

DEMYSTIFYING OPTIMIZATION AND GENERALIZATION OF DEEP LEARNING VIA FEATURE LEARNING THEORY

Andi Han (Lecturer, University of Sydney), Wei Huang (Research Scientist, RIKEN AIP & ISM)

Tutorial @ AJCAI 2025, Canberra

SIGNIFICANT SUCCESS OF DEEP LEARNING

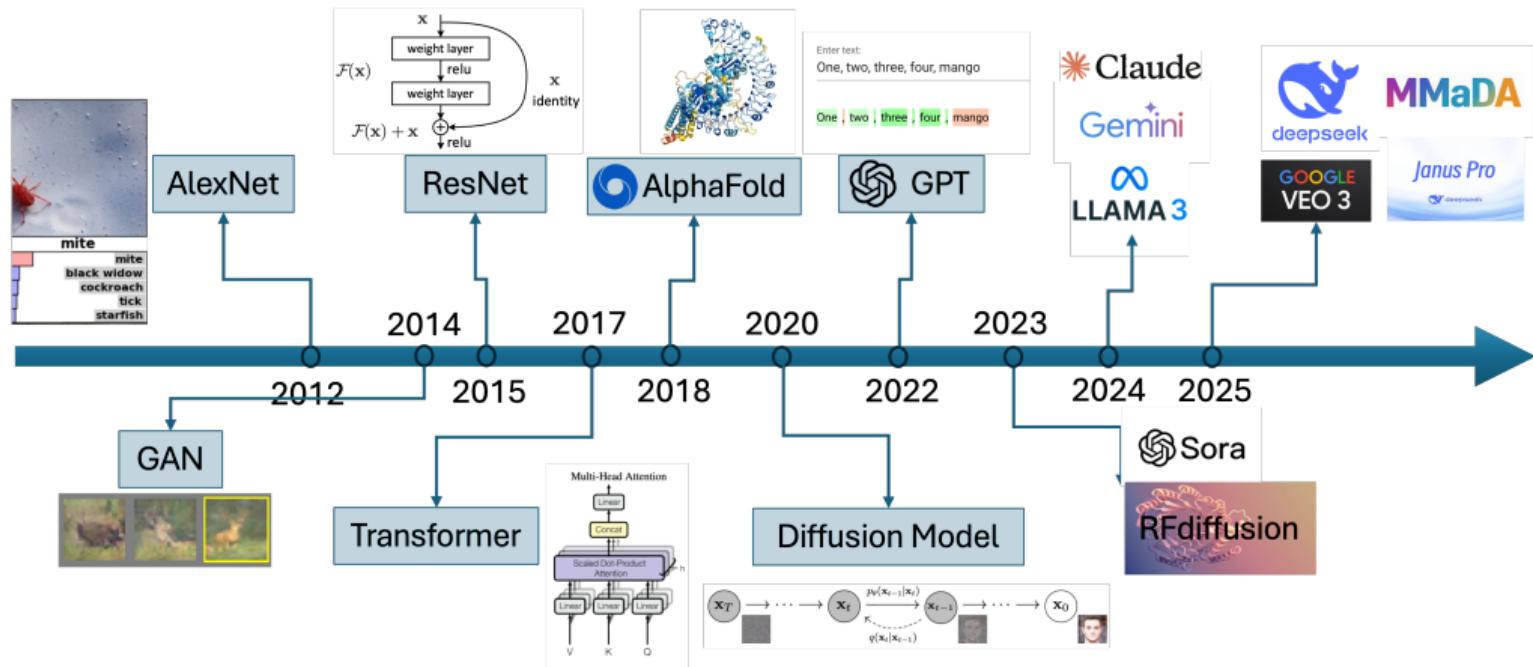
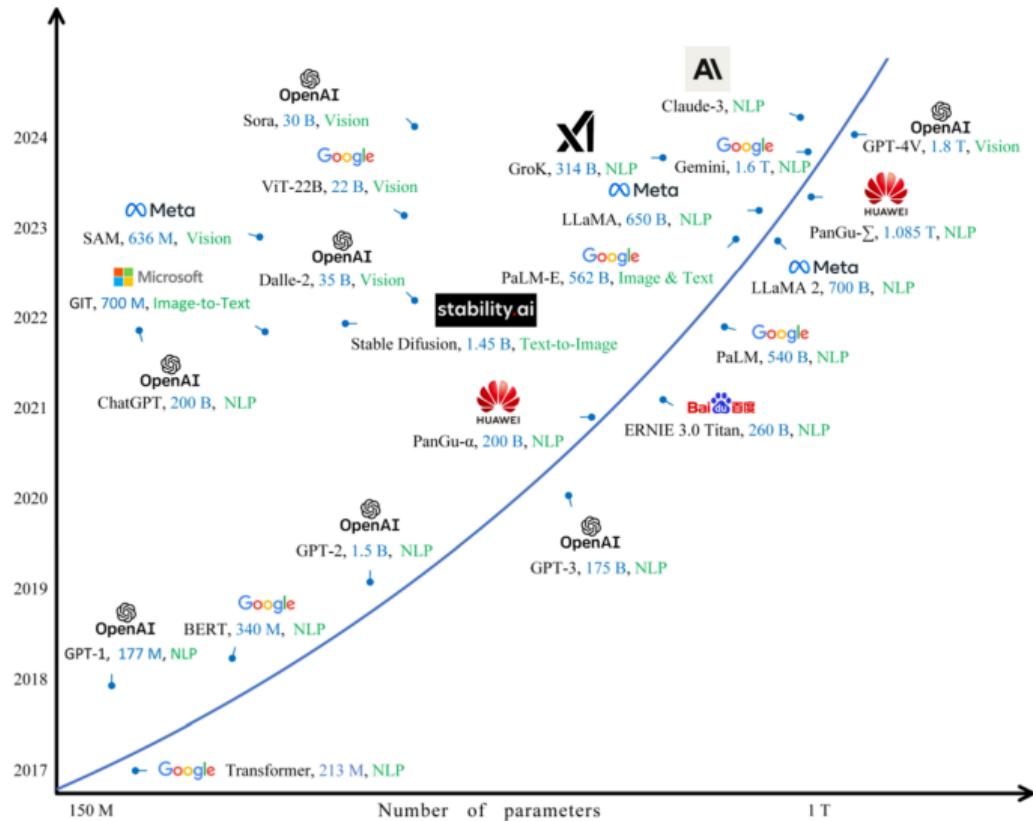


Image credits (left → right): Krizhevsky et al. 2012; Goodfellow et al. 2014; He et al. 2015; Vaswani et al. 2017;
<https://alphafold.ebi.ac.uk/>; Ho et al. 2020; <https://huggingface.co/blog/alonsosilva/nexttokenprediction>; Watson et al. 2023

MODEL SIZE AND COMPLEXITY GROWTH



YET WE UNDERSTAND LESS AND LESS...

- Transparency
- Robustness
- Privacy, Fairness, Biases



YET WE UNDERSTAND LESS AND LESS...

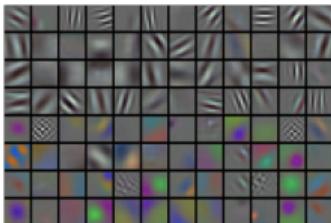
- Transparency
- Robustness
- Privacy, Fairness, Biases

Underlying principles of deep learning?

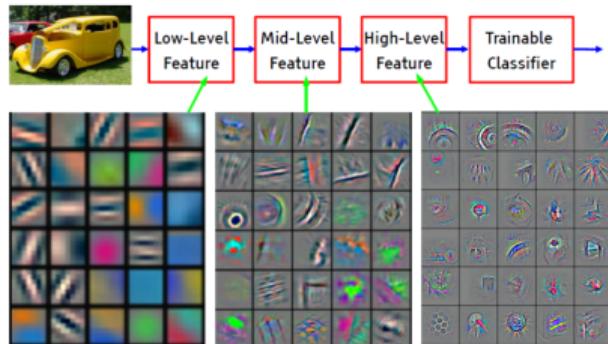


xkcd: machine learning

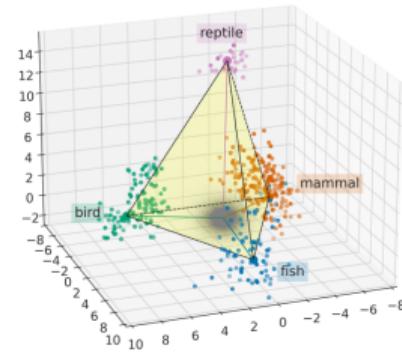
FEATURE LEARNING AT THE CORE OF DEEP LEARNING



AlexNet 1st layer (Krizhevsky et al. 2012)

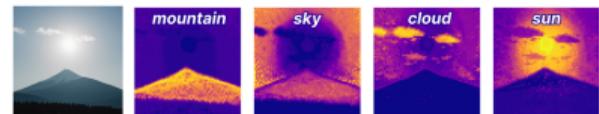


Hierarchical feature learning (LeCun 2015)



Concept features in Gemma (Park et al. 2025)

"a mountain
in the distance." →



Cross-modal features (Helbling et al. 2025)

TUTORIAL BREAKDOWN

Goal

- Introduce a theoretical sandbox to understand deep learning via feature learning
- Bridge empirical phenomena and theoretical insights on optimization and generalization

Outline

1. Deep learning benefits from feature learning
2. A signal-and-noise data model
3. Benign Overfitting with Feature Learning
4. Feature Learning under Different Training Strategies
4. Feature Learning in Foundation Generative Models
5. Conclusions & outlook

WHY DEEP LEARNING BENEFITS FROM FEATURE LEARNING

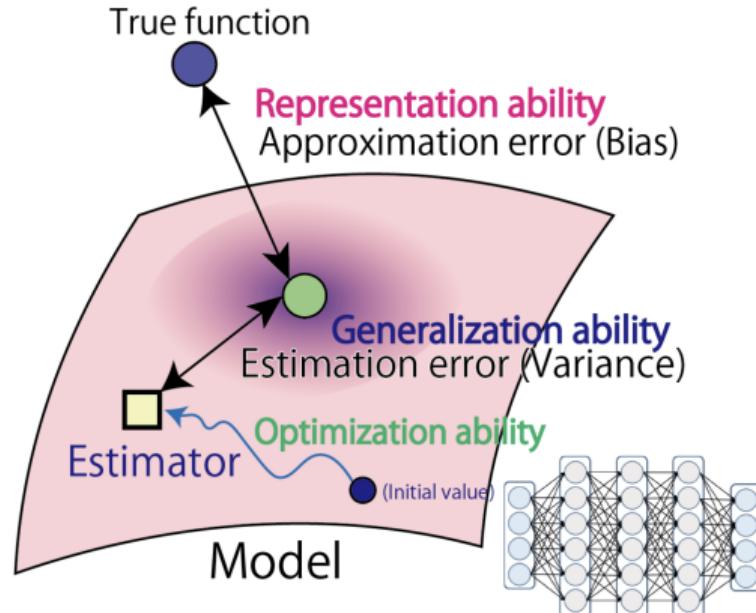
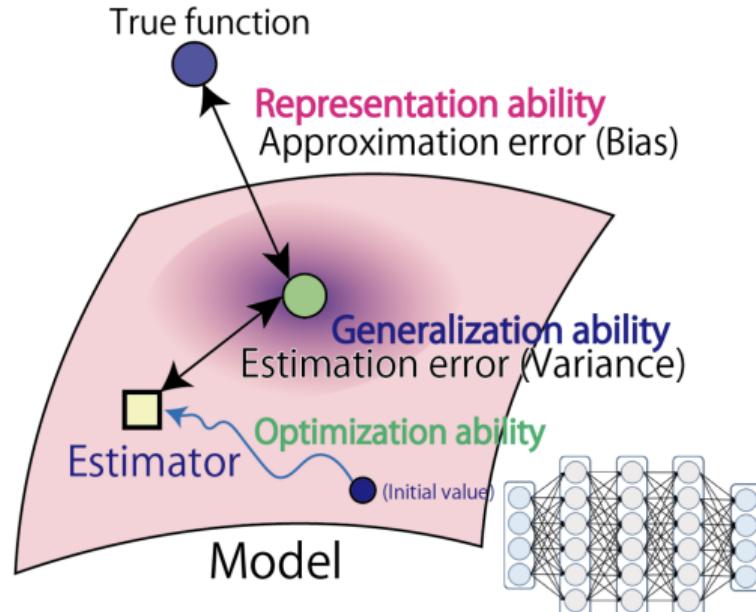


Image credits: Suzuki 2024

WHY DEEP LEARNING BENEFITS FROM FEATURE LEARNING



Feature Learning Affects ALL!

