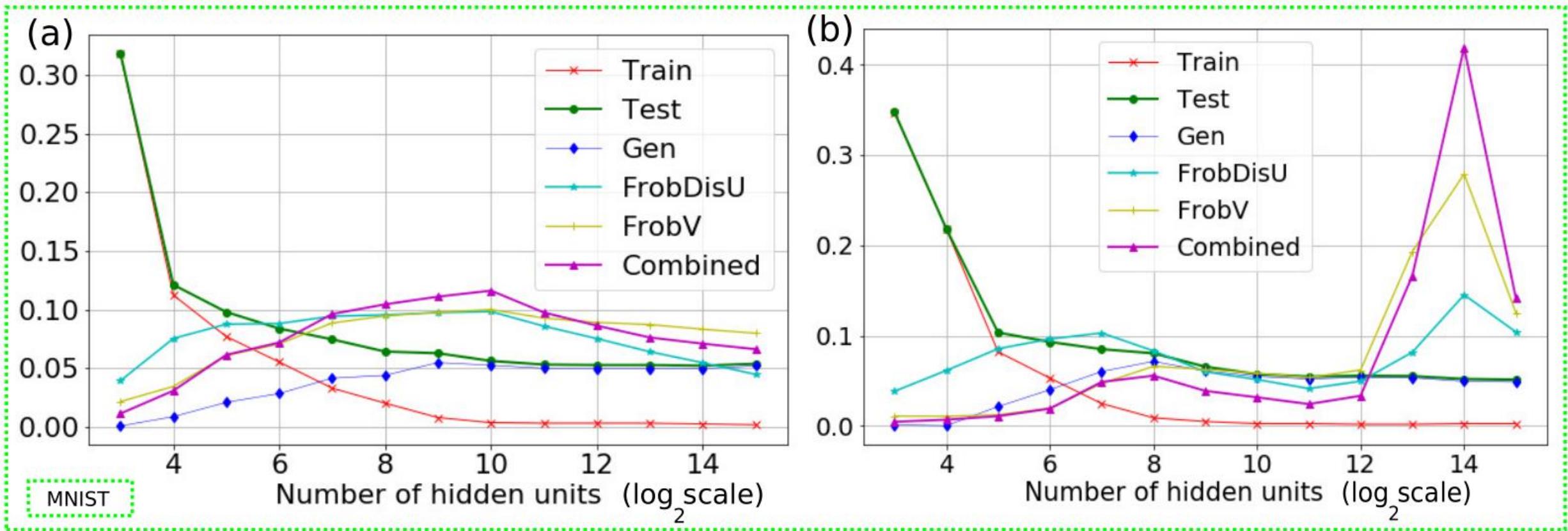


# Generalization Error



$$e_{gen,r} = \left( \mathbb{E}_{(x,y) \sim \mathcal{D}} [l o f] - \frac{1}{m} \sum_{i=1}^m l(f(x); y) \right) \times \mathcal{N}$$

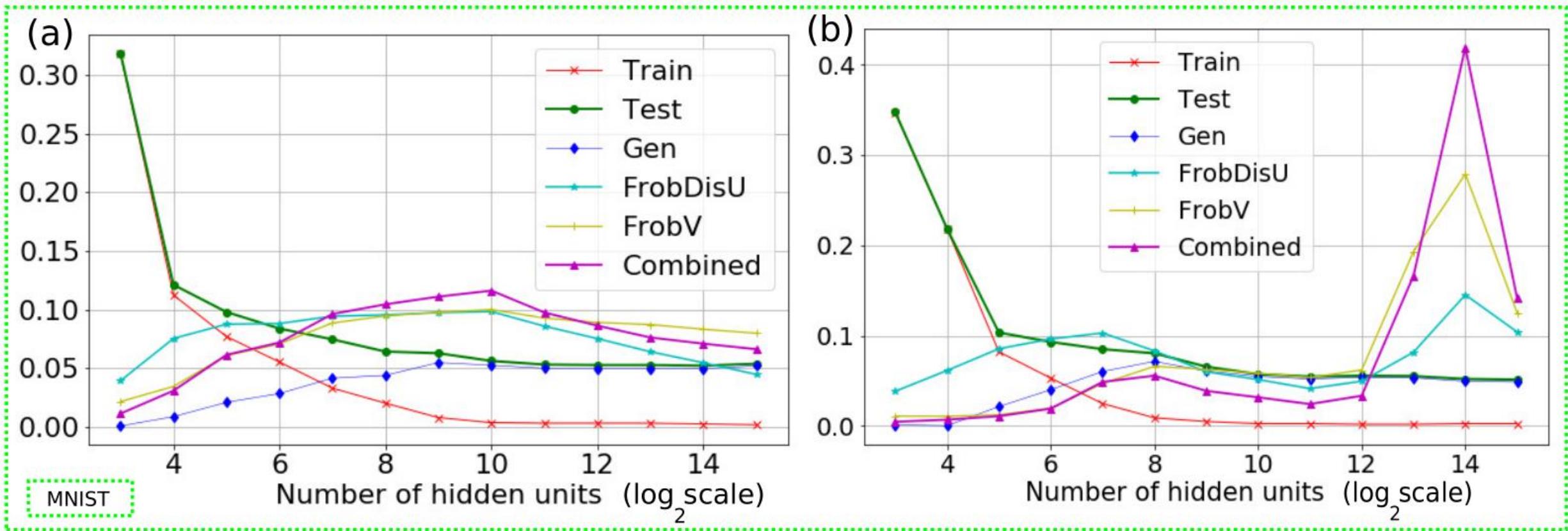
# Generalization Error



Combined Metric:  $\mathcal{O} \left( \|U^0\|_2 \|V\|_F + \|U - U^0\|_F \|V\|_F + \sqrt{h} \right)$

Neyshabur et al. ICLR ,2019

# Generalization Error



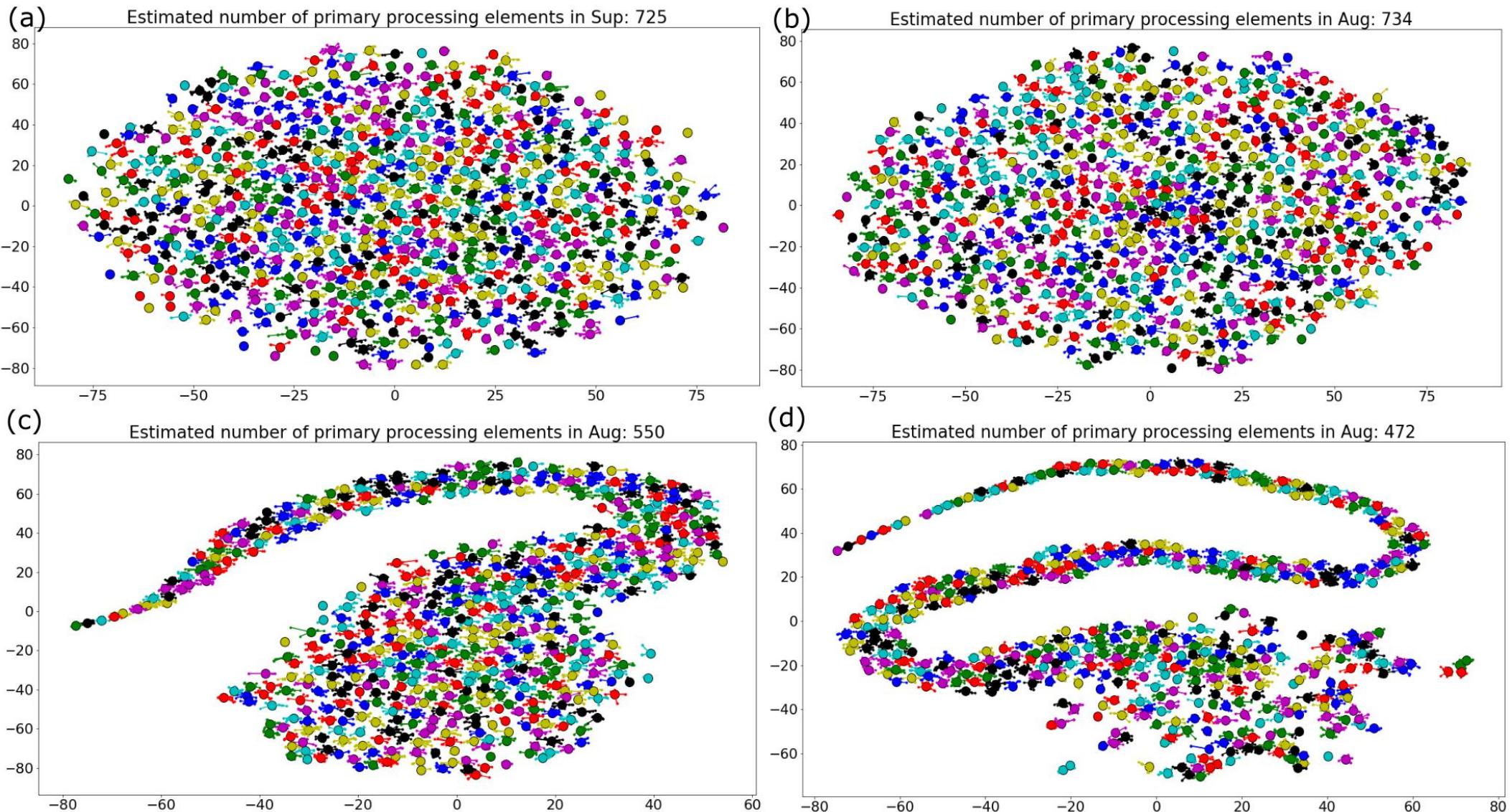
Open Problem:

How to explain the generalization behaviour of the augmented objective?

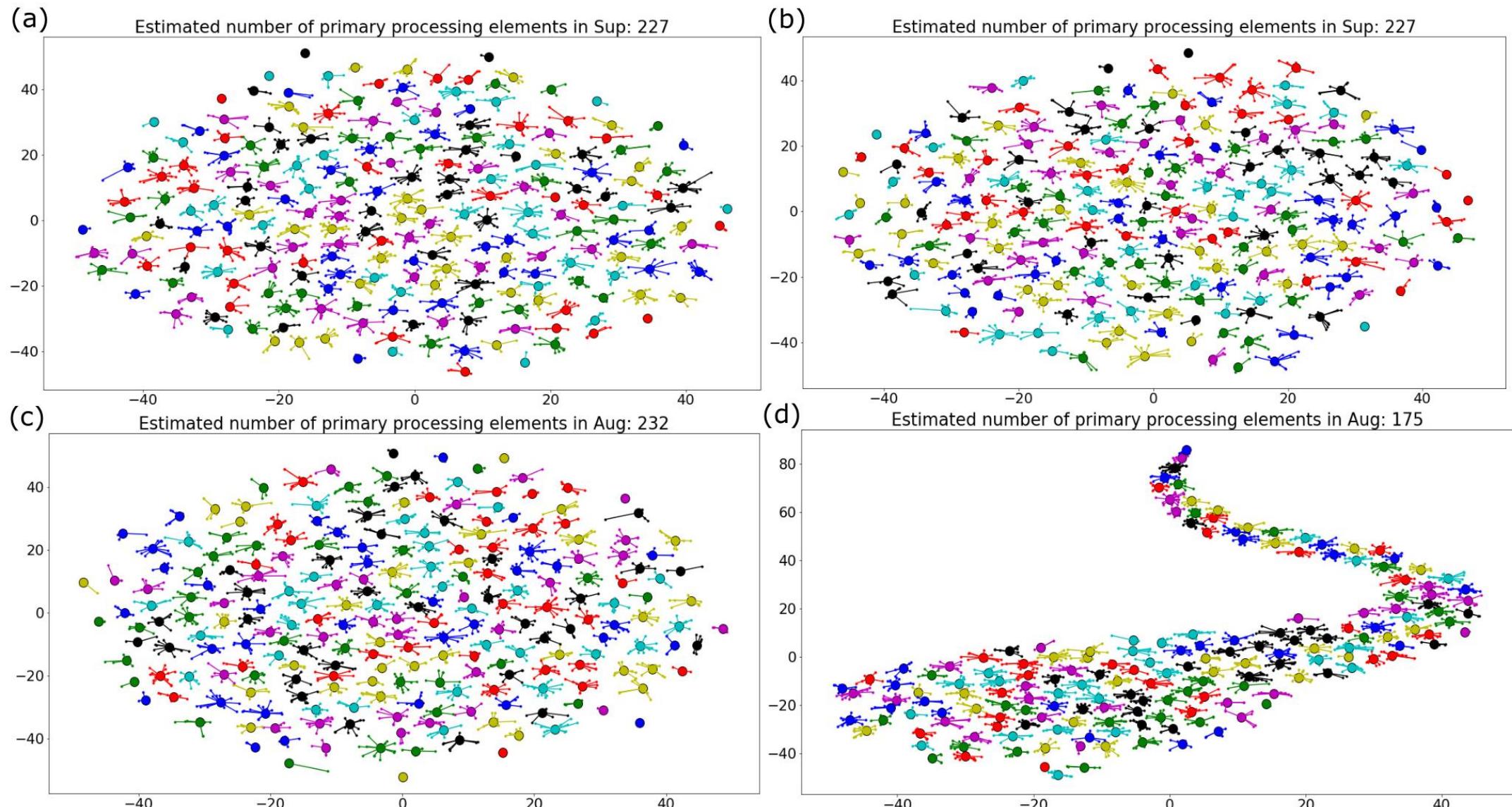
# Understanding the Role of Adversarial Regularization in Supervised Learning