Image credits: Suzuki 2024

A REPRESENTATION PERSPECTIVE

Two-layer Neural Network and Kernel

$$f(\mathbf{x}) = \frac{1}{\sqrt{N}} \sum_{i=1}^N a_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle)$$

A REPRESENTATION PERSPECTIVE

Two-layer Neural Network and Kernel

$$f(\mathbf{x}) = \frac{1}{\sqrt{N}} \sum_{i=1}^N a_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle)$$

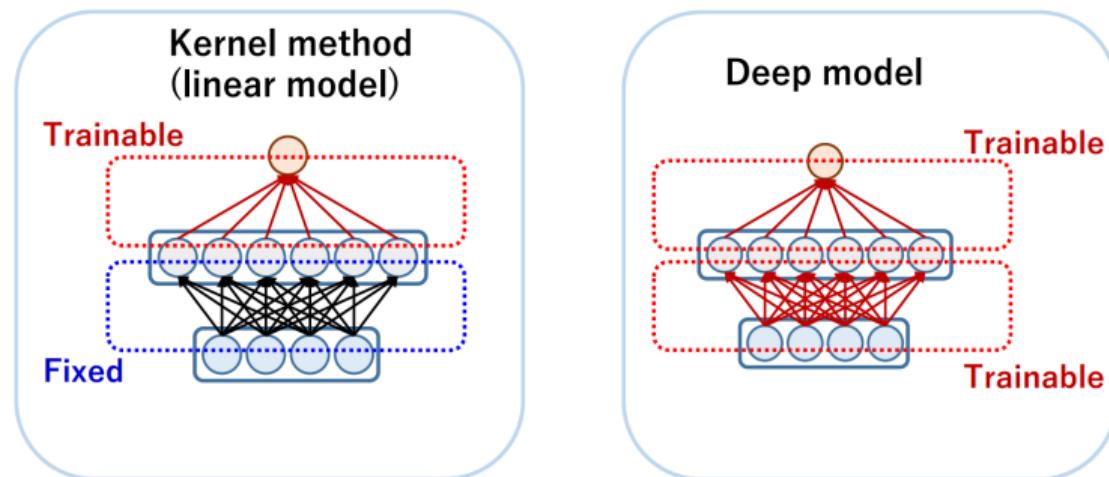


Image credits: Suzuki 2024

A REPRESENTATION PERSPECTIVE

Target function $f^*(\mathbf{x}) = \sigma^*(\langle \mathbf{x}, \boldsymbol{\beta}_* \rangle)$. $y_i = f^*(\mathbf{x}_i) + \epsilon_i$, $\mathbf{x}_i \sim \mathcal{N}(0, \mathbf{I})$, $\epsilon_i \sim \mathcal{N}(0, \sigma_\epsilon^2)$.

Let $\hat{f}_\lambda = \arg \min_f \frac{1}{n} \sum_{i=1}^n (f(\mathbf{x}_i) - y_i)^2 + \frac{\lambda}{N} \|\mathbf{a}\|^2$.

- (*Kernel*) Let \mathbf{w}_i fixed at initialization \mathbf{w}_i^0

$$\inf_{\lambda} \mathcal{R}(\hat{f}_\lambda) \geq \|\mathsf{P}_{>1} f^*\|_{L^2}^2 + o(1)$$

- (*Neural Network*) Let \mathbf{w}_i be one-step gradient update from \mathbf{w}_i^0 ,

$$\mathcal{R}(\hat{f}_\lambda) < \|\mathsf{P}_{>1} f^*\|_{L^2}^2$$

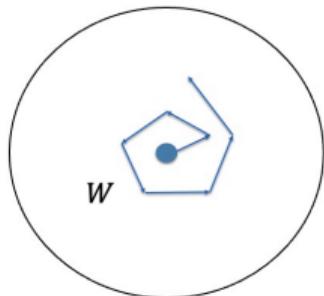
where $\mathcal{R}(\hat{f}) = \mathbb{E}_{\mathbf{x}} (\hat{f}(\mathbf{x}) - f^*(\mathbf{x}))^2$ is the prediction risk, and $f^*(\mathbf{x}) = \mu_0^* + \mu_1^* \langle \mathbf{x}, \boldsymbol{\beta}_* \rangle + \mathsf{P}_{>1} f^*$.

Ba et al. 2022. "High-dimensional Asymptotics of Feature Learning: How One Gradient Step Improves the Representation".

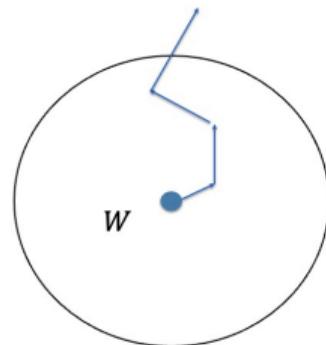
* Suppose $\sigma = \sigma^* = \tanh$.

A REPRESENTATION PERSPECTIVE

Feature Learning underlies the success of deep learning!



Lazy training



Feature learning

$$\|W(t) - W(0)\|_F = O\left(\frac{1}{\sqrt{N}}\right)$$

$$\|W(t) - W(0)\|_F = \Omega(1)$$

SIGNAL-NOISE DATA MODEL (A SANDBOX FOR FEATURE LEARNING)

Feature Decomposition: Data \approx Signal + Noise

SIGNAL-NOISE DATA MODEL (A SANDBOX FOR FEATURE LEARNING)

Feature Decomposition: Data \approx Signal + Noise

Signal-noise data model (Cao et al., 2022; Kou et al., 2023)

Data $x = [y\mu, \xi]$

- (Signal) μ is a signal vector, and y is the label
- (Noise) ξ is a random noise (commonly assumed to be Gaussian $\mathcal{N}(\mathbf{0}, \sigma_\xi^2 \mathbf{I})$)

SIGNAL-NOISE DATA MODEL (A SANDBOX FOR FEATURE LEARNING)

Feature Decomposition: Data \approx Signal + Noise

Signal-noise data model (Cao et al., 2022; Kou et al., 2023)

Data $x = [y\mu, \xi]$

- (Signal) μ is a signal vector, and y is the label
- (Noise) ξ is a random noise (commonly assumed to be Gaussian $\mathcal{N}(\mathbf{0}, \sigma_\xi^2 \mathbf{I})$)



Today is Sunday

■ Signal
■ Noise

Demo from Imagenet

BENIGN OVERFITTING WITH FEATURE LEARNING

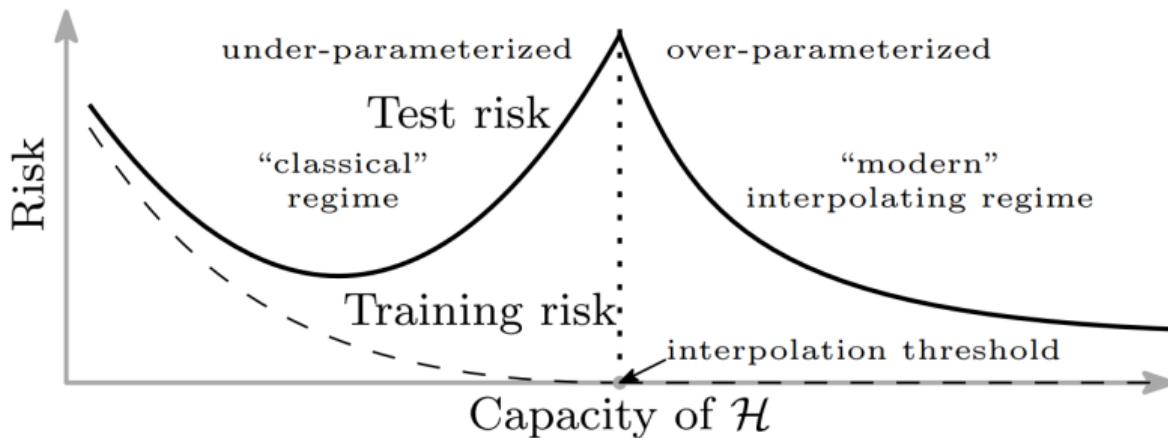


Image credit: Belkin et al. 2019. "Reconciling modern machine learning practice and the bias-variance trade-off"

BENIGN OVERFITTING WITH FEATURE LEARNING

BENIGN OVERFITTING WITH FEATURE LEARNING

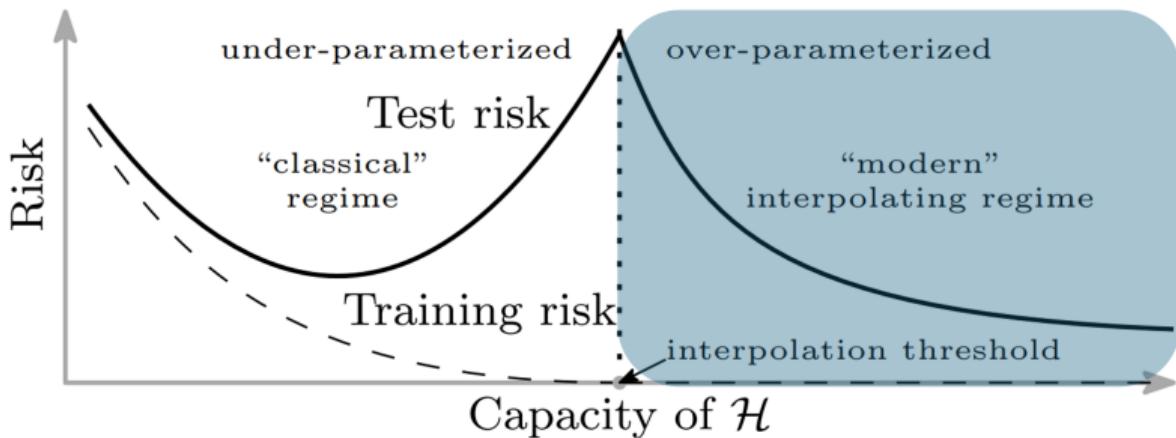


Image credit: Belkin et al. 2019. "Reconciling modern machine learning practice and the bias-variance trade-off"

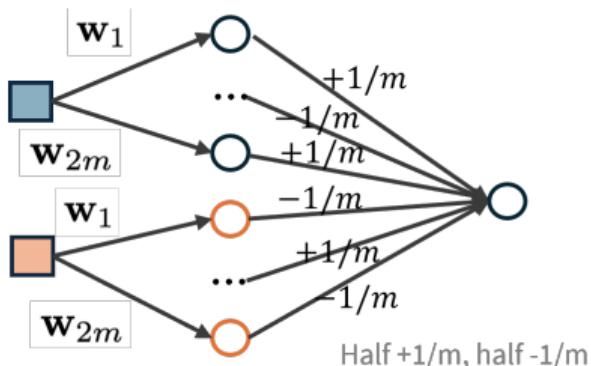
BENIGN OVERFITTING WITH FEATURE LEARNING (MODEL SETUP)

Two-layer Convolutional Neural Network

$$f(\mathbf{W}, \mathbf{x}) = F_1(\mathbf{W}_1, \mathbf{x}) - F_{-1}(\mathbf{W}_{-1}, \mathbf{x})$$

$$\text{where } F_j(\mathbf{W}_j, \mathbf{x}) = \frac{1}{m} \sum_{r=1}^m \sigma(\langle \mathbf{w}_{j,r}, \mathbf{x}^{(1)} \rangle) + \frac{1}{m} \sum_{r=1}^m \sigma(\langle \mathbf{w}_{j,r}, \mathbf{x}^{(2)} \rangle)$$

$$\text{where } \mathbf{x} = [\mathbf{x}^{(1)\top}, \mathbf{x}^{(2)\top}]^\top = [y\boldsymbol{\mu}, \boldsymbol{\xi}]$$



BENIGN OVERFITTING WITH FEATURE LEARNING (TRAINING SETUP)

- **Training Data** $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$:
 - Binary classification: $y = \pm 1$ with equal chance
 - $\mathbf{x}_i = [y_i \boldsymbol{\mu}, \boldsymbol{\xi}_i]$, with fixed signal $\boldsymbol{\mu}$ and random noise $\boldsymbol{\xi}_i \sim \mathcal{N}(\mathbf{0}, \sigma_\xi^2 \mathbf{I}_d)$
 - Define SNR = $\|\boldsymbol{\mu}\| / (\sigma_\xi \sqrt{d})$
- **Training Loss**:

$$L_S(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \ell(y_i f(\mathbf{W}, \mathbf{x}_i)), \quad \ell(z) = \log(1 + \exp(-z))$$

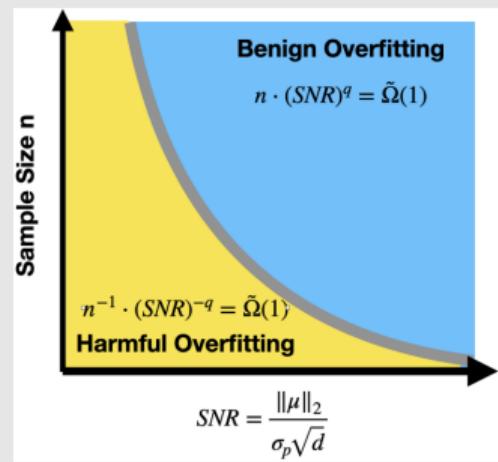
- **Test Loss**: $L_D(\mathbf{W}) = \mathbb{E}_{(\mathbf{x}, y)} [\ell(y f(\mathbf{W}, \mathbf{x}))]$.
- **Training Algorithm**: $\mathbf{W}^{(t+1)} = \mathbf{W}^{(t)} - \eta \nabla L_S(\mathbf{W}), \quad \mathbf{W}^{(0)} \sim \mathcal{N}(\mathbf{0}, \sigma_0^2 \mathbf{I})$

BENIGN OVERFITTING WITH FEATURE LEARNING (MAIN RESULTS)

Feature learning under $(\text{ReLU})^q$ (Cao et al., 2022, Theorem 4.3 & 4.4)

Suppose $\sigma = (\text{ReLU})^q$, ($q > 2$), under *over-parameterization^a* and *small initialization^b*, there exists an iterate $\mathbf{W}^{(t)}$ with $L_S(\mathbf{W}^{(t)}) \leq \varepsilon$ and

- **Benign Overfitting:** When $n \cdot \text{SNR}^q = \tilde{\Omega}(1)$,
 - $L_D(\mathbf{W}^{(t)}) \leq 6\varepsilon + \exp(-n^2)$
- **Harmful Overfitting:** When $n^{-1} \cdot \text{SNR}^{-q} = \tilde{\Omega}(1)$,
 - $L_D(\mathbf{W}^{(t)}) \geq 0.1$



^a Large d relative to n

^b Feature learning regime

BENIGN OVERFITTING WITH FEATURE LEARNING (KEY IDEA)

Signal-Noise Decomposition for feature learning

$$\mathbf{w}_{j,r}^{(t)} = \mathbf{w}_{j,r}^{(0)} + j \cdot \gamma_{j,r}^{(t)} \cdot \|\boldsymbol{\mu}\|^{-1} \cdot \boldsymbol{\mu} + \sum_{i=1}^n \rho_{j,r,i}^{(t)} \cdot \|\boldsymbol{\xi}_i\|^{-2} \cdot \boldsymbol{\xi}_i$$

such that $\gamma_{j,r}^{(t)} \approx \langle \mathbf{w}_{j,r}^{(t)}, \boldsymbol{\mu} \rangle$ (signal learning) and $\rho_{j,r,i}^{(t)} \approx \langle \mathbf{w}_{j,r}^{(t)}, \boldsymbol{\xi}_i \rangle$ (noise memorization)

- **Benign Overfitting:** learn signal and ignore noise

$$\max_r \gamma_{j,r}^{(t)} \geq C_1 > 0, \quad \max_{j,r,i} |\rho_{j,r,i}^t| \approx 0 \quad (\text{Signal dominates } \odot)$$

- **Harmful Overfitting:** memorize noise and ignore signal

$$\max_r \rho_{y_i,r,i}^{(t)} \geq C_2 > 0, \quad \max_{j,r} \gamma_{j,r}^{(t)} \approx 0 \quad (\text{Noise dominates } \odot)$$

BENIGN OVERFITTING WITH FEATURE LEARNING (RELU)

Feature learning with ReLU (Kou et al., 2023, Theorem 4.2)

Suppose $\sigma = \text{ReLU}$, under *over-parameterization*^a and *small initialization*^b, and $n \cdot \text{SNR}^2 = o(1)$, there exists an iterate $\mathbf{W}^{(t)}$ with $L_S(\mathbf{W}^{(t)}) \leq \varepsilon$:

- **Benign Overfitting:** When $n\|\boldsymbol{\mu}\|^4 \geq C_1\sigma_\xi^4 d$,
 - $L_D^{0-1}(\mathbf{W}^{(t)}) \leq \exp\left(-n\|\boldsymbol{\mu}\|^4/(C_2\sigma_\xi^4 d)\right)$
- **Harmful Overfitting:** When $n\|\boldsymbol{\mu}\|^4 \leq C_3\sigma_\xi^4 d$,
 - $L_D^{0-1}(\mathbf{W}^{(t)}) \geq 0.1$

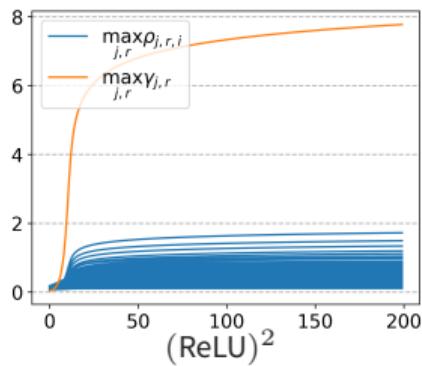
where $L_D^{0-1}(\mathbf{W}) = \mathbb{P}_{(\mathbf{x},y)}[y \neq \text{sign}(f(\mathbf{W}, \mathbf{x}))]$ is the test error.

Key differences to $(\text{ReLU})^q$: Constant separation and low-SNR regime.

BENIGN OVERRFITTING WITH FEATURE LEARNING (A COMPARISON OF FEATURE LEARNING)

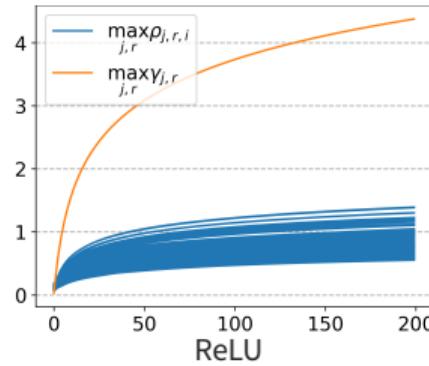
(ReLU)^q: Polynomial growth

$$\begin{aligned}\gamma_{j,r}^{(t+1)} &\approx (1 - \eta_\gamma)(\gamma_{j,r}^{(t)})^{q-1} \\ \rho_{j,r,i}^{(t+1)} &\approx (1 - \eta_\xi)(\rho_{j,r,i}^{(t)})^{q-1}\end{aligned}$$



ReLU: Linear growth

$$\begin{aligned}\gamma_{j,r}^{(t+1)} &\approx \gamma_{j,r}^{(t)} + \eta'_\gamma \\ \rho_{j,r,i}^{(t+1)} &\approx \rho_{j,r,i}^{(t)} + \eta'_\xi\end{aligned}$$



FEATURE LEARNING WITH LABEL NOISE

Label noise is common: the observed label \tilde{y} is not equal to the ground-truth label y !

Dataset	MNIST [26]	CIFAR-10 [27]	CIFAR-100 [36]
Instance			
Given Label	“8”	“cat”	“Lobstar”
Actual Label	“9”	“frog”	“Crab”
Dataset	Caltech-256 [37]	ImageNet [38]	QuickDraw [39]
Instance			
Given Label	“Dolphin”	“White stork”	“Tiger”
Actual Label	“Kayak”	“Black stork”	“eye”

Image credit: Bhatt et al. 2024

BENIGN OVERFITTING WITH FEATURE LEARNING (LABEL NOISE)

Feature learning under label noise ([Han et al., 2025b](#), Theorem 4.2 & 4.4)

Observed label $\tilde{y} \neq y$ (with $\mathbb{P}(\tilde{y} \neq y) = \tau$), and $\sigma = \text{ReLU}$, $n \cdot \text{SNR}^2 = \Theta(1)$

Two-stage behavior

BENIGN OVERFITTING WITH FEATURE LEARNING (LABEL NOISE)

Feature learning under label noise (Han et al., 2025b, Theorem 4.2 & 4.4)

Observed label $\tilde{y} \neq y$ (with $\mathbb{P}(\tilde{y} \neq y) = \tau$), and $\sigma = \text{ReLU}$, $n \cdot \text{SNR}^2 = \Theta(1)$

Two-stage behavior

- **Stage I (Model fits clean data):** there exists T_1 s.t.
 - Model learns more signal than noise, i.e., $\gamma_{j,r}^{(T_1)} > \rho_{\tilde{y}_i,r,i}^{(T_1)}$
 - For all clean samples: $\tilde{y}_i f(\mathbf{W}^{(T_1)}, \mathbf{x}_i) \geq 0$. For all noisy samples: $\tilde{y}_i f(\mathbf{W}^{(T_1)}, \mathbf{x}_i) < 0$
 - Early stopping works: $L_D^{0-1}(\mathbf{W}^{(T_1)}) \leq \exp(-\Omega(d/n))$

BENIGN OVERFITTING WITH FEATURE LEARNING (LABEL NOISE)

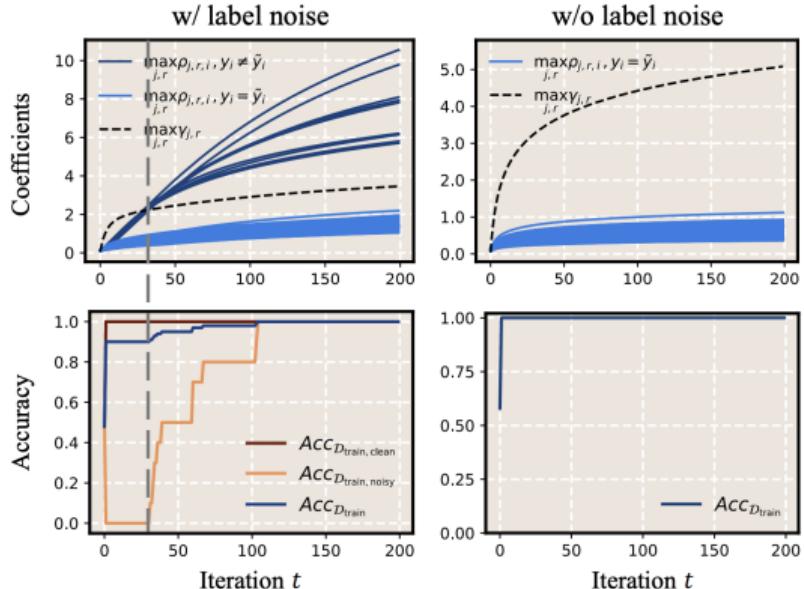
Feature learning under label noise (Han et al., 2025b, Theorem 4.2 & 4.4)

Observed label $\tilde{y} \neq y$ (with $\mathbb{P}(\tilde{y} \neq y) = \tau$), and $\sigma = \text{ReLU}$, $n \cdot \text{SNR}^2 = \Theta(1)$