- Mitigating Vanishing Gradient
- Asymptotic Iteration Complexity
- Sub-optimality Gap
- Provable Convergence
- Generalization Error (Open Problems)
- Neural Topology Analysis (Open Problems)

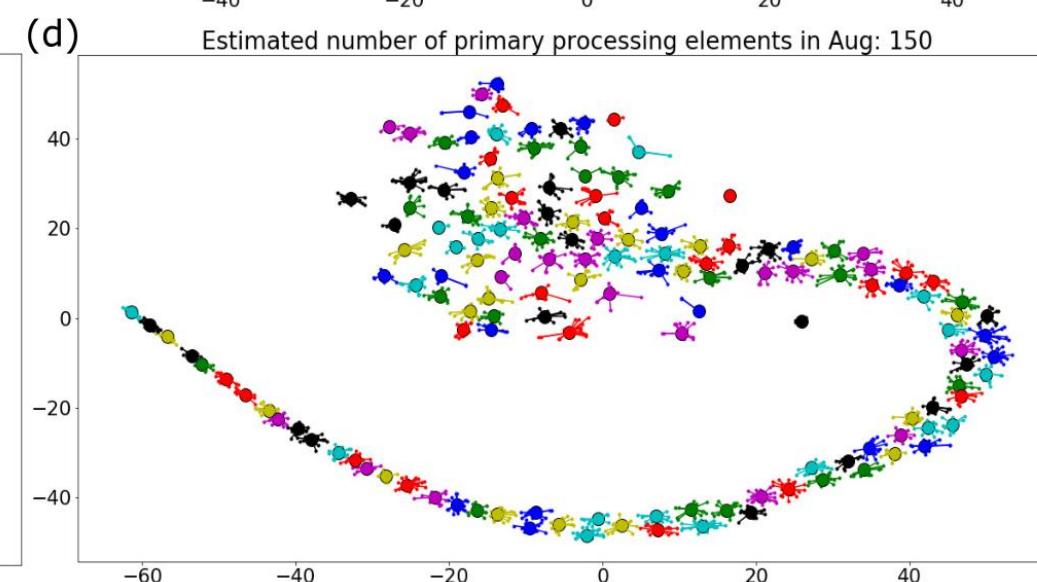
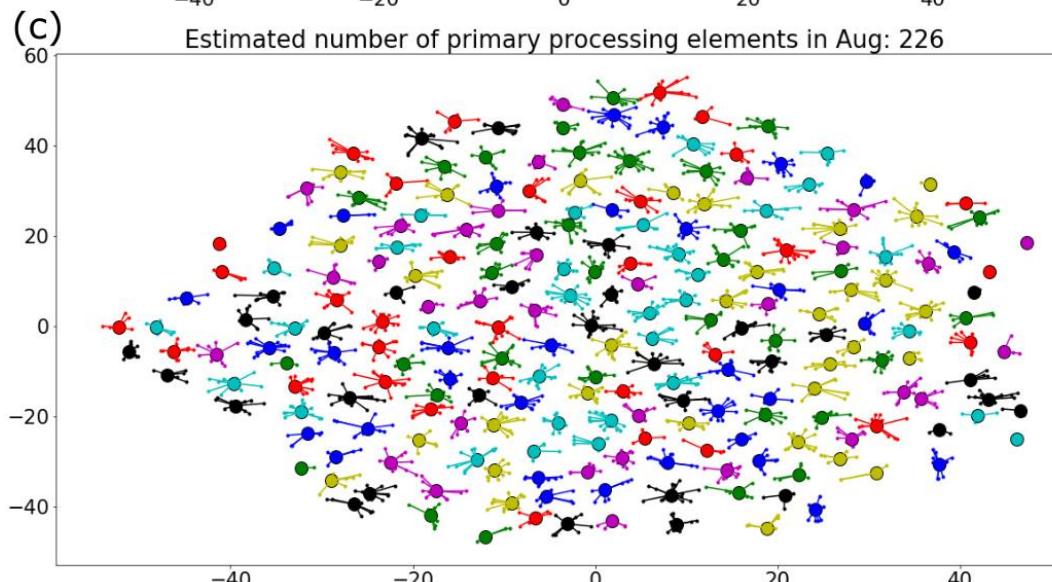
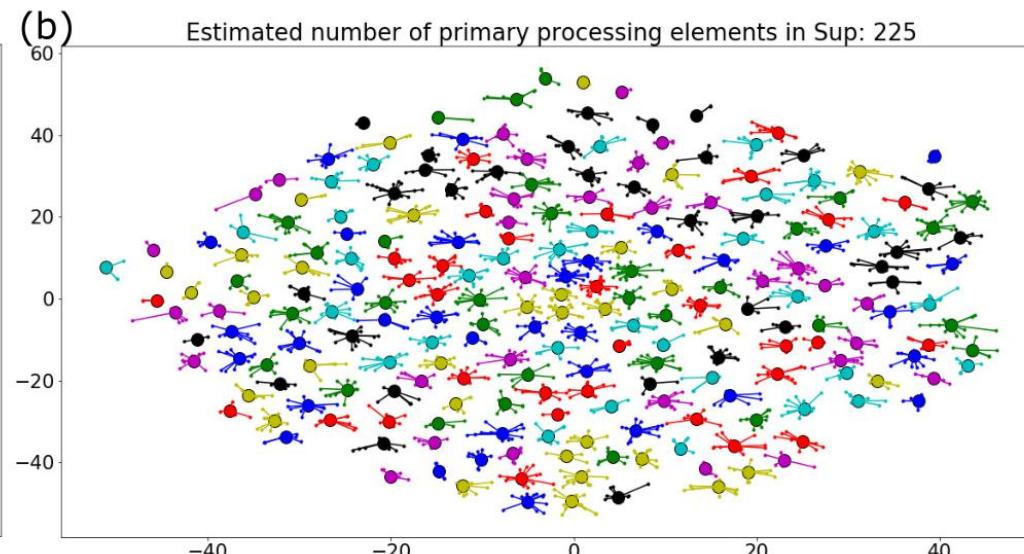
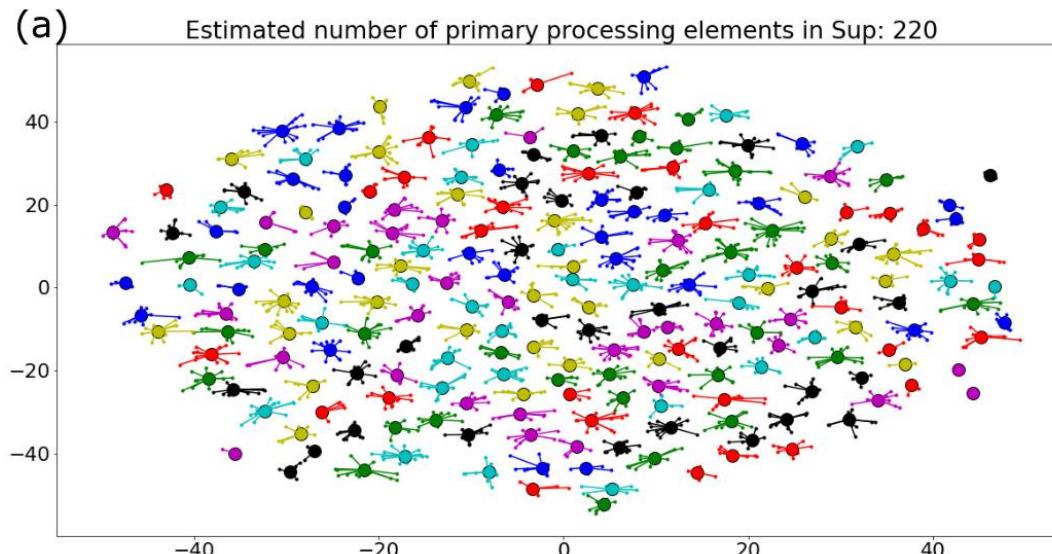
# Neural Topology Analysis



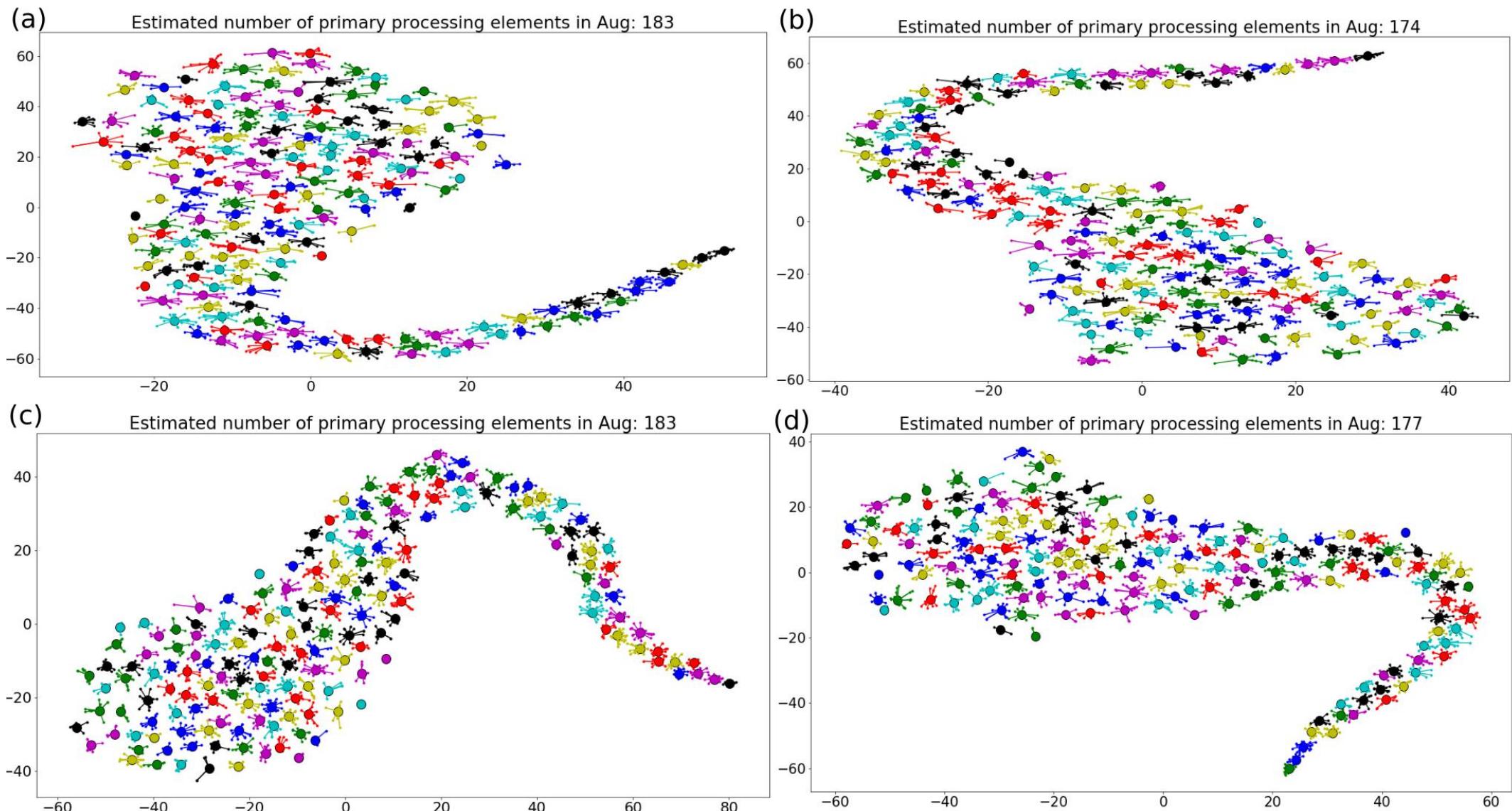
# Neural Topology Analysis: Hidden Layer