Two-stage behavior

- **Stage I (Model fits clean data):** there exists T_1 s.t.
 - Model learns more signal than noise, i.e., $\gamma_{j,r}^{(T_1)} > \rho_{\tilde{y}_i,r,i}^{(T_1)}$
 - For all clean samples: $\tilde{y}_i f(\mathbf{W}^{(T_1)}, \mathbf{x}_i) \geq 0$. For all noisy samples: $\tilde{y}_i f(\mathbf{W}^{(T_1)}, \mathbf{x}_i) < 0$
 - Early stopping works: $L_D^{0-1}(\mathbf{W}^{(T_1)}) \leq \exp(-\Omega(d/n))$
- **Stage II (Model overfits noisy data):** there exists $t^* \geq T_1$ such that
 - For most if not all samples: $\tilde{y}_i f(\mathbf{W}^{(T_1)}, \mathbf{x}_i) \geq 0$.
 - Model fails to generalize: $L_D^{0-1}(\mathbf{W}^{(t^*)}) \geq 0.5\tau$

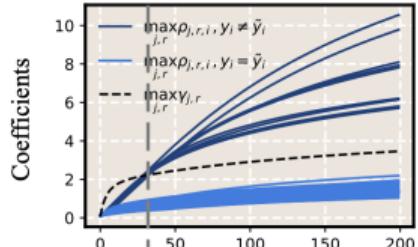
BENIGN OVERRFITTING WITH FEATURE LEARNING (LABEL NOISE)



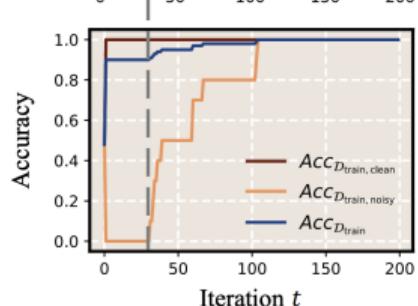
Synthetic data

BENIGN OVERRFITTING WITH FEATURE LEARNING (LABEL NOISE)

w/ label noise

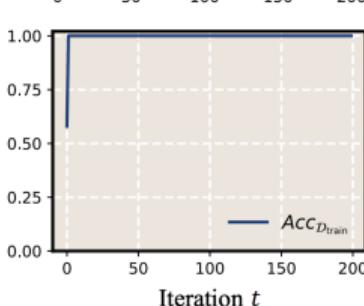
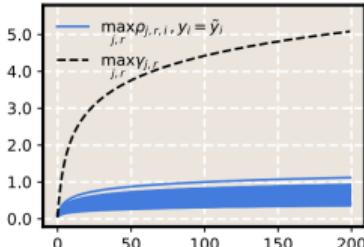


Accuracy



Iteration t

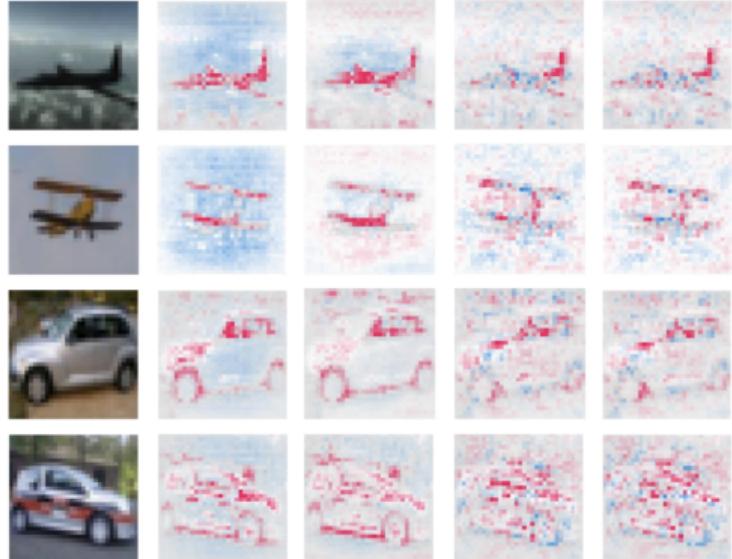
w/o label noise



Iteration t

Synthetic data

Feature Learning Process w/ Label Noise



Epoch 1

Epoch 41

Epoch 81

Epoch 121

VGG on CIFAR-10 with label noise.

BENIGN OVERFITTING WITH FEATURE LEARNING (TRANSFORMER)

The emergence of large language models (LLMs) is due to Transformers.

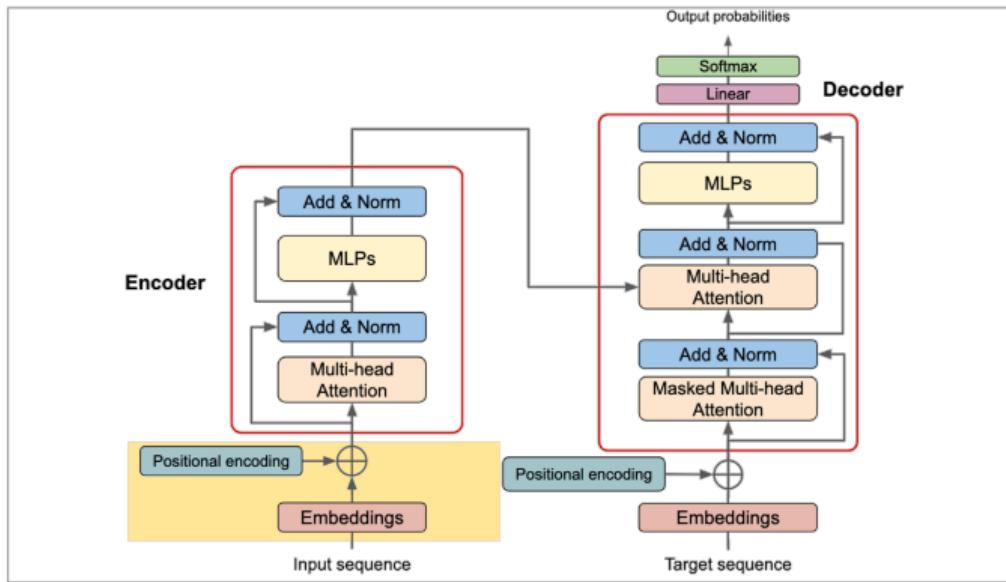


Image credit: <https://deeprunning.github.io/posts/001-transformer/>

BENIGN OVERFITTING WITH FEATURE LEARNING (TRANSFORMER)

To analyze benign overfitting, focus on the core mechanism: **attention**.

- Consider a **single-head** Transformer with *global average pooling* (Jiang et al., 2024):

$$f(\mathbf{X}) = \underbrace{\frac{1}{L} \sum_{\ell=1}^L}_{\text{Avg Pooling}} \underbrace{\text{Softmax}(\mathbf{x}^{(\ell)} \mathbf{W}_Q \mathbf{W}_K^\top \mathbf{X}^\top) \mathbf{X} \mathbf{W}_V \mathbf{w}_o}_{\text{Attention for token } l}$$

where $\mathbf{X} = [\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(L)}]^\top \in \mathbb{R}^{L \times d}$.

Training data (\mathbf{X}_i, y_i) :

- $\mathbf{x}^{(1)} = y\boldsymbol{\mu}$ (signal), $\mathbf{x}^{(2)}, \dots, \mathbf{x}^{(L)} \sim \mathcal{N}(\mathbf{0}, \sigma_\xi^2 \mathbf{I})$ (noise).

BENIGN OVERFITTING WITH FEATURE LEARNING (TRANSFORMER)

Feature learning in Transformers ([Jiang et al., 2024](#), Theorem 4.1 & 4.2)

There exists T , s.t. $L_S(\Theta^{(T)}) \approx 0$ and

BENIGN OVERFITTING WITH FEATURE LEARNING (TRANSFORMER)

Feature learning in Transformers (Jiang et al., 2024, Theorem 4.1 & 4.2)

There exists T , s.t. $L_S(\Theta^{(T)}) \approx 0$ and

- Benign overfitting: under condition $n \cdot \text{SNR}^2 = \Omega(1)$
 - Attention on signal: $\langle \mathbf{W}_Q^{(T)} \mathbf{x}^{(\ell)}, \mathbf{W}_K^{(T)} \boldsymbol{\mu} \rangle = \Omega(1)$, $\langle \mathbf{W}_Q^{(T)} \mathbf{x}^{(\ell)}, \mathbf{W}_K^{(T)} \boldsymbol{\xi} \rangle \approx 0$
 - Value focuses on signal: $\langle \mathbf{W}_V^{(T)} \mathbf{w}_o, \boldsymbol{\mu} \rangle > \langle \mathbf{W}_V^{(T)} \mathbf{w}_o, \boldsymbol{\xi} \rangle$
 - The test loss is nearly zero: $L_D(\Theta^{(T)}) \approx 0$

BENIGN OVERFITTING WITH FEATURE LEARNING (TRANSFORMER)

Feature learning in Transformers (Jiang et al., 2024, Theorem 4.1 & 4.2)

There exists T , s.t. $L_S(\Theta^{(T)}) \approx 0$ and

- **Benign overfitting:** under condition $n \cdot \text{SNR}^2 = \Omega(1)$
 - Attention on *signal*: $\langle \mathbf{W}_Q^{(T)} \mathbf{x}^{(\ell)}, \mathbf{W}_K^{(T)} \boldsymbol{\mu} \rangle = \Omega(1)$, $\langle \mathbf{W}_Q^{(T)} \mathbf{x}^{(\ell)}, \mathbf{W}_K^{(T)} \boldsymbol{\xi} \rangle \approx 0$
 - Value focuses on *signal*: $\langle \mathbf{W}_V^{(T)} \mathbf{w}_o, \boldsymbol{\mu} \rangle > \langle \mathbf{W}_V^{(T)} \mathbf{w}_o, \boldsymbol{\xi} \rangle$
 - The test loss is nearly zero: $L_D(\Theta^{(T)}) \approx 0$
- **Harmful overfitting:** under condition $n^{-1} \cdot \text{SNR}^{-2} = \Omega(1)$
 - Attention on *noise*: $\langle \mathbf{W}_Q^{(T)} \mathbf{x}^{(\ell)}, \mathbf{W}_K^{(T)} \boldsymbol{\xi} \rangle = \Omega(1)$, $\langle \mathbf{W}_Q^{(T)} \mathbf{x}^{(\ell)}, \mathbf{W}_K^{(T)} \boldsymbol{\mu} \rangle \approx 0$
 - Value focuses on *noise*: $\langle \mathbf{W}_V^{(T)} \mathbf{w}_o, \boldsymbol{\xi} \rangle > \langle \mathbf{W}_V^{(T)} \mathbf{w}_o, \boldsymbol{\mu} \rangle$
 - The test loss is high: $L_D(\Theta^{(T)}) = \Theta(1)$

OTHER RELEVANT WORKS

Meng et al. 2024. Benign overfitting in two-layer ReLU convolutional neural networks for XOR data. *International Conference on Machine Learning (ICML 2025)*.

Huang et al. 2025. Quantifying the Optimization and Generalization Advantages of Graph Neural Networks Over Multilayer Perceptrons. *International Conference on Artificial Intelligence and Statistics (AISTATS 2025)*.

Karhadkar et al. 2024. Benign overfitting in leaky relu networks with moderate input dimension. *Advances in Neural Information Processing Systems (NeurIPS 2024)*.

Shang et al. 2024. Initialization Matters: On the Benign Overfitting of Two-Layer ReLU CNN with Fully Trainable Layers. *arXiv preprint arXiv:2410.19139*.

Sakamoto & Sato. 2025. Benign Overfitting in Token Selection of Attention Mechanism. *International Conference on Machine Learning (ICML 2025)*.

Frei et al. 2022. Benign overfitting without linearity: Neural network classifiers trained by gradient descent for noisy linear data. *Conference on Learning Theory (COLT 2022)*.

FEATURE LEARNING UNDER DIFFERENT TRAINING STRATEGIES

FEATURE LEARNING UNDER DIFFERENT TRAINING STRATEGY

Optimization and Generalization are **disentangled** in Deep Learning

FEATURE LEARNING UNDER DIFFERENT TRAINING STRATEGY

Optimization and Generalization are **disentangled** in Deep Learning

- A tweak in training strategy can *drastically* affect convergence in training and test error

FEATURE LEARNING UNDER DIFFERENT TRAINING STRATEGY

Optimization and Generalization are **disentangled** in Deep Learning

- A tweak in training strategy can *drastically* affect convergence in training and test error

Adaptive Gradient

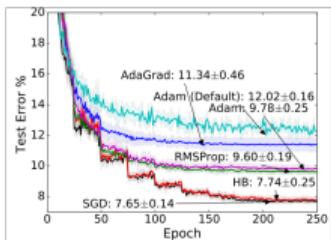


Image credit: Wilson et al. (2017)

Sharpness Aware Minimization

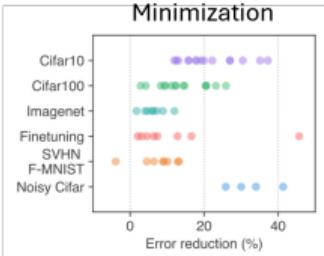


Image credit: Foret et al. (2021)

Momentum

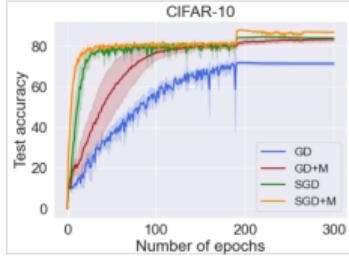


Image credit: Jelassi & Li (2022)

Mixup

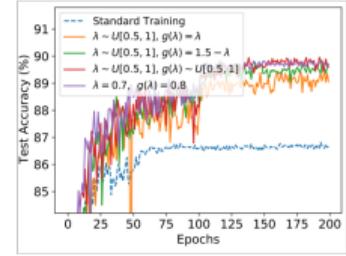


Image credit: Zou et al. (2022)

Large Learning Rate

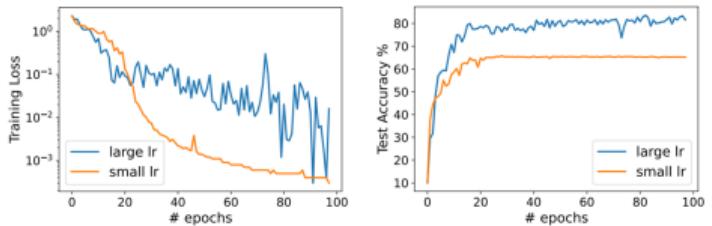


Image credit: Lu et al. (2023)

Label Noise Regularization

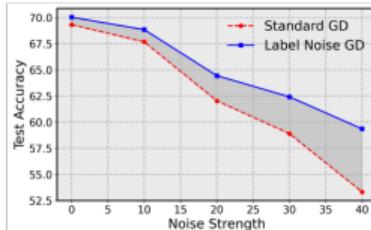


Image credit: Huang et al. (2025)

Matrix Normalization

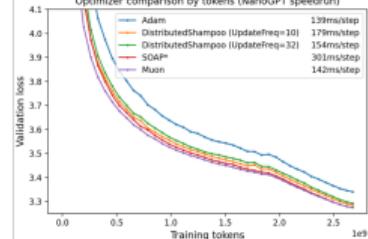


Image credit: Jordan (2025)

FEATURE LEARNING WITH SIGN GD AND ADAM

Adam (Kigma & Ba 2015)

$$\mathbf{m}_t = \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \mathbf{g}_t$$

$$\mathbf{v}_t = \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) \mathbf{g}_t \odot \mathbf{g}_t$$

$$\theta_{t+1} = \theta_t - \eta_t \frac{\mathbf{m}_t}{\sqrt{\mathbf{v}_t} + \epsilon}$$

Sign-GD (Adam when $\beta_1 = \beta_2 = 0$)

$$\theta_{t+1} = \theta_t - \eta_t \text{sign}(\mathbf{g}_t)$$

Sign GD close to Adam

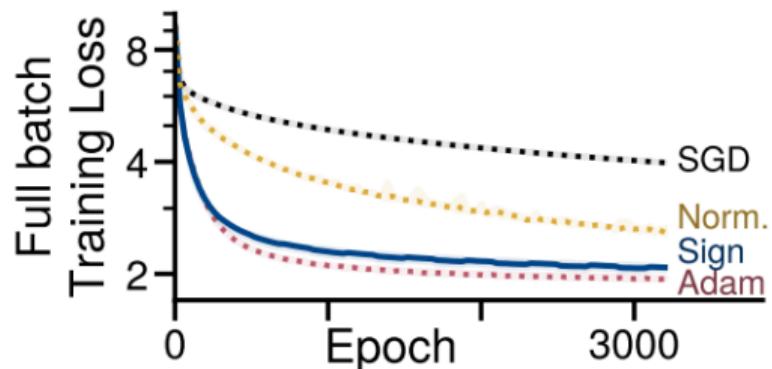


Image Credit: Kunstner et al. 2023.

FEATURE LEARNING WITH SIGN GD AND ADAM

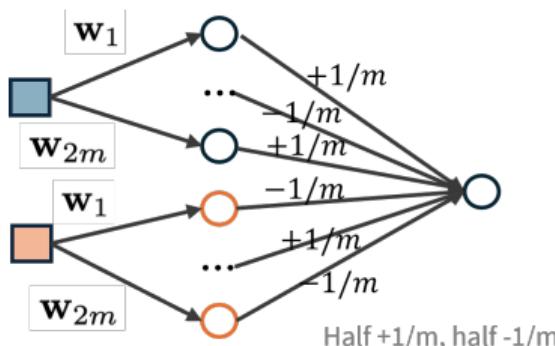
A sparse signal-noise model ([Zou et al., 2023](#))

Data $\mathbf{x} = [y\boldsymbol{\mu}, \boldsymbol{\xi}]$

- (Signal) $\boldsymbol{\mu} = [1, 0, 0, \dots, 0]^\top$
- (Noise) $\boldsymbol{\xi} = \tilde{\boldsymbol{\xi}} \odot \mathbf{s}$, where $\tilde{\boldsymbol{\xi}} \sim \mathcal{N}(\mathbf{0}, \sigma_{\boldsymbol{\xi}}^2 \mathbf{I}_d)$, $\mathbf{s} \in \{0, 1\}^d$ is a random binary mask^a

^aFurther adversarial feature noise is added

Consider the same two-layer CNN (with $\sigma = \text{ReLU}^q$ ($q \geq 3$))



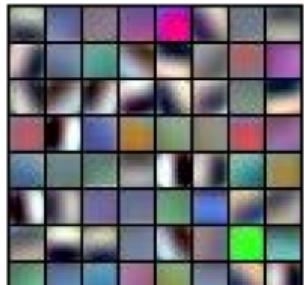
FEATURE LEARNING WITH SIGN GD AND ADAM

Adam generalizes worse than GD ([Zou et al., 2023, Theorem 4.1](#))

- Adam can output a stationary point \mathbf{W}_{adam} in L_1 norm with
 - $L_S(\mathbf{W}_{\text{adam}}) \approx 0, L_D^{0-1}(\mathbf{W}_{\text{adam}}) \geq 0.5$
- GD can output a point \mathbf{W}_{gd} in L_2 norm with
 - $L_S(\mathbf{W}_{\text{gd}}) \approx 0, L_D^{0-1}(\mathbf{W}_{\text{gd}}) \leq 1/\text{poly}(n)$



Adam



GD

Adam learns **more noisy features** than GD
(AlexNet on CIFAR-10).

FEATURE LEARNING WITH SIGN GD AND ADAM

Adam/Sign-GD learns noise faster

$$\begin{aligned}\langle \mathbf{w}_{j,r}^{(t+1)}, j \cdot \mathbf{v} \rangle &\leq \langle \mathbf{w}_{j,r}^{(t)}, j \cdot \mathbf{v} \rangle + \eta \\ \langle \mathbf{w}_{y_i,r}^{(t+1)}, \boldsymbol{\xi}_i \rangle &\approx \langle \mathbf{w}_{y_i,r}^{(t)}, \boldsymbol{\xi}_i \rangle + \eta s \sigma_{\xi}\end{aligned}$$

noise dominates as $s\sigma_{\xi} \gg 1$

GD learns signal faster

$$\begin{aligned}\langle \mathbf{w}_{j,r}^{(t+1)}, j \cdot \mathbf{v} \rangle &\geq \langle \mathbf{w}_{j,r}^{(t)}, j \cdot \mathbf{v} \rangle + \eta \langle \mathbf{w}_{j,r}^{(t)}, j \cdot \mathbf{v} \rangle^{q-1} \\ \langle \mathbf{w}_{j,r}^{(t+1)}, \boldsymbol{\xi}_i \rangle &\leq \langle \mathbf{w}_{j,r}^{(t)}, \boldsymbol{\xi}_i \rangle + \eta s \sigma_{\xi}^2 / n \langle \mathbf{w}_{j,r}^{(t)}, \boldsymbol{\xi}_i \rangle^{q-1}\end{aligned}$$

signal dominates as $s\sigma_{\xi}^2/n \ll 1$

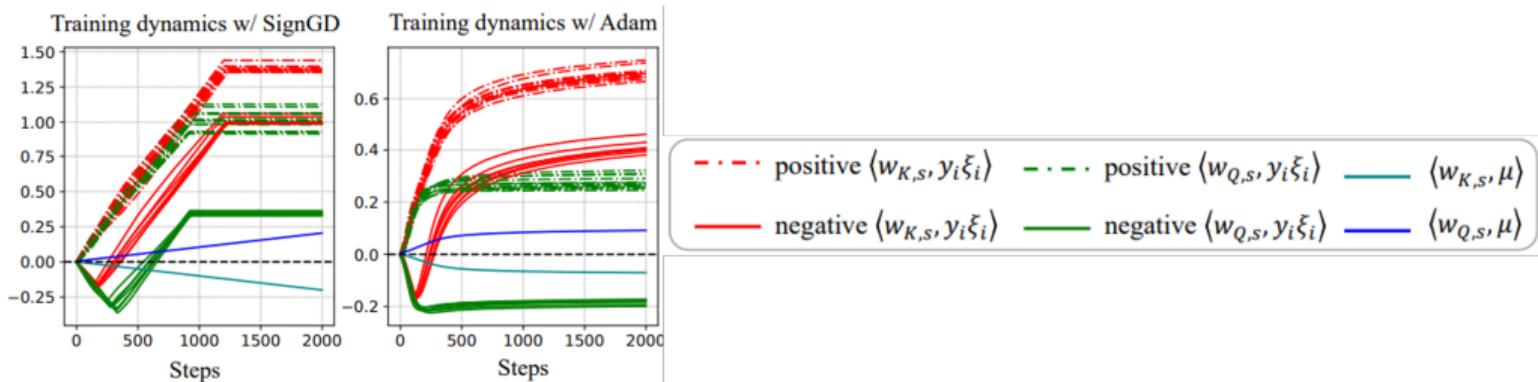
FEATURE LEARNING WITH SIGN GD FOR TRANSFORMER

For a two-layer transformer: a similar result holds

Sign GD converges fast but generalize poorly (Li et al., 2025)

There exists T such that

- Training converges but test loss remains large: $L_S^{(T)}(\mathbf{W}^{(T)}) \leq \epsilon$, and $L_D^{(T)}(\mathbf{W}^{(T)}) = \Theta(1)$.
- Value, query and key matrices memorizes noise.



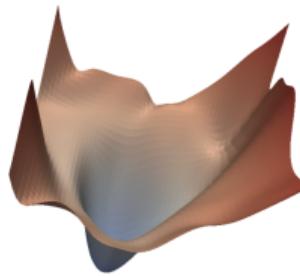
FEATURE LEARNING UNDER DIFFERENT OPTIMIZERS (SAM)

For deep learning, loss landscape is **highly nonconvex**.