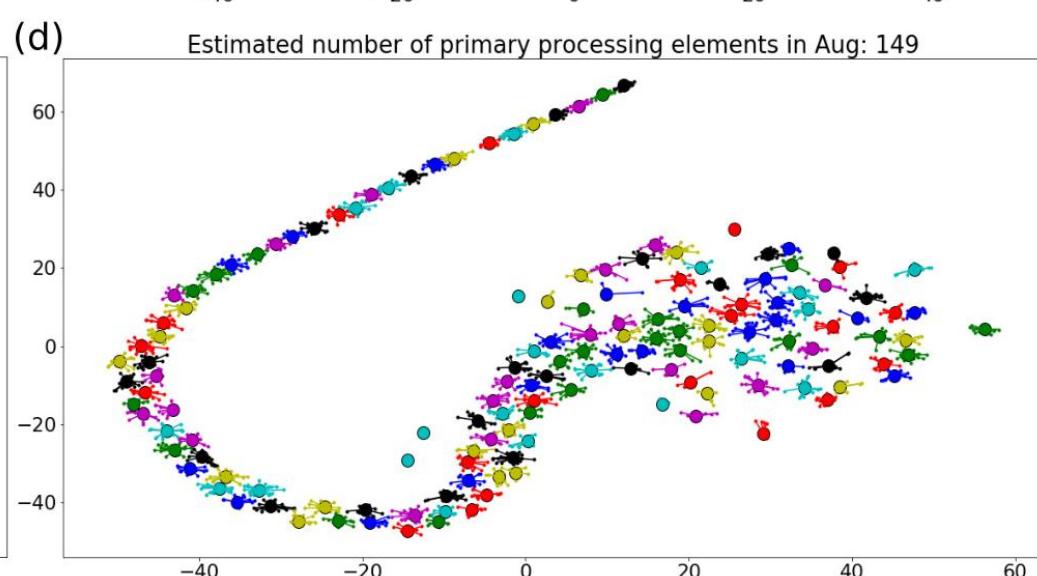
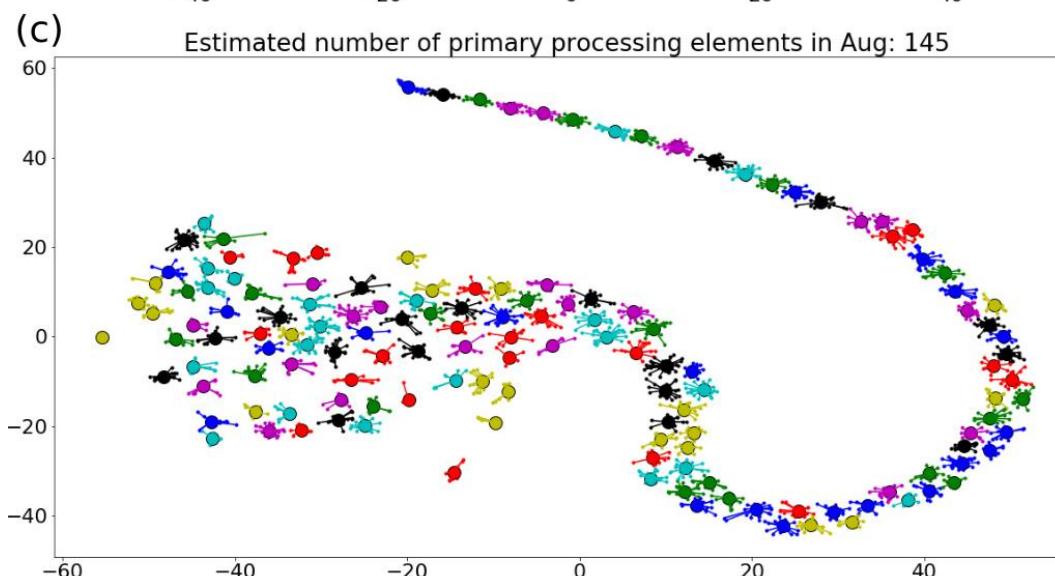
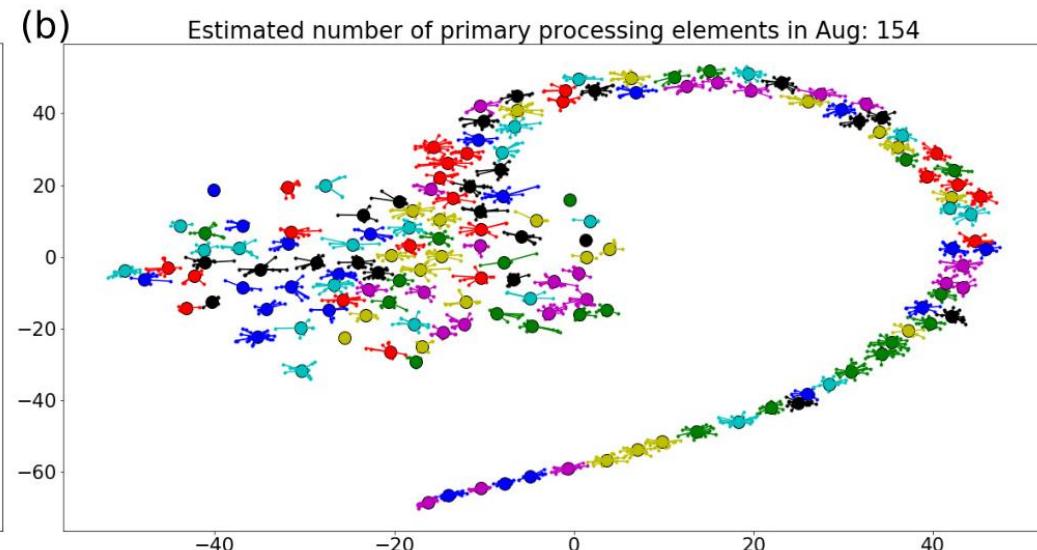
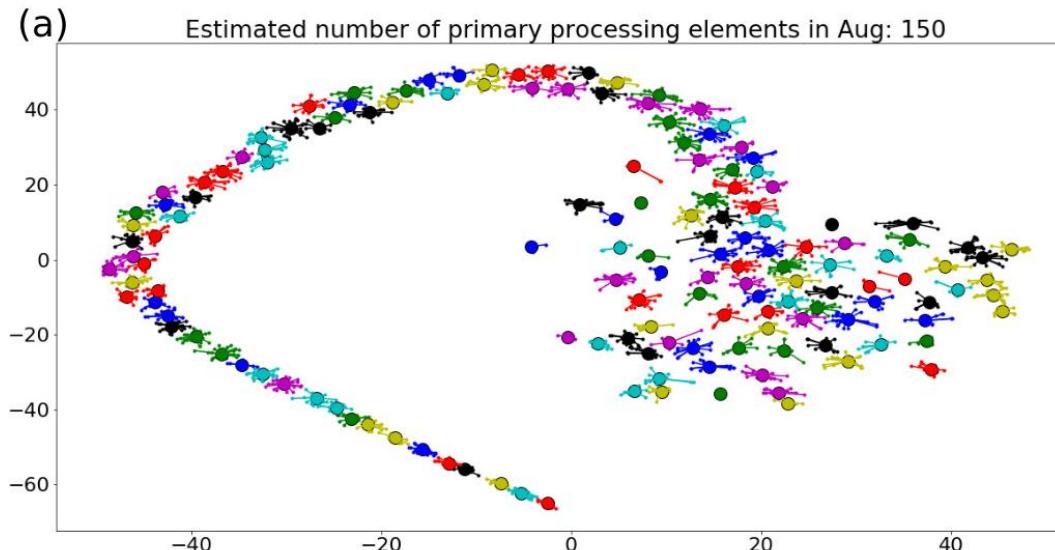
# Neural Topology Analysis: Top Layer



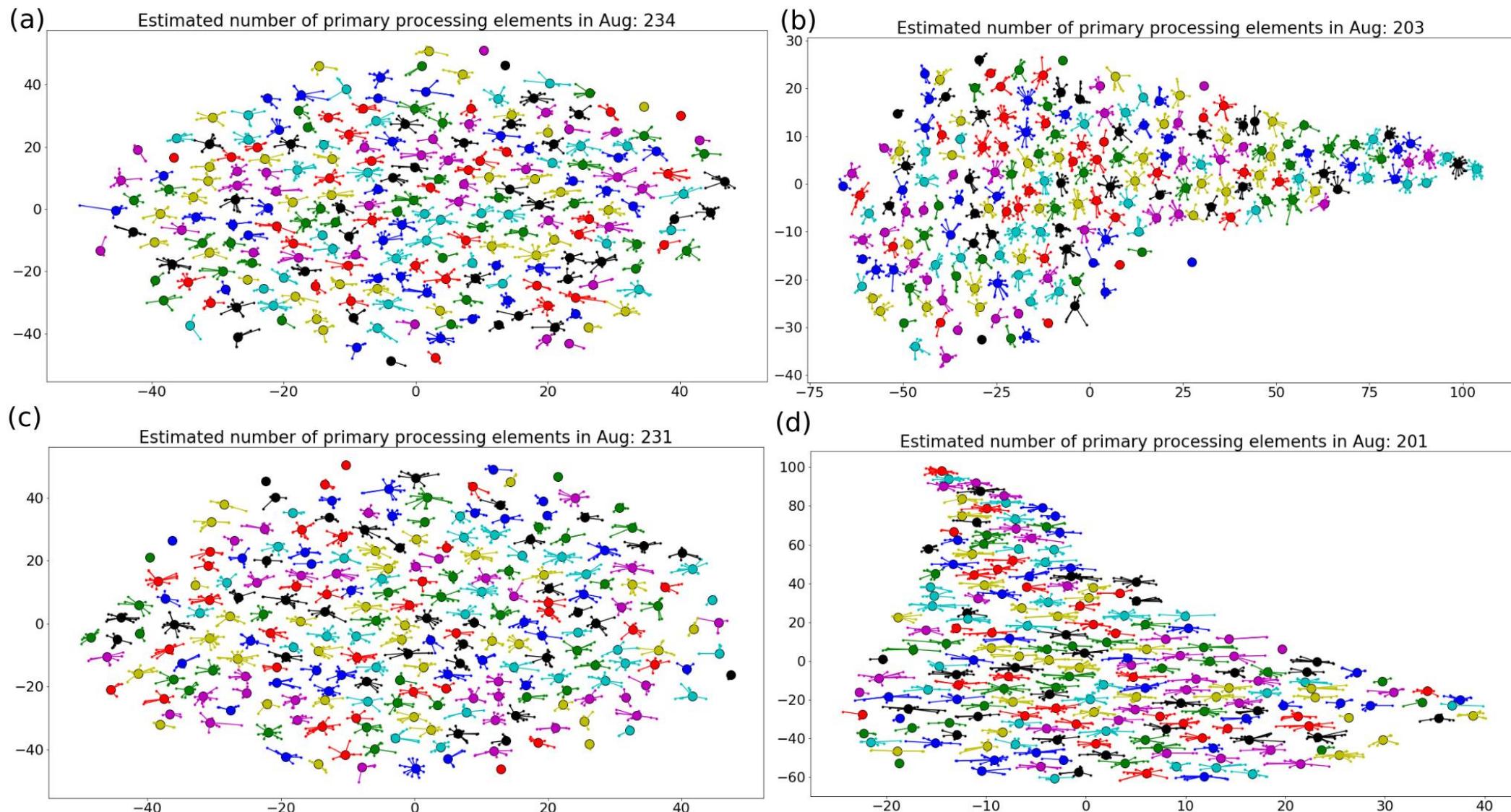
# Neural Topology Analysis: Hidden Layer Subsets