- Minimum found by GD (**sharp**)



- Minimum found by SAM (**flat**)



$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \nabla_{\mathbf{w}} L(\mathbf{w}^{(t)})$$

$$\mathbf{g}^{(t)} = \nabla_{\mathbf{w}} L \left(\mathbf{w}^{(t)} + \tau \frac{\nabla_{\mathbf{w}} L(\mathbf{w}^{(t)})}{\|\nabla_{\mathbf{w}} L(\mathbf{w}^{(t)})\|} \right)$$

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \mathbf{g}^{(t)}$$

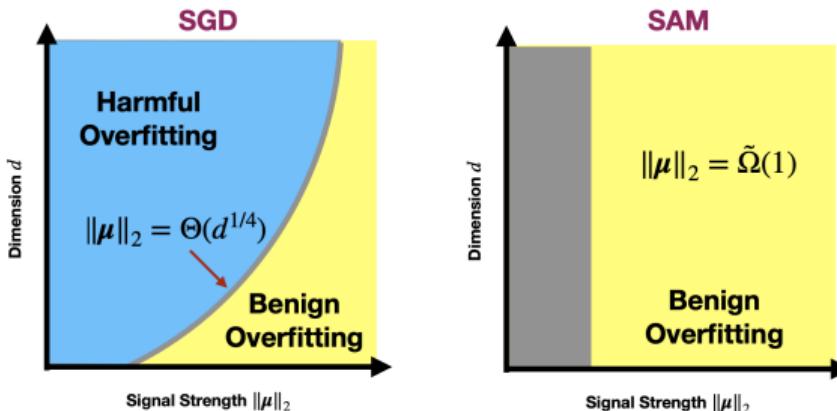
FEATURE LEARNING UNDER DIFFERENT OPTIMIZERS (SAM)

Benign Overfitting of SAM (Chen et al., 2023)

Under Signal-noise data model $\mathbf{x} = [y\boldsymbol{\mu}, \xi]$ and two-layer CNN model with $\sigma = \text{ReLU}$, suppose $\|\boldsymbol{\mu}\| = \tilde{\Omega}(1)$,^a then neural network first trained with SAM, then with SGD can find $\mathbf{W}^{(T)}$ with

- small training loss $L_S(\mathbf{W}^{(T)}) \approx 0$ and small test error $L_D^{0-1}(\mathbf{W}^{(T)}) \approx 0$.

^athis is milder compared to GD Kou et al. (2023), requiring $\|\boldsymbol{\mu}\|^4 = \tilde{\Omega}(d/n)$.



FEATURE LEARNING UNDER DIFFERENT OPTIMIZERS (LABEL NOISE SGD)

Label Noise SGD: A simple regularization by introducing randomness to labels during training.

For each step t and sample (\mathbf{x}_i, y_i) :

1. Sample a random variable $\epsilon_i^{(t)}$:

$$\epsilon_i^{(t)} = \begin{cases} 1 & \text{prob } 1 - p \quad (\text{Keep}) \\ -1 & \text{prob } p \quad (\text{Flip}) \end{cases}$$

2. Update weights using the **noisy gradient**:

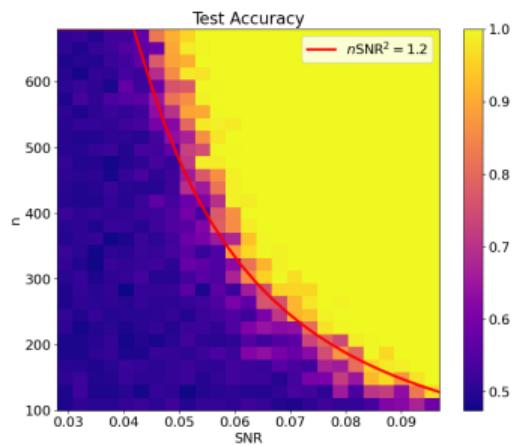
$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \frac{1}{n} \sum_{i=1}^n \nabla \ell(\epsilon_i^{(t)} y_i, f(\mathbf{x}_i))$$

FEATURE LEARNING UNDER DIFFERENT OPTIMIZERS (LABEL NOISE SGD)

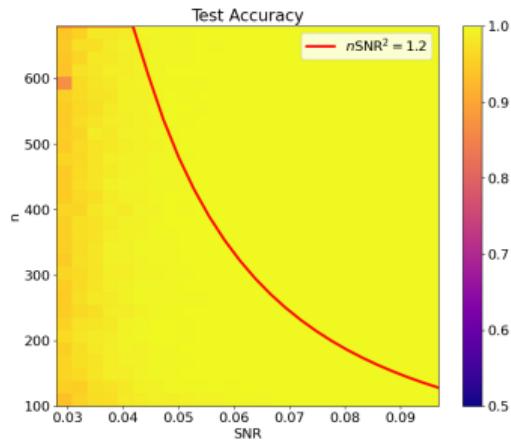
Improved Generalization of Label Noise GD ([Huang et al., 2025](#))

Under Signal-noise data model $x = [y\mu, \xi]$ and two-layer CNN model with $\sigma = \text{ReLU}^2$, then neural network trained with Label Noise GD can find $\mathbf{W}^{(T)}$ with

- constant training loss $L_S(\mathbf{W}^{(T)}) = \Theta(1)$ and small test error $L_D^{0-1}(\mathbf{W}^{(T)}) \approx 0$.



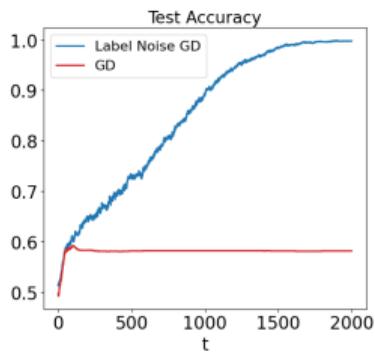
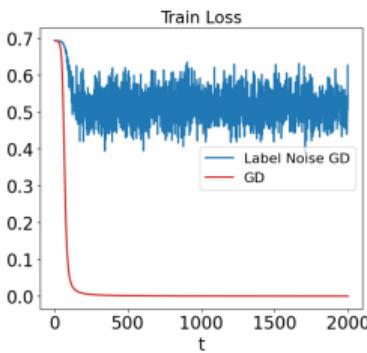
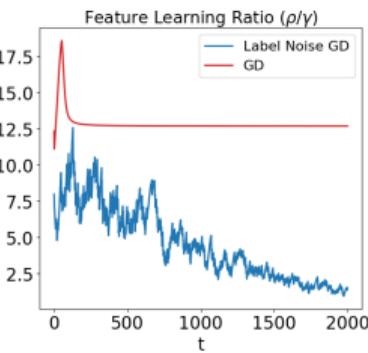
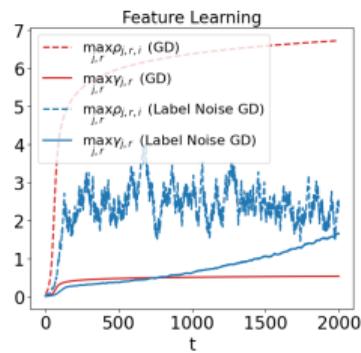
Fails in Low SNR



Robust across regimes

FEATURE LEARNING UNDER DIFFERENT OPTIMIZERS (LABEL NOISE SGD)

Component	Standard GD	Label Noise GD
Signal ($\gamma^{(t)}$)	Grows until loss ≈ 0	Grows exponentially (Stage II)
Noise ($\rho^{(t)}$)	Dominate & Unbounded	Suppressed & Bounded



OTHER RELEVANT WORKS [1]

Jelassi & Li. 2022. Towards understanding how momentum improves generalization in deep learning. *International Conference on Machine Learning (ICML 2022)*.

Chen et al. 2022. Towards understanding the mixture-of-experts layer in deep learning. *Advances in Neural Information Processing Systems (NeurIPS 2022)*.

Zou et al. 2023. The benefits of mixup for feature learning. *International Conference on Machine Learning (ICML 2023)*.

Chen et al. 2023. Understanding and improving feature learning for out-of-distribution generalization. *Advances in Neural Information Processing Systems (NeurIPS 2023)*.

Pan et al. 2024. Federated learning from vision-language foundation models: Theoretical analysis and method. *Advances in Neural Information Processing Systems (NeurIPS 2024)*.

Oh & Yun. 2024. Provable benefit of cutout and cutmix for feature learning. *Advances in Neural Information Processing Systems (NeurIPS 2024)*.

OTHER RELEVANT WORKS [2]

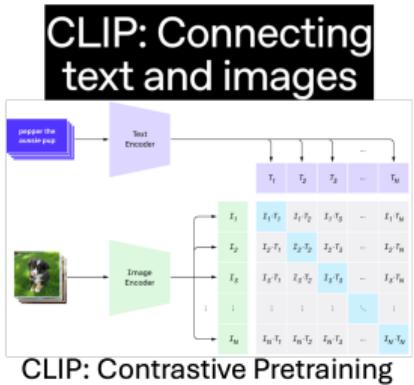
Huang et al. 2024. Understanding convergence and generalization in federated learning through feature learning theory. *International Conference on Learning Representations (ICLR 2024)*.

Lu et al. 2024. Benign Oscillation of Stochastic Gradient Descent with Large Learning Rate. *International Conference on Learning Representations (ICLR 2024)*.

Oh et al. 2025. From linear to nonlinear: Provable weak-to-strong generalization through feature learning. *Advances in Neural Information Processing Systems (NeurIPS 2025)*.

FEATURE LEARNING IN FOUNDATION GENERATIVE MODELS

FEATURE LEARNING IN FOUNDATION GENERATIVE MODELS



Stable Diffusion 3.5
Our most powerful image model yet.

Diffusion Model



Reasoning LLM

Janus Pro AI
Janus Pro AI Unified Multimodal Understanding and Generation Models Build by Deepseek

Multimodal LLM

Image credit: <https://openai.com/index/clip/>, <https://stability.ai/stable-image>, <https://janusai.pro/>, <https://chat-deep.ai/>

MULTI-MODEL CONTRASTIVE LEARNING

Contrastive learning

Draw similar objects (**positive**) closer. Repel dissimilar objects (**negative**)

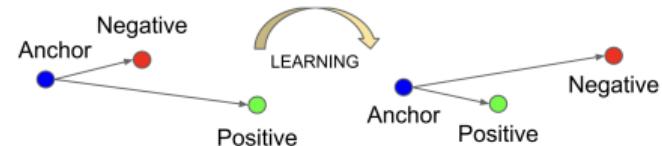


Image credit: Schroff et al. 2015.

MULTI-MODEL CONTRASTIVE LEARNING

Contrastive learning

Draw similar objects (**positive**) closer. Repel dissimilar objects (**negative**)

- *Single-Modal:* positive pairs from **data augmentation**
- *Multi-Modal:* positive pairs from **other modalities**

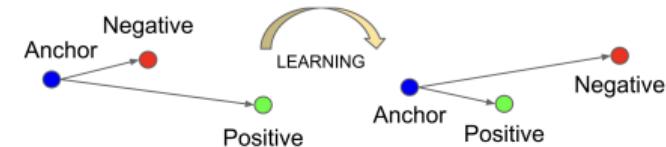
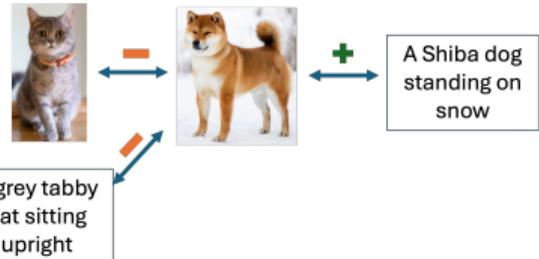


Image credit: Schroff et al. 2015.

Single-Modal Contrastive Learning



Multi-Modal Contrastive Learning



WHAT IS THE DIFFERENCE?

MULTI-MODEL CONTRASTIVE LEARNING [1]

Multi-modal Signal-Noise Model ([Huang et al., 2024](#))

- Two modalities:

$$\mathbf{x} = [\mathbf{x}^{(1)}, \mathbf{x}^{(2)}] = [y\boldsymbol{\mu}, \boldsymbol{\xi}], \quad \tilde{\mathbf{x}} = [\tilde{\mathbf{x}}^{(1)}, \tilde{\mathbf{x}}^{(2)}] = [y\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\xi}}] \text{ (label sharing).}$$

- Nonlinear Embedding:

Let $\text{Emb}(\mathbf{x}) = \sigma(\mathbf{W}\mathbf{x})$ and $\text{Emb}(\tilde{\mathbf{x}}) = \sigma(\widetilde{\mathbf{W}}\tilde{\mathbf{x}})$ where $\sigma = \text{ReLU}$.

- Data Augmentation:

$$\hat{\mathbf{x}} = [\hat{\mathbf{x}}^{(1)}, \hat{\mathbf{x}}^{(2)}] = [y\boldsymbol{\mu}, \boldsymbol{\xi} + \boldsymbol{\epsilon}], \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma_\epsilon^2 \mathbf{I})$$

- Patch-wise Similarity:

$$\text{Sim}(\mathbf{x}, \mathbf{x}') = \left\langle \text{Emb}(\mathbf{x}^{(1)}), \text{Emb}(\mathbf{x}'^{(1)}) \right\rangle + \left\langle \text{Emb}(\mathbf{x}^{(2)}), \text{Emb}(\mathbf{x}'^{(2)}) \right\rangle$$

MULTI-MODEL CONTRASTIVE LEARNING [2]

- Contrastive Loss:

$$\ell(\mathbf{x}, \mathbf{x}^+, \{\mathbf{x}_j^-\}_{j=1}^M) = -\log \left(\frac{e^{\text{Sim}(\mathbf{x}_i, \hat{\mathbf{x}}_i)/\tau}}{e^{\text{Sim}(\mathbf{x}, \hat{\mathbf{x}}')/\tau} + \sum_{j \neq i}^M e^{\text{Sim}(\mathbf{x}, \hat{\mathbf{x}}')/\tau}} \right)$$

- Single-Modal:

$$L = \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{x}_i, \hat{\mathbf{x}}_i, \{\mathbf{x}_j\}_{j \neq i}^M)$$

- Multi-Modal:

$$L = \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{x}_i, \tilde{\mathbf{x}}_i, \{\tilde{\mathbf{x}}_j\}_{j \neq i}^M)$$

Focus on the feature learning of the first modality, i.e., μ, ξ_i .

MULTI-MODEL CONTRASTIVE LEARNING

Multi-modal benefits from *cooperation* between modalities (Huang et al., 2024)

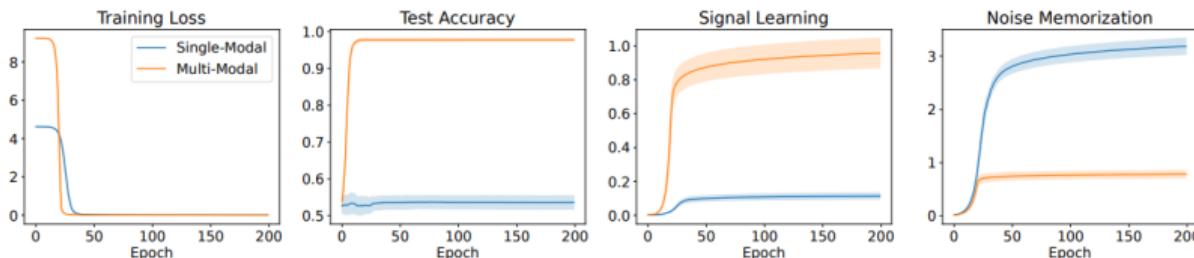
Suppose $n \cdot \text{SNR}^2 = \Theta(1)$ and $\|\tilde{\mu}\| = C_\mu \|\mu\| > \|\mu\|$.

- Single-Modal: memorize noise (data augmentation does not change SNR)

$$\langle \mathbf{w}_r^{(T)}, \boldsymbol{\mu} \rangle \approx 0, \quad \langle \mathbf{w}_r^{(T)}, \boldsymbol{\xi}_i \rangle \geq C$$

- Multi-Modal: learn signal ($\|\tilde{\mu}\| > \|\mu\|$)

$$\langle \mathbf{w}_r^{(T)}, \boldsymbol{\mu} \rangle \geq C', \quad \langle \mathbf{w}_r^{(T)}, \boldsymbol{\xi}_i \rangle \approx 0$$



DIFFUSION MODEL

$$d\mathbf{x}_t = -\mathbf{x}_t dt + \sqrt{2} d\mathbf{W}_t, \quad \mathbf{x}_0 \sim p_0$$

Forward process



Reverse process

(Sohl-Dickstein et al 2015, Ho et al. 2020, Song et al 2021)

$$d\mathbf{z}_t = -(\mathbf{z}_t + 2\nabla \log p_t(\mathbf{z}_t))dt + \sqrt{2} d\bar{\mathbf{W}}_t, \quad \mathbf{z}_T \sim p_T$$

Diffusion model learns the score function via denoising score matching ([Ho et al., 2020](#))

$$\min_{\mathbf{W}} \mathbb{E}_{\mathbf{x}_0 \sim p_0, \epsilon \sim \mathcal{N}(0, \mathbf{I}), t \in [0, T]} \|f(\mathbf{W}, \alpha_t \mathbf{x}_0 + \beta_t \epsilon, t) - \epsilon\|^2$$

DIFFUSION MODEL FEATURE LEARNING

Question: What is the feature learning process of diffusion model? Why do we care?

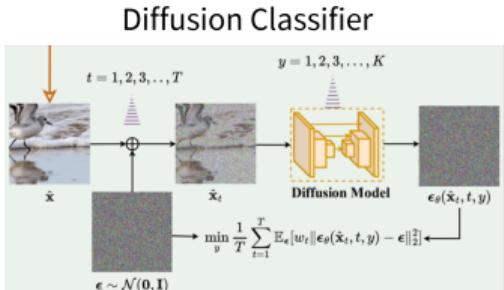


Image credit: Chen, et al. 2024

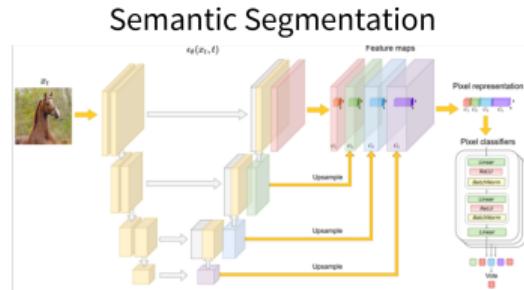


Image credit: Baranchuk, et al. 2022

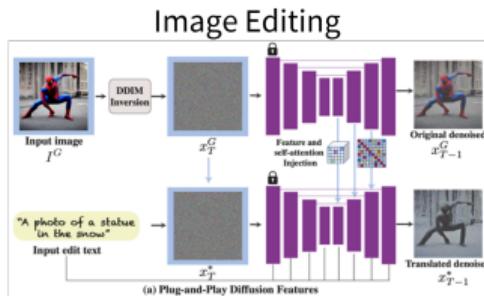


Image credit: Tumanyan et al. 2023

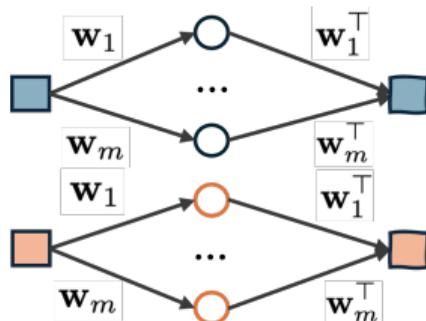
DIFFUSION MODEL FEATURE LEARNING

[Han et al. \(2025a\)](#) compares feature learning of diffusion model with classification models

Diffusion Model

$$\boldsymbol{f}(\mathbf{W}, \boldsymbol{x}) = \left[\boldsymbol{f}_1(\mathbf{W}, \boldsymbol{x}^{(1)}), \boldsymbol{f}_2(\mathbf{W}, \boldsymbol{x}^{(2)}) \right]^\top,$$

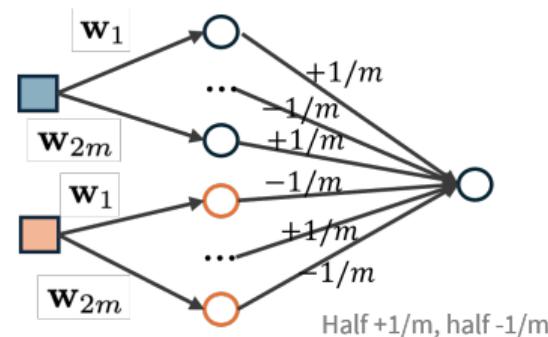
$$f_p(\mathbf{W}, \mathbf{x}^{(p)}) = \frac{1}{\sqrt{m}} \sum_{r=1}^m \sigma(\langle \mathbf{w}_r, \mathbf{x}^{(p)} \rangle) \mathbf{w}_r$$



Classification Model

$$f(\mathbf{W}, x) = F_1(\mathbf{W}_1, x) - F_{-1}(\mathbf{W}_{-1}, x),$$

$$F_j(\mathbf{W}, \mathbf{x}) = \frac{1}{m} \sum_{r=1}^m \sum_{p=1,2} \sigma(\langle \mathbf{w}_{j,r}, \mathbf{x}^{(p)} \rangle)$$



DIFFUSION MODEL FEATURE LEARNING

(Han et al., 2025a, Theorem 3.1 & 3.2)

Diffusion Model learns *balanced* features

- There exists a stationary point \mathbf{W}^* such that

$$|\langle \mathbf{w}_r^*, \boldsymbol{\mu} \rangle| / |\langle \mathbf{w}_r^*, \boldsymbol{\xi} \rangle| = \Theta(n \cdot \text{SNR}^2)$$

Classification learns *dominate* features

- There exists \mathbf{W}^* with $L_S(\mathbf{W}^*) \approx 0$:

– When $n \cdot \text{SNR}^2 \geq \bar{C}$, then $|\langle \mathbf{w}_r^*, \boldsymbol{\mu} \rangle| \geq C$, $|\langle \mathbf{w}_r^*, \boldsymbol{\xi}_i \rangle| \approx 0$ (Signal dominates)

– When $n \cdot \text{SNR}^2 \leq \underline{C}$, then $|\langle \mathbf{w}_r^*, \boldsymbol{\mu} \rangle| \approx 0$, $|\langle \mathbf{w}_r^*, \boldsymbol{\xi}_i \rangle| \geq C'$ (Noise dominates)

IN-CONTEXT LEARNING

In-context learning is the ability of LLMs that learn new rules with few examples.

Fill in the blank with one word: Apple - red, Watermelon - ____



Apple - red, Watermelon - green

Fill in the blank: 311 - 5, 4569 - 24, 12 - ____



311 - 5, 4569 - 24, 12 - 3

Question: Can we understand in-context learning from feature learning?

IN-CONTEXT LEARNING FEATURE LEARNING

Bu et al. (2024): Each prompt contains a **shared concept/task**, with the input

$$\mathbf{H} = \begin{pmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \cdots & \mathbf{x}_L & \mathbf{x}_q \\ \mathbf{y}_1 & \mathbf{y}_2 & \cdots & \mathbf{y}_L & \mathbf{0} \end{pmatrix}, \quad \text{Goal: predict } y_q$$

Each concept k encodes binary semantics $y = \pm 1$:

$$\mathbf{x}_\ell \in \{\mathbf{a}_k + y\mathbf{b}_k\}, \quad \mathbf{y}_\ell \in \{\mathbf{c}_k + y\mathbf{d}_k\}$$

Training data $(\mathbf{H}_n, y_n)_{n=1}^N$:

- Sample $k \in [K]$ and $y \in \{\pm 1\}$
- Construct query $\mathbf{x}_q = \mathbf{a}_k + y\mathbf{b}_k + \boldsymbol{\xi}$, $\mathbf{y}_q = \mathbf{c}_k + y\mathbf{d}_k + \boldsymbol{\xi}'$
- Sample prompt examples: $y_\ell \in \{\pm 1\}$, $\mathbf{x}_\ell = \mathbf{a}_k + y_\ell \mathbf{b}_k + \boldsymbol{\xi}_\ell$, $\mathbf{y}_\ell = \mathbf{c}_k + y_\ell \mathbf{d}_k + \boldsymbol{\xi}'_\ell$, $\ell \in [L]$

IN-CONTEXT LEARNING FEATURE LEARNING [1]

Suppose we train a two-layer transformer on (\mathbf{H}_n, y_n) with expected *cross entropy loss*

$$f(\mathbf{H}) = \mathbf{r}^\top \text{ReLU}(\mathbf{W}_O \text{attn}(\mathbf{H})), \quad \text{attn}(\mathbf{H}) = \sum_{\ell=1}^L \mathbf{W}_V \mathbf{h}_\ell \text{smax} \left(\mathbf{h}_\ell^\top \mathbf{W}_K^\top \mathbf{W}_Q \mathbf{h}_q \right)$$

with the following parameterization

$$\mathbf{W}_Q = \begin{pmatrix} \mathbf{W}_Q^x & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}, \quad \mathbf{W}_K = \begin{pmatrix} \mathbf{W}_K^x & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}, \quad \mathbf{W}_V = \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{pmatrix}, \quad \mathbf{W}_O = \begin{pmatrix} \mathbf{0} & w_O^y \end{pmatrix}$$

Attention only attends to demo inputs and output only depends on demo output.