# Neural Topology Analysis: Top Layer Subsets

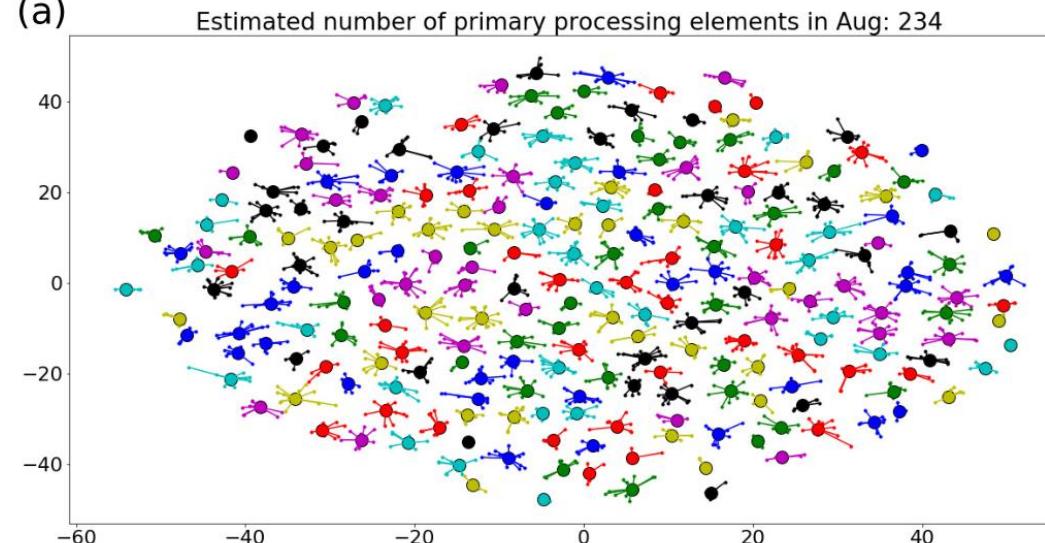


# Neural Topology Analysis: Hidden Layer FMNIST

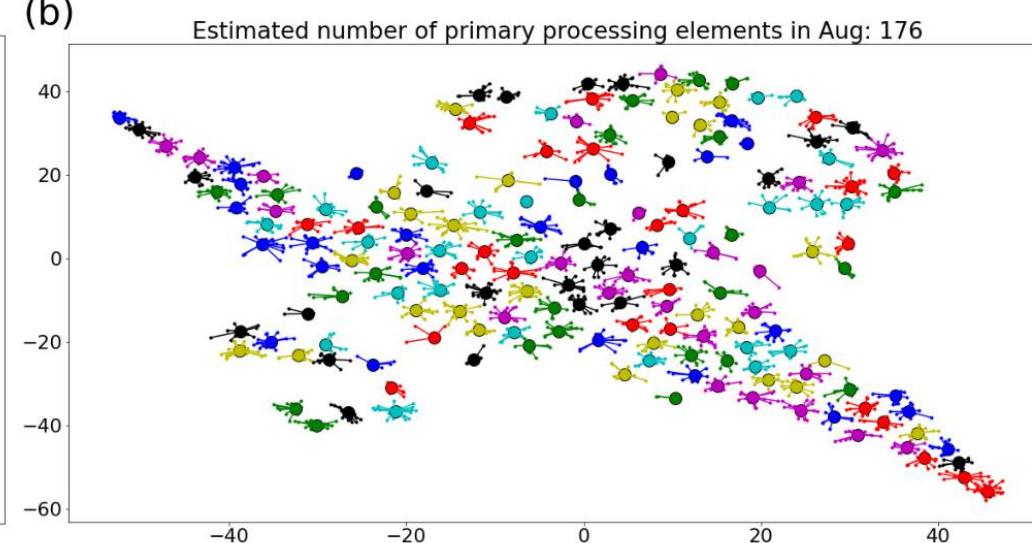


# Neural Topology Analysis: Top Layer FMNIST

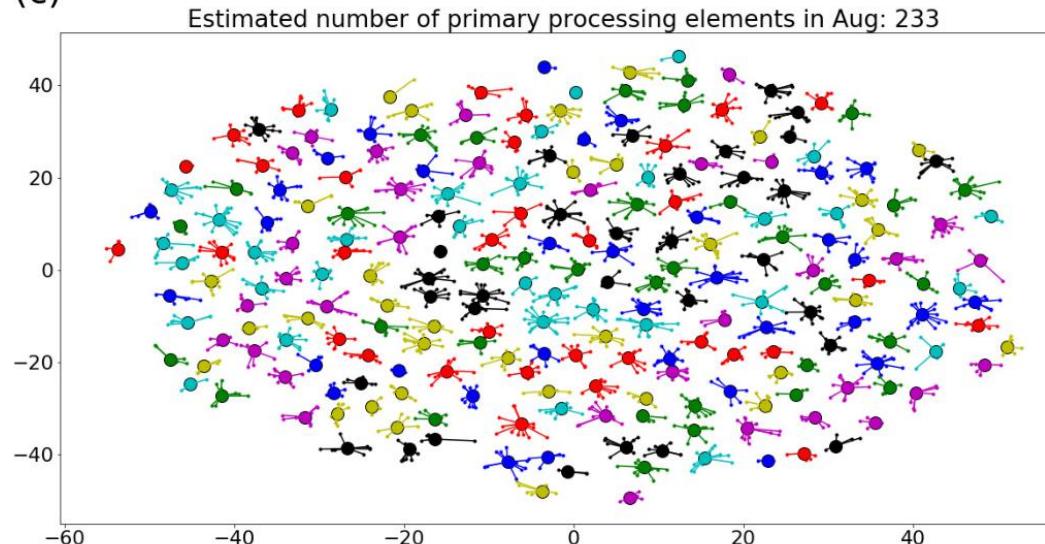
(a)



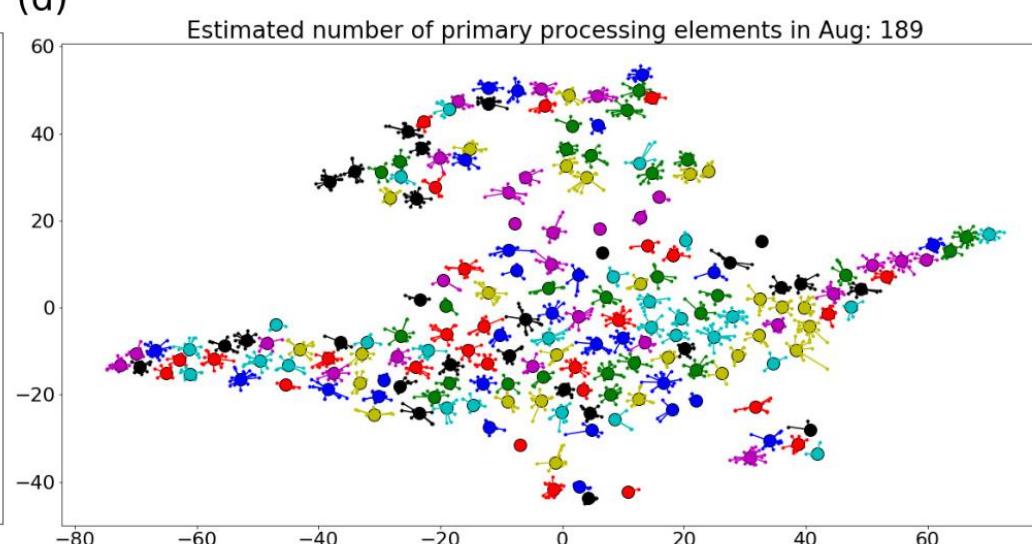
(b)



(c)



(d)



# Interesting Problems

- Why does adversarial interaction create non-homogeneous patterns?
- How do we predict the shape of the non-homogeneous patterns?
- What happens to the parameters if the supervised cost is augmented in the discriminator objective?
- Is gradient descent or backpropagation the optimal learning algorithm?

# Interesting Problems

- Why does adversarial interaction create non-homogeneous patterns?
- How do we predict the shape of the non-homogeneous patterns?
- What happens to the parameters if the supervised cost is augmented in the discriminator objective?
- Is gradient descent or backpropagation the optimal learning algorithm?

# Why Adversarial Interaction Creates Non-Homogeneous Patterns: A Pseudo-Reaction Diffusion Model for Turing Instability

- Adversarial Interaction
  - Generative Adversarial Networks (GANs)
  - Application of conditional GANs
- Non-Homogeneous Patterns
  - Homogeneous patterns
  - Supervised learning
- Reaction-Diffusion
  - Turing's RD model (1952)
  - Gray-Scott RD model (1984)
- Turing Instability
  - Reaction dynamics
  - Diffusion dynamics

# Why Adversarial Interaction Creates Non-Homogeneous Patterns: A Pseudo-Reaction Diffusion Model for Turing Instability

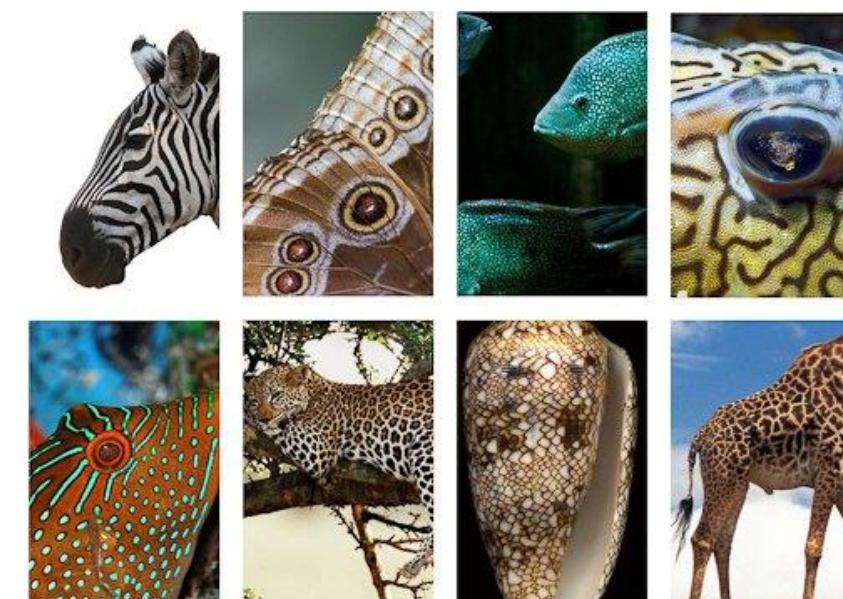
- Adversarial Interaction
  - Generative Adversarial Networks (GANs)
  - Application of conditional GANs
- Non-Homogeneous Patterns
  - Homogeneous patterns
  - Supervised learning
- Reaction-Diffusion
  - Turing's RD model (1952)
  - Gray-Scott RD model (1984)
- Turing Instability
  - Reaction dynamics
  - Diffusion dynamics

# Turing's Reaction-Diffusion Model



$$\frac{d\mathbf{u}_j}{dt} = \mathfrak{R}_j^{\mathbf{u}} (\mathbf{u}_j, \mathbf{v}_j) + \mathfrak{D}_j^{\mathbf{u}} (\nabla^2 \mathbf{u}_j)$$

$$\frac{d\mathbf{v}_j}{dt} = \mathfrak{R}_j^{\mathbf{v}} (\mathbf{u}_j, \mathbf{v}_j) + \mathfrak{D}_j^{\mathbf{v}} (\nabla^2 \mathbf{v}_j)$$