IN-CONTEXT LEARNING FEATURE LEARNING [2]

Transformer learns *concepts* and *semantics* for in-context learning ([Bu et al., 2024](#))

Upon convergence

- $\mathbf{W}_Q, \mathbf{W}_K$ learn *semantics* rather than concept:

$$\mathbf{a}_k^\top \mathbf{W}_Q^x \mathbf{a}_k \approx 0, \quad \mathbf{b}_k^\top \mathbf{W}_Q^x \mathbf{b}_k = \Omega(1),$$

$$\mathbf{a}_k^\top \mathbf{W}_K^x \mathbf{a}_k \approx 0, \quad \mathbf{b}_k^\top \mathbf{W}_K^x \mathbf{b}_k = \Omega(1),$$

- \mathbf{W}_O learn *both concept and semantics*

$$\langle \mathbf{w}_O^y, \mathbf{c}_k \rangle, \langle \mathbf{w}_O^y, \mathbf{d}_k \rangle = \Omega(1)$$

⇒ this allows to leverage label contains in semantics of x_q for output prediction

NONPARAMETRIC IN-CONTEXT FEATURE LEARNING

Each task k is defined via a task function F_k° (Kim et al., 2024; Kim and Suzuki, 2024)

$$F_k^\circ(\mathbf{x}) = \boldsymbol{\beta}_k^\top \mathbf{f}^\circ(\mathbf{x})$$

where $\boldsymbol{\beta}_k$ is (linear) task-specific and $\mathbf{f}^\circ(\mathbf{x})$ is (nonlinear) task-common features.

Key idea:

- *Pretraining*: learn \mathbf{f}°
- *In-context*: adapt to $\boldsymbol{\beta}_k$

NONPARAMETRIC IN-CONTEXT FEATURE LEARNING

Given pretraining-data (K tasks)

$$\left\{ \begin{pmatrix} \mathbf{x}_{1,k} & \cdots & \mathbf{x}_{L,k} & \mathbf{x}_{q,k} \\ y_{1,k} & \cdots & y_{L,k} & 0 \end{pmatrix}, y_{q,k} \right\}_{k=1}^K$$

Kim and Suzuki (2024) considers linear transformer

$$\frac{1}{L} \sum_{\ell=1}^L y_{\ell,k} \mathbf{h}_\mu(\mathbf{x}_{\ell,k})^\top \boldsymbol{\Gamma} \mathbf{h}_\mu(\mathbf{x}_{q,k}) \xrightarrow{\text{predict}} y_{q,k}$$

with a **mean-field neural network** as feature embedding (*infinite limit of two-layer MLP*):

$$\underbrace{\mathbf{h}_{\theta_m}(\mathbf{x}) = \frac{1}{m} \sum_{r=1}^m \mathbf{a}_r \sigma(\mathbf{w}_r^\top \mathbf{x})}_{\text{two-layer MLP}} \xrightarrow{\text{as } m \rightarrow \infty} \underbrace{\mathbf{h}_\mu(\mathbf{x}) = \int \mathbf{a} \sigma(\mathbf{w}^\top \mathbf{x}) d\mu(\mathbf{a}, \mathbf{w})}_{\text{mean-field limit}}$$

NONPARAMETRIC IN-CONTEXT FEATURE LEARNING

Minimize **expected ICL risk** ($K \rightarrow \infty, L \rightarrow \infty$) w.r.t. μ and Γ

$$\mathcal{L}(\mu, \Gamma) = \mathbb{E}_{\mathbf{x}_q} \left[\|\mathbf{f}^\circ(\mathbf{x}_q) - \mathbb{E}_{\mathbf{x}}[\mathbf{f}^\circ(\mathbf{x}) \mathbf{h}_\mu(\mathbf{x})^\top] \mathbf{\Gamma} \mathbf{h}_\mu(\mathbf{x}_q)\|^2 \right]$$

Key idea: pretraining to $\mathcal{L} = 0$, so that for unseen task β_{new} ,

$$\tilde{y}_q = \mathbb{E}_{\mathbf{x}}[\beta_{\text{new}}^\top \mathbf{f}^\circ(\mathbf{x}) \mathbf{h}_\mu(\mathbf{x})^\top] \mathbf{\Gamma} \mathbf{h}_\mu(\tilde{\mathbf{x}}_q) = \beta_{\text{new}}^\top \mathbf{f}^\circ(\tilde{\mathbf{x}}_q)$$

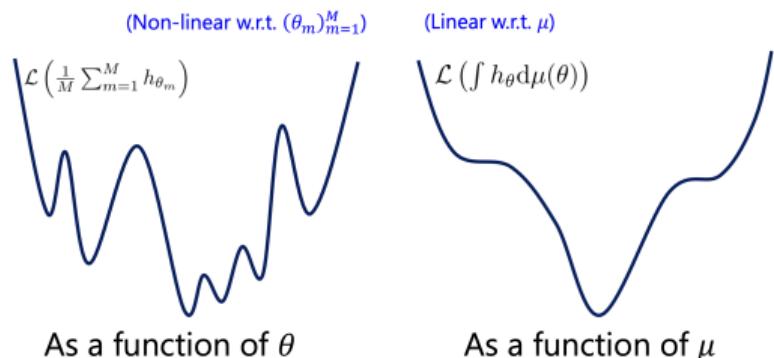


Image credit: Suzuki 2024.

NONPARAMETRIC IN-CONTEXT FEATURE LEARNING

Feature learning under two time-scale dynamics (Γ converges first)

$$\mathcal{L}(\mu) = \min_{\Gamma} \mathcal{L}(\mu, \Gamma) = \mathbb{E}_{\mathbf{x}_q} [\|\mathbf{f}^\circ(\mathbf{x}_q) - \boldsymbol{\Sigma}_{\mu^\circ, \mu} \boldsymbol{\Sigma}_{\mu, \mu}^{-1} \mathbf{h}_\mu(\mathbf{x}_q)\|^2]$$

where $\boldsymbol{\Sigma}_{\mu, \nu} = \mathbb{E}_{\mathbf{x}} [\mathbf{h}_\mu(\mathbf{x}) \mathbf{h}_\nu^\top(\mathbf{x})]$ is the feature covariance, and μ° satisfies $\mathbf{h}_{\mu^\circ} = \mathbf{f}^\circ$.

Wasserstein gradient flow escapes strict saddles and converges to global minimum (Kim and Suzuki, 2024):

- Nonlinear feature learning: $\mathbf{h}_\mu \rightarrow \mathbf{R}\mathbf{f}^\circ$
(for some invertible matrix \mathbf{R} with bounded norm).

OTHER RELEVANT WORKS [1]

Li et al. 2024. How do nonlinear transformers learn and generalize in in-context learning? *International Conference on Machine Learning (ICML 2024)*.

Nichani et al. 2024. How transformers learn causal structure with gradient descent. *International Conference on Machine Learning (ICML 2024)*.

Chen et al. 2024. Unveiling induction heads: Provable training dynamics and feature learning in transformers. *Advances in Neural Information Processing Systems (NeurIPS 2024)*.

Oko et al. 2024. Pretrained transformer efficiently learns low-dimensional target functions in-context. *Advances in Neural Information Processing Systems (NeurIPS 2024)*.

Bu et al. 2025. Provable In-Context Vector Arithmetic via Retrieving Task Concepts. *International Conference on Machine Learning (ICML 2025)*.

Nishikawa et al. 2025. Nonlinear transformers can perform inference-time feature learning. *International Conference on Machine Learning (ICML 2025)*.

OTHER RELEVANT WORKS [2]

Yang et al. 2025. Multi-head Transformers Provably Learn Symbolic Multi-step Reasoning via Gradient Descent.
Advances in Neural Information Processing Systems (NeurIPS 2025).

CONCLUSION AND OUTLOOK

CONCLUSION AND OUTLOOK

- **Feature learning** underlies the *success of deep learning* and provides a *theoretical framework* for **understanding**, **controlling** and **improving** deep learning
 - Benign overfitting (CNN, Transformer)
 - Training strategies (Adam, Sign-GD, SAM, Label noise)
 - Foundation models (contrastive pre-training, diffusion models, in-context learning)
 - and many more

Understanding: Unbox the black-box to study internal representation

Controlling: Manipulate the latent features for controlled model output

Improving: Leverage learned features for model safety, privacy, and robustness.

THANK YOU!

REFERENCES [1]

- Bu, D., Huang, W., Han, A., Nitanda, A., Suzuki, T., Zhang, Q., and Wong, H.-S. (2024). Provably transformers harness multi-concept word semantics for efficient in-context learning. *Advances in Neural Information Processing Systems*, 37:63342–63405.
- Cao, Y., Chen, Z., Belkin, M., and Gu, Q. (2022). Benign overfitting in two-layer convolutional neural networks. *Advances in Neural Information Processing Systems*, 35:25237–25250.
- Chen, Z., Zhang, J., Kou, Y., Chen, X., Hsieh, C.-J., and Gu, Q. (2023). Why does sharpness-aware minimization generalize better than sgd? *Advances in Neural Information Processing Systems*, 36:72325–72376.
- Han, A., Huang, W., Cao, Y., and Zou, D. (2025a). On the feature learning in diffusion models. In *The Thirteenth International Conference on Learning Representations*.

REFERENCES [2]

- Han, A., Huang, W., Zhou, Z., Niu, G., Chen, W., Yan, J., Takeda, A., and Suzuki, T. (2025b). On the role of label noise in the feature learning process. In *International Conference on Machine Learning*.
- Ho, J., Jain, A., and Abbeel, P. (2020). Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851.
- Huang, W., Han, A., Chen, Y., Cao, Y., Xu, Z., and Suzuki, T. (2024). On the comparison between multi-modal and single-modal contrastive learning. *Advances in Neural Information Processing Systems*, 37:81549–81605.
- Huang, W., Han, A., Song, Y., Chen, Y., Wu, D., Zou, D., and Suzuki, T. (2025). How does label noise gradient descent improve generalization in the low snr regime? *Advances in Neural Information Processing Systems*.

REFERENCES [3]

- Jiang, J., Huang, W., Zhang, M., Suzuki, T., and Nie, L. (2024). Unveil benign overfitting for transformer in vision: Training dynamics, convergence, and generalization. *Advances in Neural Information Processing Systems*, 37:135464–135625.
- Kim, J., Nakamaki, T., and Suzuki, T. (2024). Transformers are minimax optimal nonparametric in-context learners. *Advances in Neural Information Processing Systems*, 37:106667–106713.
- Kim, J. and Suzuki, T. (2024). Transformers learn nonlinear features in context: nonconvex mean-field dynamics on the attention landscape. In *International Conference on Machine Learning*, pages 24527–24561.
- Kou, Y., Chen, Z., Chen, Y., and Gu, Q. (2023). Benign overfitting in two-layer relu convolutional neural networks. In *International Conference on Machine Learning*, pages 17615–17659. PMLR.

REFERENCES [4]

- Li, B., Huang, W., Han, A., Zhou, Z., Suzuki, T., Zhu, J., and Chen, J. (2025). On the optimization and generalization of two-layer transformers with sign gradient descent. In *International Conference on Learning Representations*.
- Zou, D., Cao, Y., Li, Y., and Gu, Q. (2023). Understanding the generalization of adam in learning neural networks with proper regularization. In *The Eleventh International Conference on Learning Representations*.