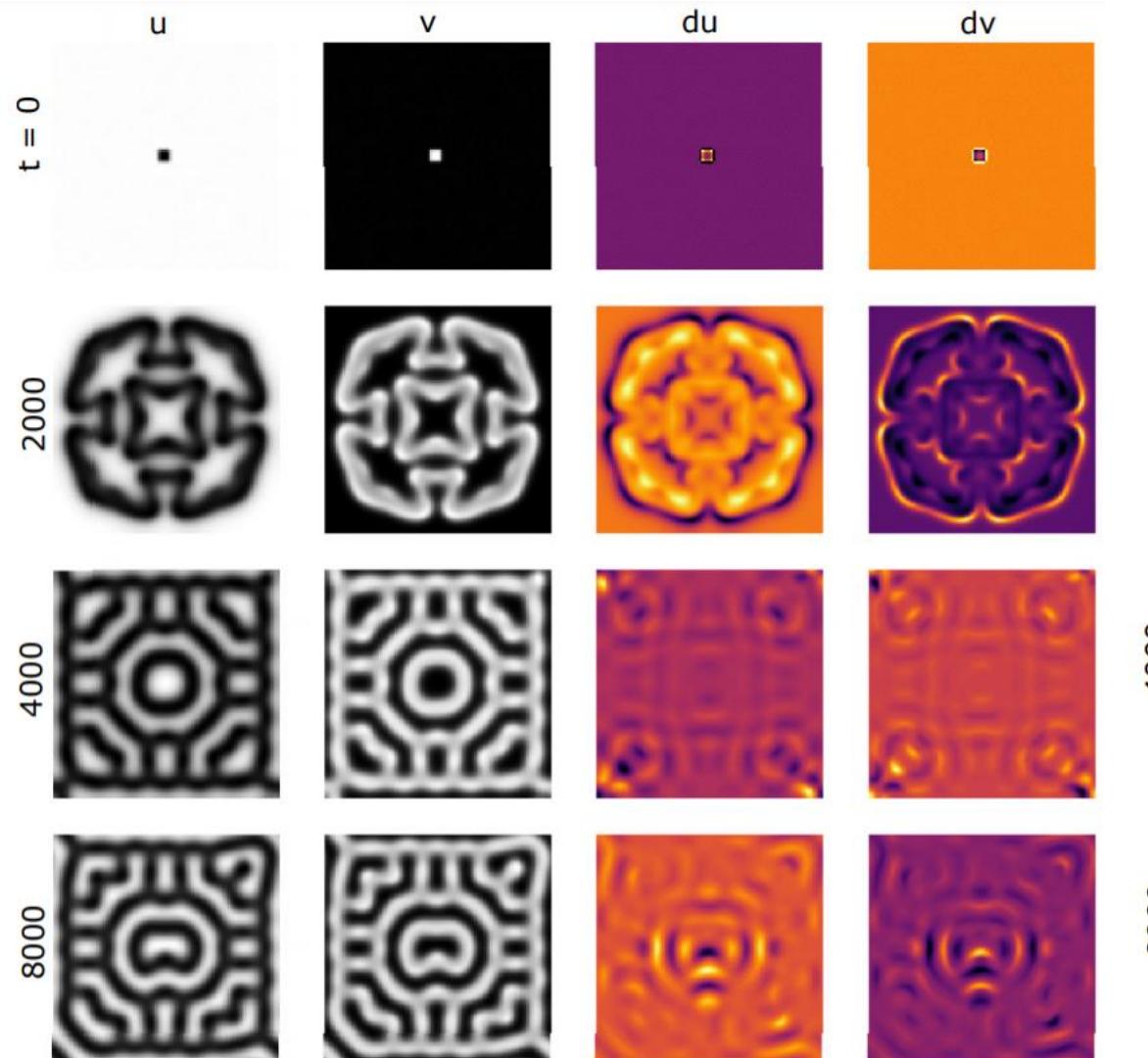


A. M. Turing, 1952

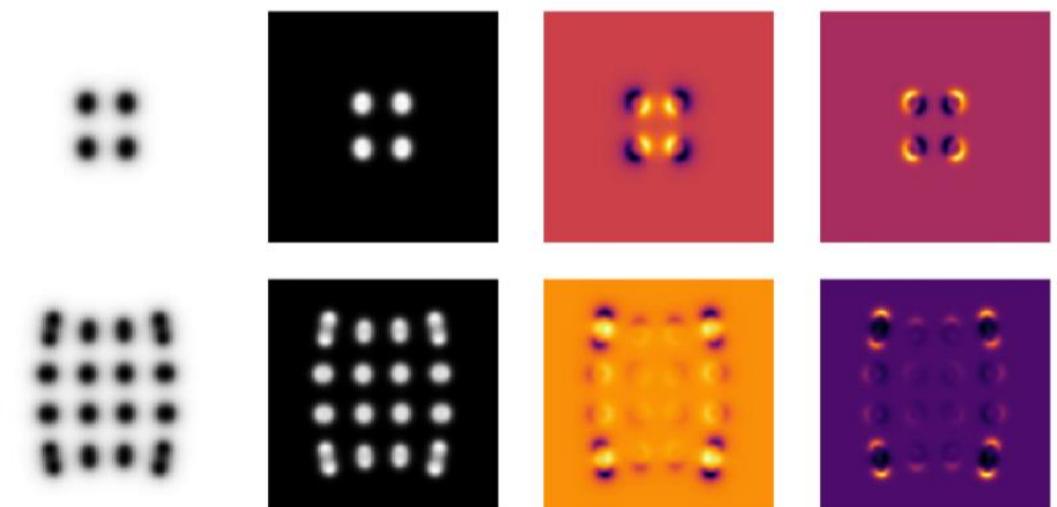
# Why Adversarial Interaction Creates Non-Homogeneous Patterns: A Pseudo-Reaction Diffusion Model for Turing Instability

- Adversarial Interaction
  - Generative Adversarial Networks (GANs)
  - Application of conditional GANs
- Non-Homogeneous Patterns
  - Homogeneous patterns
  - Supervised learning
- Reaction-Diffusion
  - Turing's RD model (1952)
  - Gray-Scott RD model (1984)
- Turing Instability
  - Reaction dynamics
  - Diffusion dynamics

# Gray-Scott Reaction-Diffusion Model

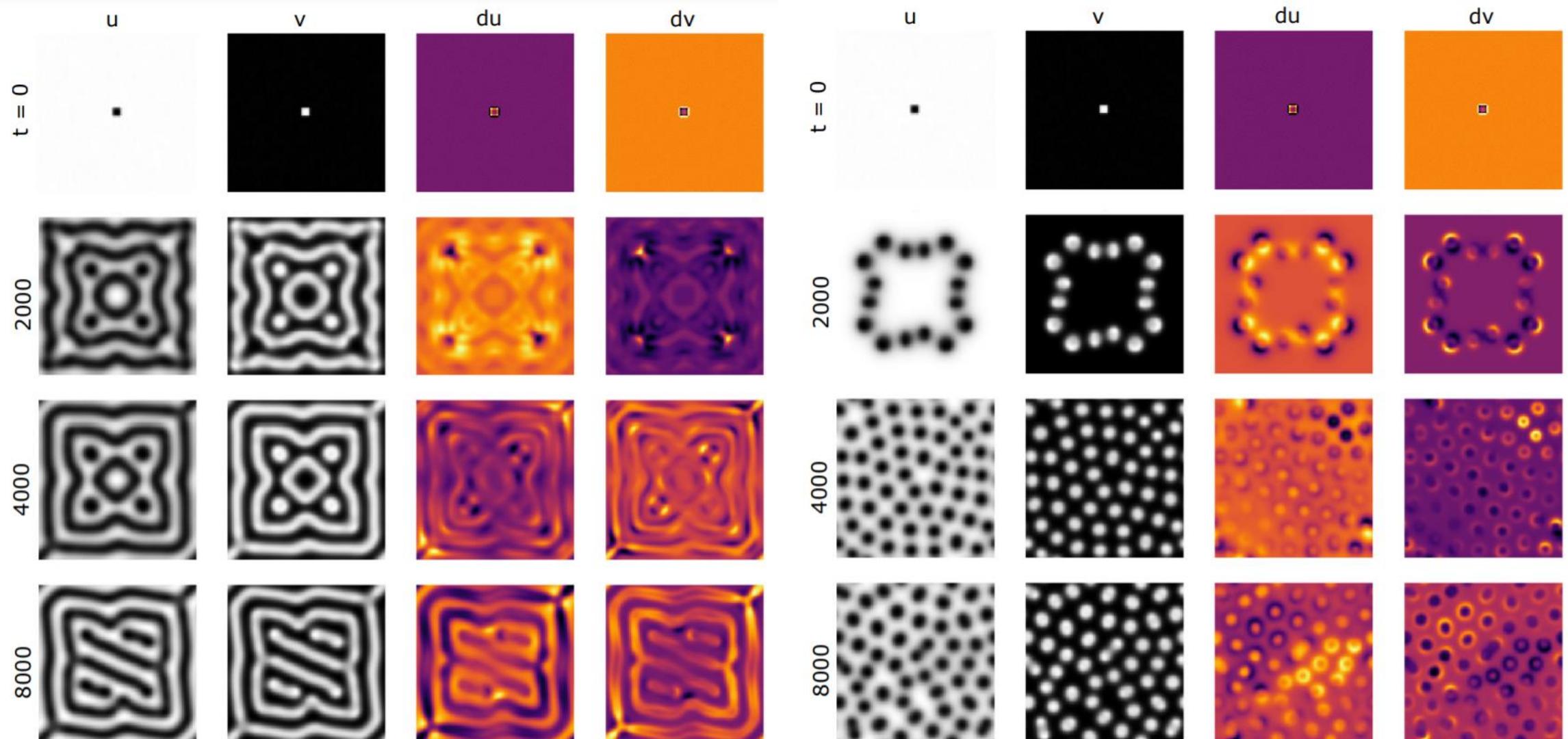


$$\frac{\partial u}{\partial t} = F(1-u) - uv^2 + \mu' \nabla^2 u$$
$$\frac{\partial v}{\partial t} = -(F+k)v + uv^2 + \nu' \nabla^2 v$$



Gray and Scott, 1984

# Gray-Scott Reaction-Diffusion Model



Gray and Scott, 1984

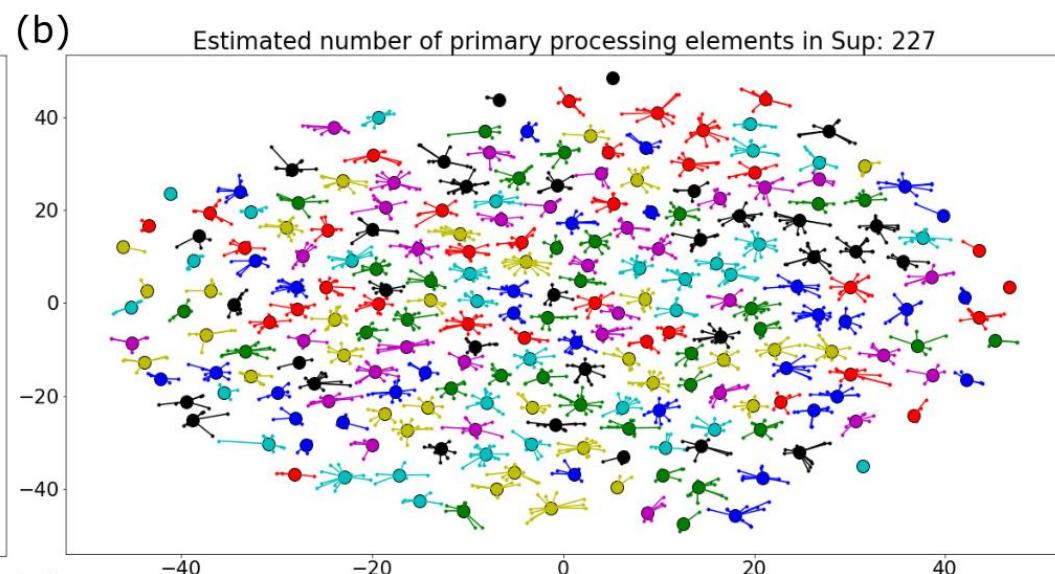
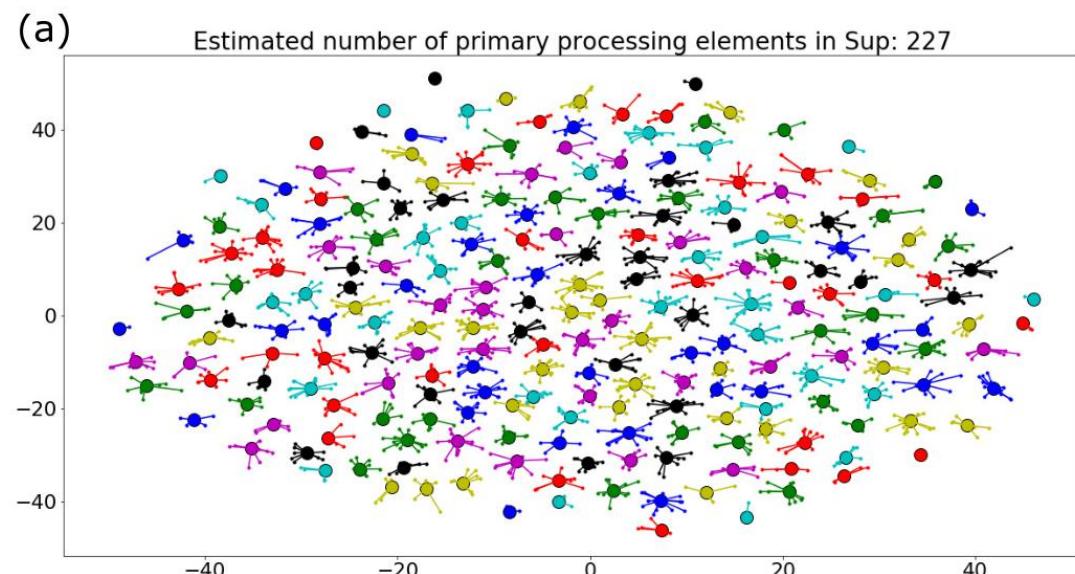
# Why Adversarial Interaction Creates Non-Homogeneous Patterns: A Pseudo-Reaction Diffusion Model for Turing Instability

- Adversarial Interaction
  - Generative Adversarial Networks (GANs)
  - Application of conditional GANs
- Non-Homogeneous Patterns
  - Homogeneous patterns
  - Supervised learning
- Reaction-Diffusion
  - Turing's RD model (1952)
  - Gray-Scott RD model (1984)
- Turing Instability
  - Reaction dynamics
  - Diffusion dynamics

# Symmetry and Homogeneity

## Supervised Learning

$$\mathcal{L}_{sup} (\mathbf{U}, \mathbf{V}) = \frac{1}{2} \sum_{p=1}^n \left\| \frac{1}{\sqrt{d_{out}m}} \mathbf{V} \sigma (\mathbf{U} \mathbf{x}_p) - \mathbf{y}_p \right\|_2^2$$

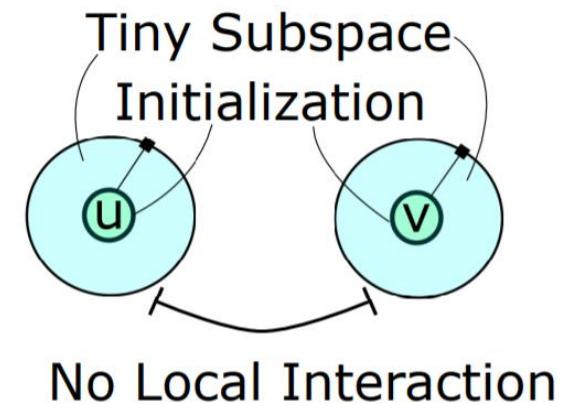


# Symmetry and Homogeneity

## Theorem (Symmetry and Homogeneity)

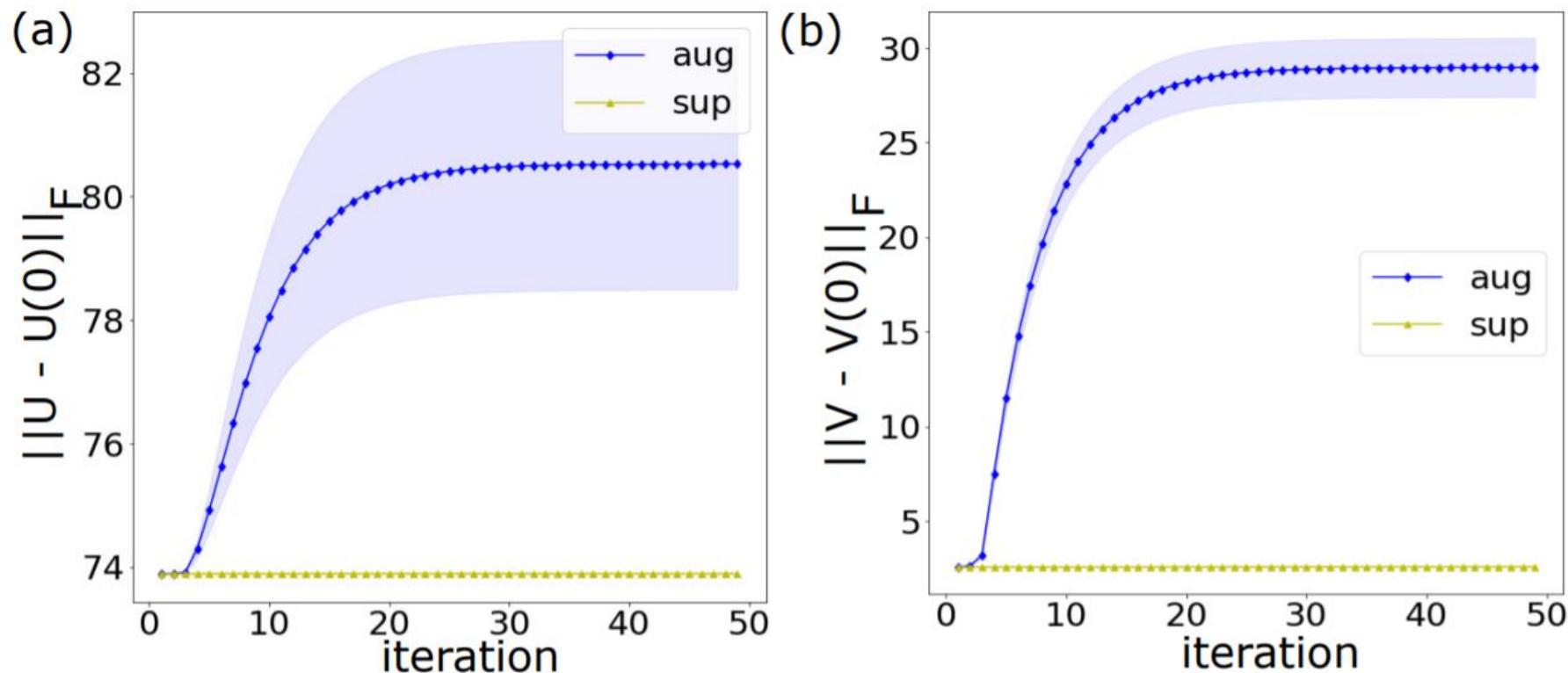
Suppose **Assumption 1** holds. Let us i.i.d. initialize  $u_j \sim \mathcal{N}(0, I)$  and sample  $v_j$  uniformly from  $\{+1, -1\}$  for all  $j \in [m]$ . If we choose  $\|x_p\|_2 = 1$  for  $p \in [n]$ , then we obtain the following with probability at least  $1 - \delta$ :

$$\|\mathbf{u}_j(t) - \mathbf{u}_j(0)\|_2 \leq \mathcal{O}\left(\frac{n^{3/2}}{\sqrt{m}\lambda_0\delta}\right),$$
$$\|\mathbf{U}(t) - \mathbf{U}(0)\|_F \leq \mathcal{O}\left(\frac{n^{3/2}}{\lambda_0\delta}\right).$$



# Symmetry and Homogeneity

## Supervised Learning



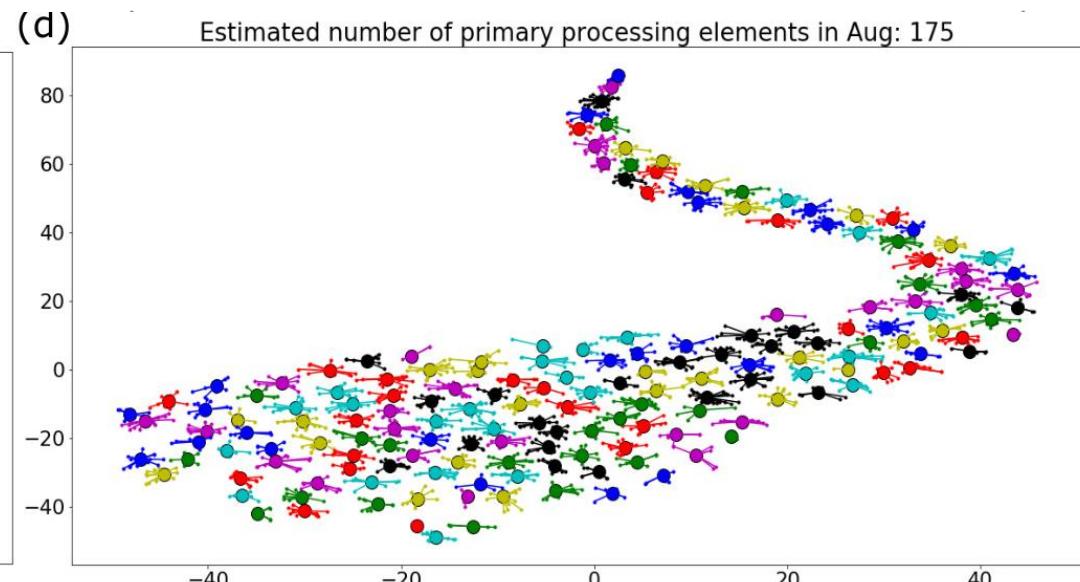
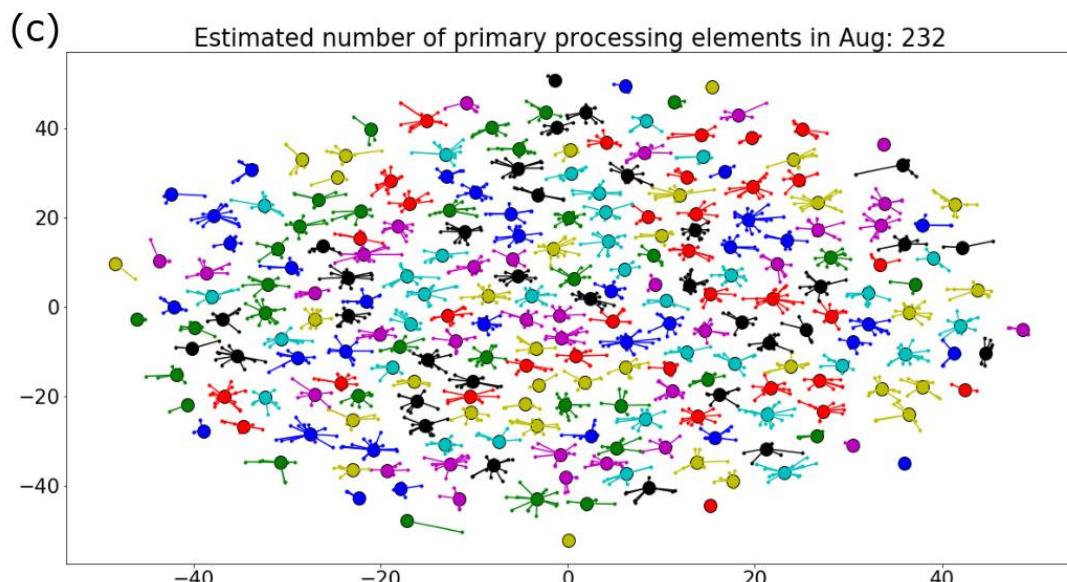
# Why Adversarial Interaction Creates Non-Homogeneous Patterns: A Pseudo-Reaction Diffusion Model for Turing Instability

- Adversarial Interaction
  - Generative Adversarial Networks (GANs)
  - Application of conditional GANs
- Non-Homogeneous Patterns
  - Homogeneous patterns
  - Supervised learning
- Reaction-Diffusion
  - Turing's RD model (1952)
  - Gray-Scott RD model (1984)
- Turing Instability
  - Reaction dynamics
  - Diffusion dynamics

# Breakdown of Symmetry and Homogeneity

## Adversarial Regularization

$$\begin{aligned}\mathcal{L}_{aug} (\mathbf{U}, \mathbf{V}, \mathbf{W}, \mathbf{a}) = & \underbrace{\frac{1}{2} \left\| \frac{1}{\sqrt{d_{out}m}} \mathbf{V} \sigma(\mathbf{U} \mathbf{X}) - \mathbf{Y} \right\|_F^2}_{\mathcal{L}_{sup}} \\ & - \frac{1}{m \sqrt{d_{out}}} \sum_{p=1}^n \mathbf{a}^T \sigma(\mathbf{W} \mathbf{V} \sigma(\mathbf{U} \mathbf{x}_p))\end{aligned}$$



# Breakdown of Symmetry and Homogeneity

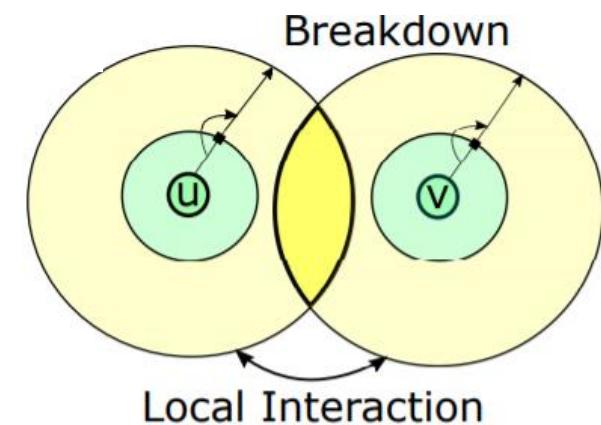
Theorem (Breakdown of Symmetry and Homogeneity)

Suppose **Assumption 1** holds. Let us i.i.d. initialize  $u_j, w_r \sim \mathcal{N}(0, I)$  and sample  $v_j, a_r$  uniformly from  $\{+1, -1\}$  for  $j, r \in [m]$ . Let  $\|x_p\|_2 = 1$  for all  $p \in [n]$ . If we choose  $\|\mathbf{w}\|_2 \leq L \leq \mathcal{O} \left( \frac{\epsilon \sqrt{m}}{\kappa n \sqrt{2 \log(2/\delta)}} \right)$ ,  $\kappa = \mathcal{O}(\kappa^\infty)$  where

$\kappa^\infty$  denotes the condition number of  $\mathcal{H}^\infty$ , and define  $\mu \triangleq \frac{Ln \sqrt{2 \log(2/\delta)}}{\sqrt{m}}$ , then with probability at least  $1 - \delta$ , we obtain the following:

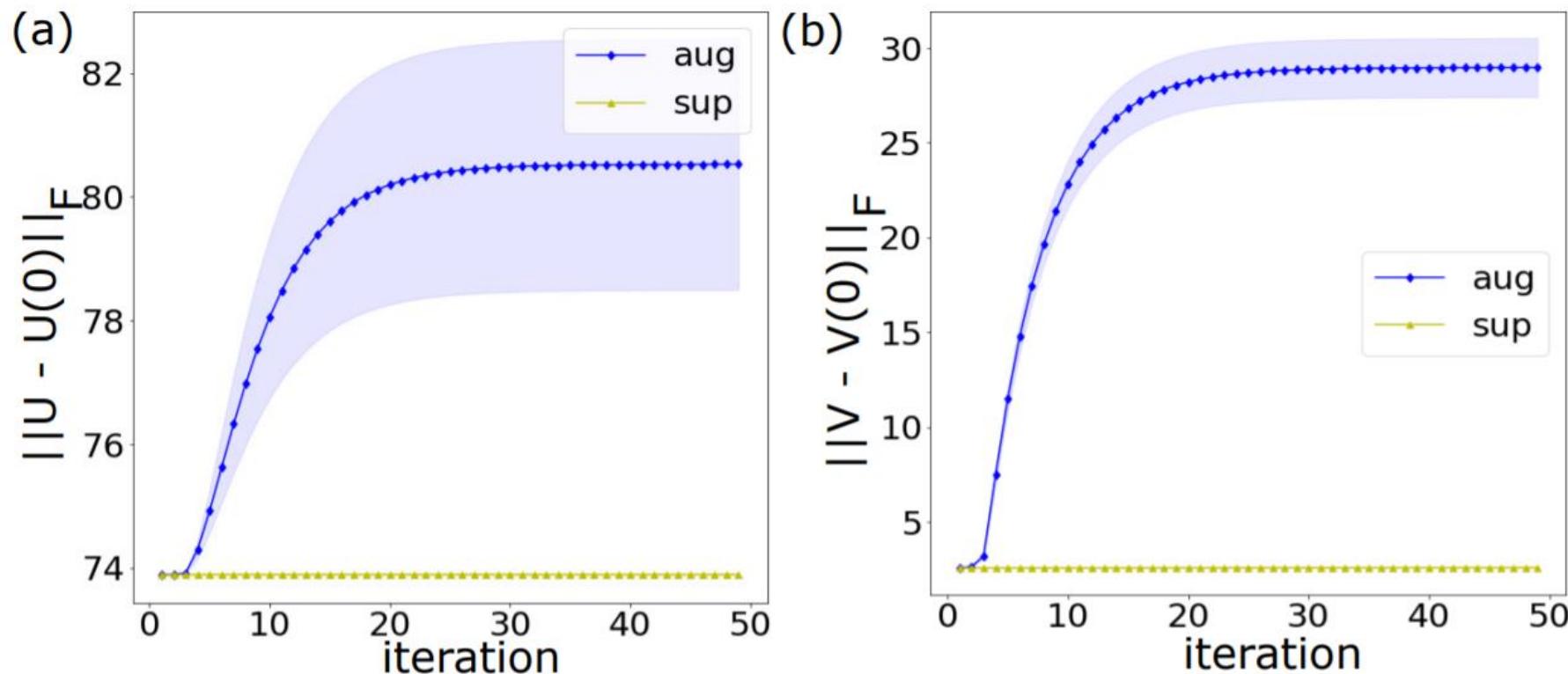
$$\|\mathbf{u}_j(t) - \mathbf{u}_j(0)\|_2 \leq \mathcal{O} \left( \frac{n^{3/2}}{\sqrt{m} \lambda_0 \delta} + \left( \frac{\mu (1 + \kappa \sqrt{n})}{\sqrt{m}} \right) t \right),$$

$$\|\mathbf{U}(t) - \mathbf{U}(0)\|_F \leq \mathcal{O} \left( \frac{n^{3/2}}{\lambda_0 \delta} + \mu (1 + \kappa \sqrt{n}) t \right).$$



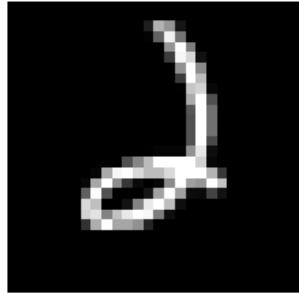
# Breakdown of Symmetry and Homogeneity

## Adversarial Regularization

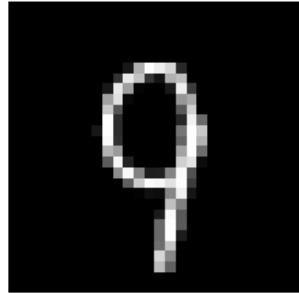


# Interpretability of Supervised Learning

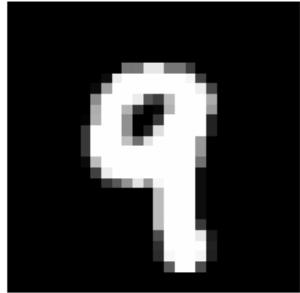
Label: 2



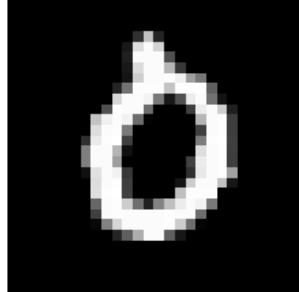
Label: 9



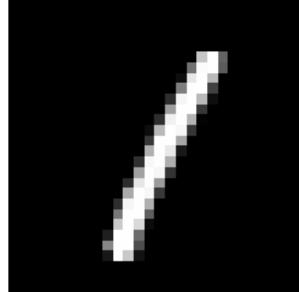
Label: 9



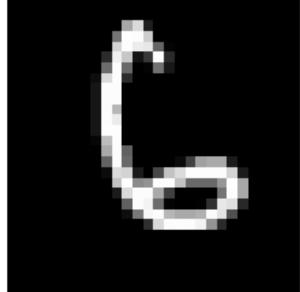
Label: 0



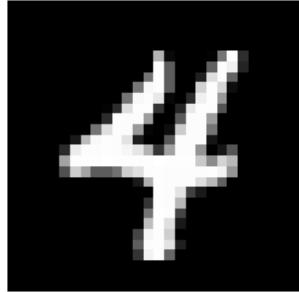
Label: 1



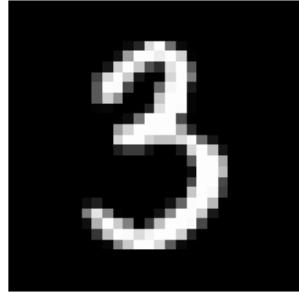
Label: 6



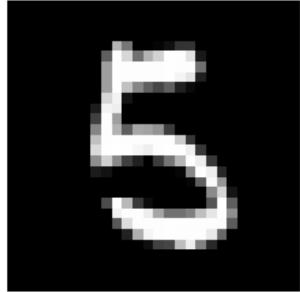
Label: 4



Label: 3

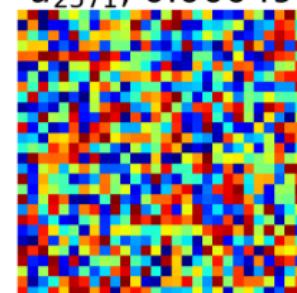


Label: 5

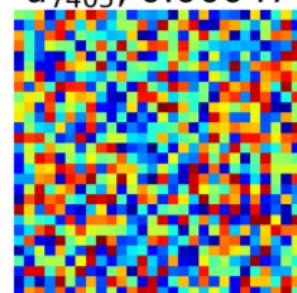


(a)

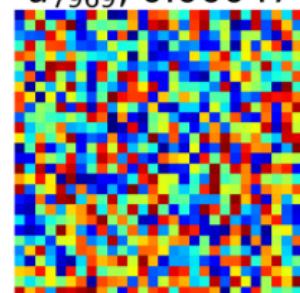
$u_{2571}, 0.00049$



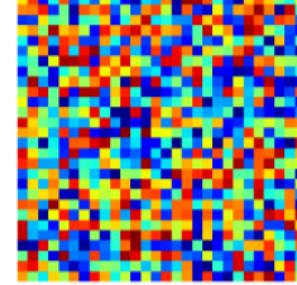
$u_{7405}, 0.00047$



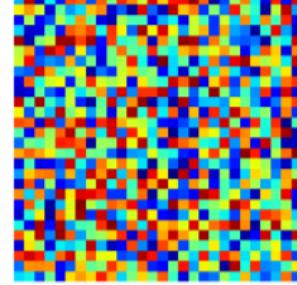
$u_{7969}, 0.00047$



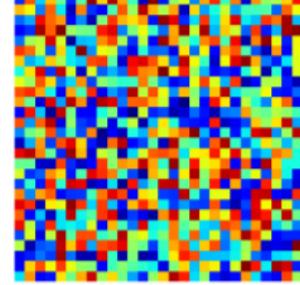
$u_{1076}, 0.00047$



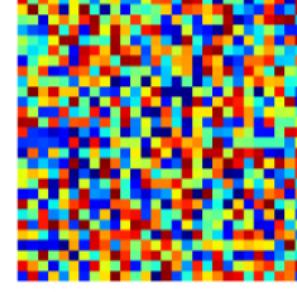
$u_{55}, 0.00047$



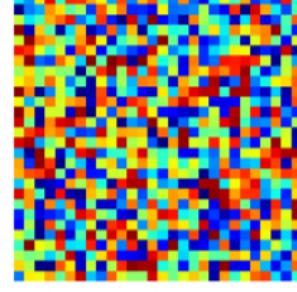
$u_{6704}, 0.00047$



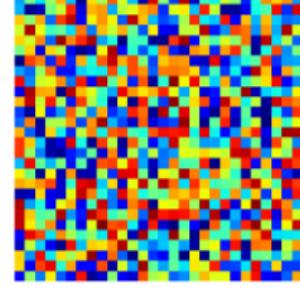
$u_{3595}, 0.00047$



$u_{177}, 0.00047$

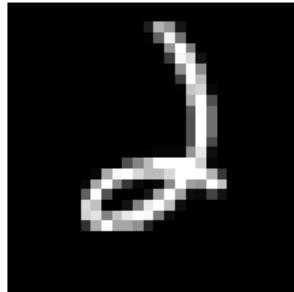


$u_{567}, 0.00047$

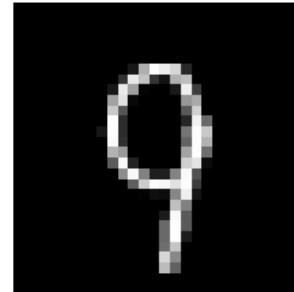


# Interpretability of Adversarial Regularization

Label: 2



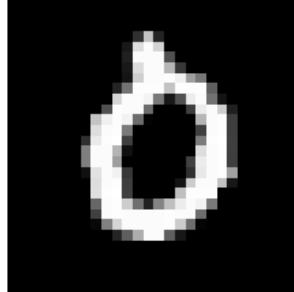
Label: 9



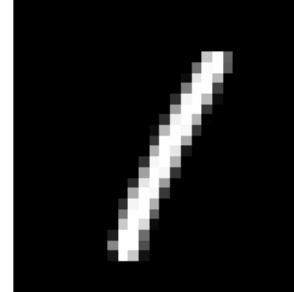
Label: 9



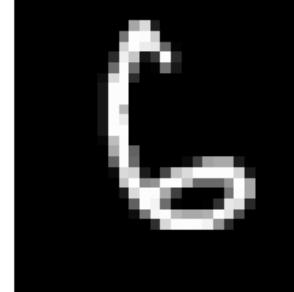
Label: 0



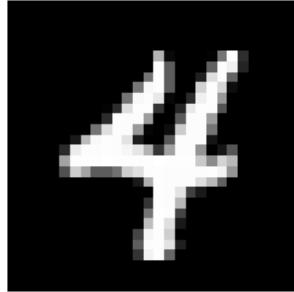
Label: 1



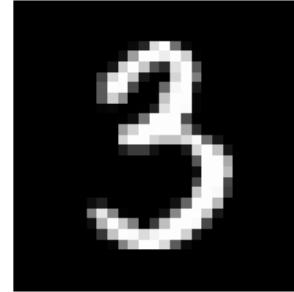
Label: 6



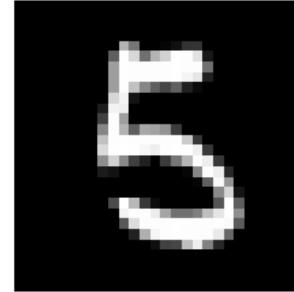
Label: 4



Label: 3

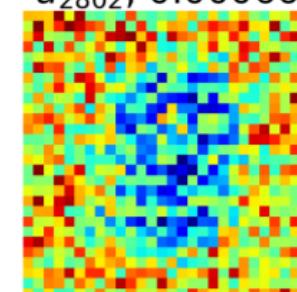


Label: 5

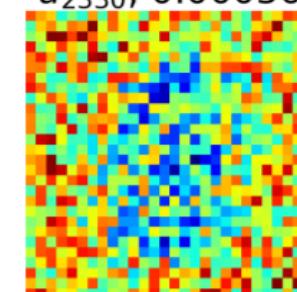


(b)

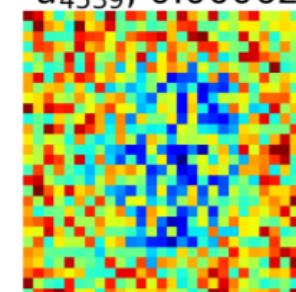
$u_{2802}, 0.00066$



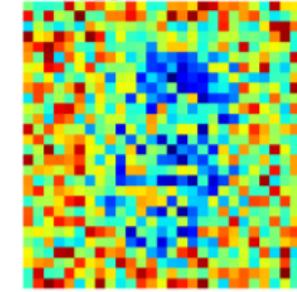
$u_{2330}, 0.00056$



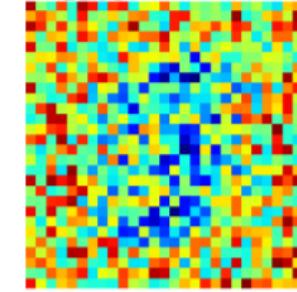
$u_{4539}, 0.00062$



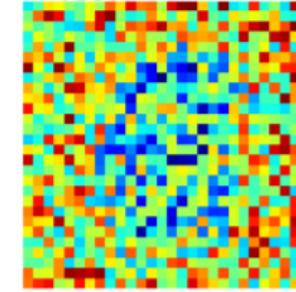
$u_{7615}, 0.00058$



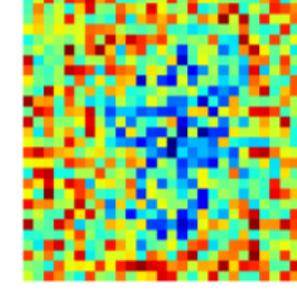
$u_{1368}, 0.00058$



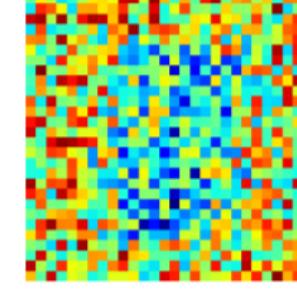
$u_{7597}, 0.00056$



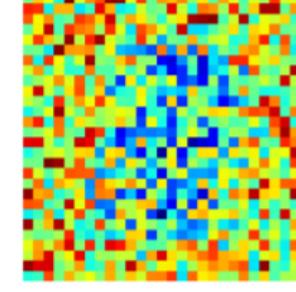
$u_{3462}, 0.00058$



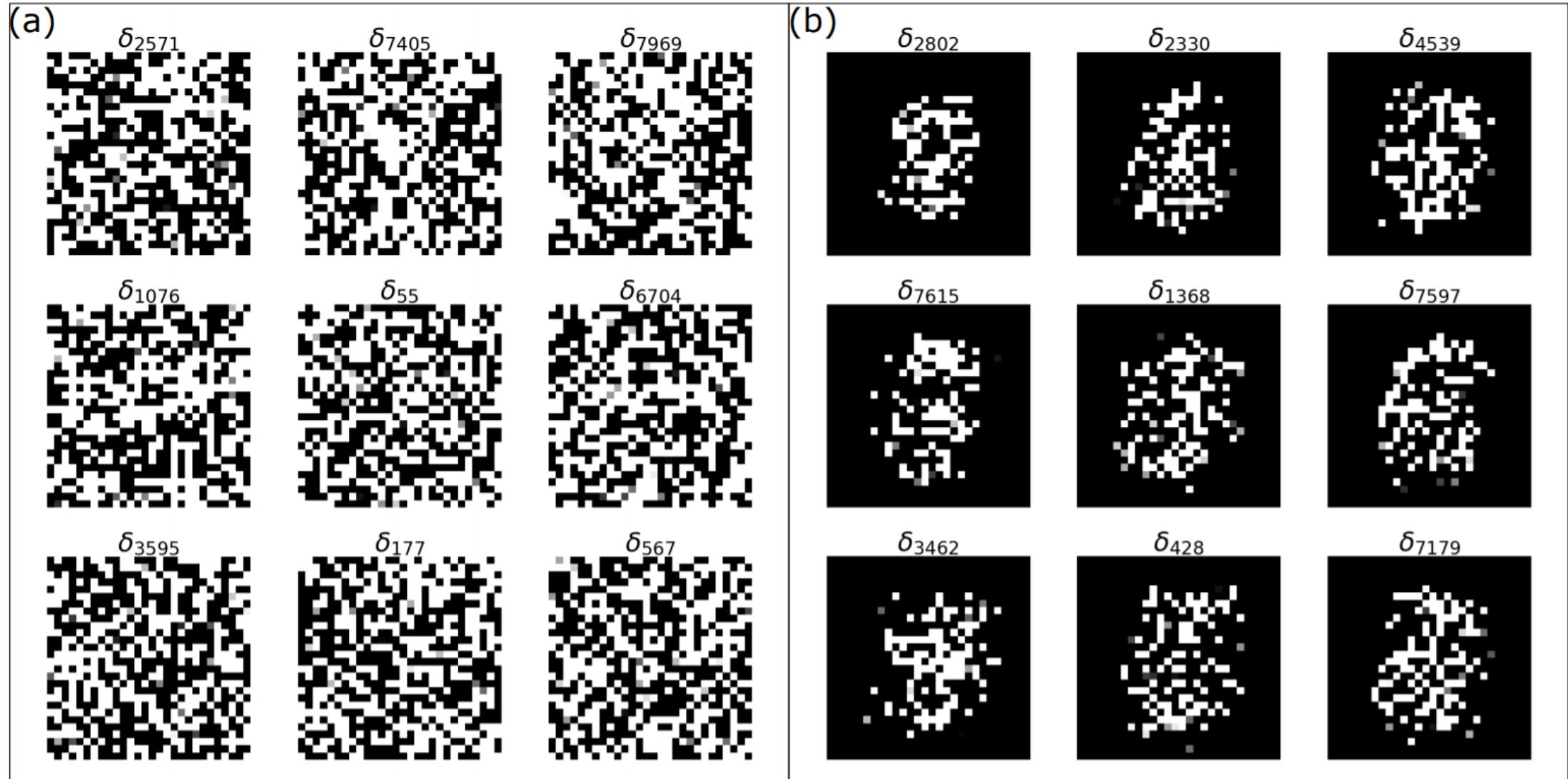
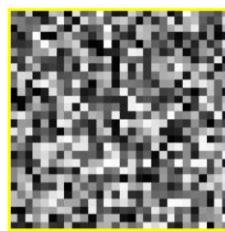
$u_{428}, 0.00054$



$u_{7179}, 0.00054$



# Interpretability of Adversarial Regularization



# Jointly Training Both Layers

$$\mathfrak{R}_j^{\mathbf{u}}(\mathbf{u}_j, \mathbf{v}_j) = \mathcal{O}\left(nd_{in}\sqrt{\frac{d_{out}}{m}}\right)$$

$$\mathfrak{D}_j^{\mathbf{u}}(\nabla^2 \mathbf{u}_j) = \mathcal{O}\left(nm^2 d_{in} d_{out}^{3/2}\right)$$

$$\mathfrak{R}_j^{\mathbf{v}}(\mathbf{u}_j, \mathbf{v}_j) = \mathcal{O}\left(nd_{in}\sqrt{\frac{d_{out}}{m}}\right)$$

$$\mathfrak{D}_j^{\mathbf{v}}(\nabla^2 \mathbf{v}_j) = \mathcal{O}\left(nm^2 d_{in} d_{out}^{1/2}\right)$$

# Bernoulli Differential Equation

$$\frac{dx(t)}{dt} + P(t)x(t) = Q(t)x^n(t) \text{ for } n \in \mathbb{R} \setminus \{0, 1\}$$

## Natural Phenomena

- Population Growth
- Tumour Growth
- Fermi-Dirac Statistics
- Diffusion of Innovation

$$\frac{dP}{dt} = rP \left(1 - \frac{P}{K}\right)$$

# Bernoulli Differential Equation

$$\frac{dx(t)}{dt} + P(t)x(t) = Q(t)x^n(t) \text{ for } n \in \mathbb{R} \setminus \{0, 1\}$$

## Natural Phenomena

- Population Growth
- Tumour Growth
- Fermi-Dirac Statistics
- Diffusion of Innovation

$$\frac{dP}{dt} = rP \left(1 - \frac{P}{K}\right)$$

## Adversarial Regularization

$$\frac{d\psi}{dt} \leq r\psi^{1/2} \left(1 - \frac{\psi^{1/2}}{K}\right)$$

# Why Adversarial Interaction Creates Non-Homogeneous Patterns: A Pseudo-Reaction Diffusion Model for Turing Instability

- Adversarial Interaction
  - Generative Adversarial Networks (GANs)
  - Application of conditional GANs
- Non-Homogeneous Patterns
  - Homogeneous patterns
  - Supervised learning
- Reaction-Diffusion
  - Turing's RD model (1952)
  - Gray-Scott RD model (1984)
- Turing Instability
  - Reaction dynamics
  - Diffusion dynamics

# Interesting Problems

- Why does adversarial interaction create non-homogeneous patterns?
- How do we predict the shape of the non-homogeneous patterns?
- What happens to the parameters if the supervised cost is augmented in the discriminator objective?
- Is gradient descent or backpropagation the optimal learning algorithm?